From A to B via a synthesis of rare-event sampling and machine learning



York Universit

NEW YORK UNIVERSITY MRSEC

MATERIALS RESEARCH SCIENCE AND

ENGINEERING CENTERS







<u>External</u>

- Serdal Kirmizialtin (NYU AD)
- Charlles Abreu (Redesign Science)

Postdocs

• Nawavi Naleem (NYU AD)

<u>Students</u>

- Muchen Tong (NYU)
- Shitanshu Bajpai (IIT Kanpur)
- Luke Dai (UC Berkeley)

Funding

جامعـة نيويورك أبوظـي NYU ABU DHABI



@GroupTuckerman



Exemplifying rare events: Alanine decamer





Rare-event MD simulation targeting collective variables [Cuendet et al. *JCP* (2018)]

$$\Delta A_{\text{hairpin} \rightarrow \text{helix}} \approx -40 \text{ kJ/mol}$$

Using deactivated morphing [Park et al. JCP (2008)]:

$$\Delta A_{\text{hairpin} \rightarrow \text{helix}} = -42 \pm 4 \text{ kJ/mol}$$







L. Rosso, P. Minary, Z. Zhu, MET J. Chem. Phys. 116, 4389 (2002), Margliano and Vanden-Eijnden, Chem. Phys. Lett. 426, 168 (2006); J. B. Abrams, MET, J. Phys. Chem. B 112, 14752 (2008). For prediction of crystal structures: Yu and MET Phys. Rev. Lett. 107, 015701 (2011); Yu et al. J. Chem. Phys. 150, 214109 (2014)

Suppose *n* collective variables characterize a free energy landscape of interest

$$q_{\alpha}(\mathbf{R}): \Box^{dN} \to \Box, \qquad \alpha = 1, ..., n < dN$$

Canonical probability distribution and free energy surface:

$$P(s_1, \dots, s_n, T) = \int d\mathbf{R} \ e^{-\beta U(\mathbf{R})} \prod_{\alpha=1}^n \delta(q_\alpha(\mathbf{R}) - s_\alpha)$$
$$A(s_1, \dots, s_n, T) = -k_B T \ln P(s_1, \dots, s_n, T)$$

Write δ -functions as product of Gaussians:

$$\delta(q_{\alpha}(\mathbf{R}) - s_{\alpha}) = \lim_{\{\kappa_{\alpha} \to \infty\}} \left(\frac{\beta\kappa_{\alpha}}{2\pi}\right)^{1/2} \exp\left[-\frac{\beta\kappa_{\alpha}}{2}(q_{\alpha}(\mathbf{R}) - s_{\alpha})^{2}\right]$$

Dynamically sample Gaussian centers s_1, \ldots, s_n by introducing a kinetic energy for them And extending the phase space, which gives the following Hamiltonian:

$$\mathcal{H}(\mathbf{R}, \mathbf{P}, s, p_s) = H(\mathbf{R}, \mathbf{P}) + \sum_{\alpha=1}^{n} \frac{p_{s_\alpha}^2}{2m_\alpha} + \sum_{\alpha=1}^{n} \frac{\kappa_\alpha}{2} \left(q_\alpha(\mathbf{R}) - s_\alpha\right)^2$$

Driven Adiabatic free energy (temperature-accelerated molecular) dynamics (d-AFED/TAMD)



Ensure "on the fly" sampling of the marginal: $m_{\alpha} \square m_i$

Accelerate sampling via temperature: $T_s \square T$ for extended variables

Adiabatically decoupled equations of motion:

$$M_{i}\ddot{\mathbf{R}}_{i} = -\frac{\partial U}{\partial \mathbf{R}_{i}} + \sum_{\alpha} \kappa_{\alpha} \left(s_{\alpha} - q_{\alpha}(\mathbf{R}) \right) \frac{\partial q_{\alpha}}{\partial \mathbf{R}_{i}} + \text{heat bath}(T)$$
$$m_{\alpha} \ddot{s}_{\alpha} = -\kappa_{\alpha} \left(s_{\alpha} - q_{\alpha}(\mathbf{R}) \right) + \text{heat bath}(T_{s})$$



Under adiabatic conditions, we generate a distribution $P_{adb}^{(\{\kappa\})}(s_1,...,s_n,T_s)$

$$\lim_{\{\kappa \to \infty\}} P_{adb}^{(\{\kappa\})}(s_1, ..., s_n, T_s, T) = e^{-\beta_s A(s_1, ..., s_n)} = \left[e^{-\beta A(s_1, ..., s_n)} \right]^{\beta_s / \beta} = \left[P(s_1, ..., s_n, T) \right]^{T/T_s}$$

$$\lim_{\{\kappa \to \infty\}} \left[\ln P_{adb}^{(\{\kappa\})}(s_1, ..., s_n, T_s, T) \right] = \frac{T}{T_s} \ln P(s_1, ..., s_n, T)$$

$$\lim_{\{\kappa \to \infty\}} \left[-k_B T_s \ln P_{adb}^{(\{\kappa\})}(s_1, ..., s_n, T_s, T) \right] = -k_B T_s \frac{T}{T_s} \ln P(s_1, ..., s_n, T) = A(s_1, ..., s_n, T)$$





M. Chen, M. Cuendet, and MET J. Chem. Phys. 137, 024102 (2012); S. Paul, N. N. Nair, H. Vashisth Mol. Simulat. 45, 1273 (2019)

Let $s_1, \tilde{\mathbf{s}} \in \mathbf{s}$

Apply metadynamics-like bias in the \tilde{s} subspace:

$$U_G(\tilde{\mathbf{s}},t) = \sum_i h e^{-\|\tilde{\mathbf{s}}-\tilde{\mathbf{s}}_G(t_i)\|/2\sigma^2}$$



Hamiltonian: ("UFED" method)

$$\mathcal{H}(\mathbf{R}, \mathbf{P}, \mathbf{s}, \mathbf{p}_s) = H(\mathbf{R}, \mathbf{P}) + \sum_{\alpha=1}^n \frac{p_{s_\alpha}^2}{2m_\alpha} + \sum_{\alpha=1}^n \frac{\kappa_\alpha}{2} \left(q_\alpha(\mathbf{R}) - s_\alpha\right)^2 + U_G(\tilde{\mathbf{s}}, t)$$

Hamiltonian: ("TASS" method)

$$\mathcal{H}(\mathbf{R},\mathbf{P},\mathbf{s},\mathbf{p}_s) = H(\mathbf{R},\mathbf{P}) + \sum_{\alpha=1}^n \frac{p_{s_\alpha}^2}{2m_\alpha} + \sum_{\alpha=1}^n \frac{\kappa_\alpha}{2} \left(q_\alpha(\mathbf{R}) - s_\alpha\right)^2 + U_G(\tilde{\mathbf{s}},t) + \frac{1}{2}k_{ub}\left(s_1 - \sigma\right)^2$$

All now available in OpenMM! https://ufedmm.readthedocs.io/en/latest/ Publication: Bajpai et al. J. Comp. Chem. (submitted)



Folding free energy of the Trp-cage

20-residue miniprotein:

Asn-Leu-Tyr-Ile-Gln-Trp-Leu -Lys-Asp-Gly-Gly-Pro-Ser-Ser -Gly-Arg-Pro-Pro-Pro-Ser

Simulated using TASS with 8 CVs: UFED: RMSD of the alpha carbons. UFED: R_g of alpha carbons. D-AFED/TAMD coordinates:

- 1. RMSD of helical residues (2-8).
- 2. RMSD of hydrophobic core (6, 17-19).
- 3. Salt-bridge distance between Asp9-Arg16.
- 4. End-to-end distance.
- 5. Dihedral angle correlation

$$q(\mathbf{R}) = \sum_{i} \left[1 + \cos\left(\phi_{i} - \psi_{i}\right) \right]$$

6. α -helical similarity:

$$q(\mathbf{R}) = \sum_{i} \left[1 + \cos\left(\phi_{i} - \phi_{i}^{ref}\right) \right]$$

CV set from Juraszek and Bolhuis PNAS (2006), Biophys J. (2008)







Folding free energy of the Trp-cage in water



20-residue miniprotein:

Asn-Leu-Tyr-Ile-Gln-Trp-Leu -Lys-Asp-Gly-Gly-Pro-Ser-Ser -Gly-Arg-Pro-Pro-Pro-Ser

Simulated using TASS with 8 CVs: UB(25): RMSD of the alpha carbons UFED: R_g of alpha carbons. D-AFED/TAMD coordinates:

- 1. RMSD of helical residues (2-8).
- 2. RMSD of hydrophobic core (6, 17-19).
- 3. Salt-bridge distance between Asp9-Arg16.
- 4. End-to-end distance.
- 5. Dihedral angle correlation

$$q(\mathbf{R}) = \sum_{i} \left[1 + \cos\left(\phi_{i} - \psi_{i}\right) \right]$$

6. α -helical similarity:

$$q(\mathbf{R}) = \sum_{i} \left[1 + \cos\left(\phi_{i} - \phi_{i}^{ref}\right) \right]$$



(c)



Exp: Streicher & Makhatadze Biochem (2007)









The trained networks can be used to compute observables



Finding optimal reaction coordinates aided machine learning



Paths in classification space:

J. Rogal et al. Phys. Rev. Lett. **123** 245701 (2019)



Path CVs in classification space.





$p_{\mathcal{B}} = 1/2$ Isocommittor surface



Committor learning:

 $p_{\mathcal{B}}(\mathbf{R})$ is the probability of a trajectory initiated at \mathbf{R} reaching \mathcal{B} before \mathcal{A} .

 $p_{\mathcal{B}}(\mathbf{R})$ is an ideal reaction coordinate for \mathcal{A} to \mathcal{B} transition.

Ma and Dinner (2005), Peters and Trout (2006), Peters et al. (2007), Lechner et al. (2010), Jung et al. (2019), Mori et al. (2020), Jung et al. (2021), Kikutusuji et al. (2022),







 $p_{\mathcal{B}} = 1/2$ Isocommittor surface Suppose we have a large set of putative CVs $q_{\alpha}(\mathbf{R})$. Represent $p_{\beta}(\mathbf{R})$ as a linear combination of these:

$$p_{\mathcal{B}}(\mathbf{R};\mathbf{w}) = \sum_{\alpha=1}^{n} w_{\alpha} q_{\alpha}(\mathbf{R}) + w_{0}$$

(Method 1)

Alternatively, $p_{\beta}(\mathbf{R})$ can be approximated via a function $\pi_{\beta}(r(\mathbf{R}))$, where $r(\mathbf{R})$ is a putative reaction coordinate that leads to sigmoidal behavior of $\pi_{\beta}(r(\mathbf{R}))$. A model form for $\pi_{\beta}(r(\mathbf{R}))$ is

$$\pi_{\mathcal{B}}(r(\mathbf{R})) = \frac{1}{1 + e^{-r(\mathbf{R})}}$$

Represent $r(\mathbf{R})$ as a linear combination of the CVs:

$$r(\mathbf{R};\mathbf{w}) = \sum_{\alpha=1}^{n} w_{\alpha} q_{\alpha}(\mathbf{R}) + w_{0}$$

(Method 2)





Learning models for finding an optimal reaction coordinate: Assume N training values of $p_B^*(\mathbf{R}_k)$

• <u>Binary classification (cross-entropy) model (CREM):</u>

$$E(\mathbf{w}) = -\sum_{k=1}^{N} p_{\mathcal{B}}^{*}(\mathbf{R}_{k}) \ln p_{\mathcal{B}}(\mathbf{R}_{k}, \mathbf{w}) - \sum_{k=1}^{N} \left(1 - p_{\mathcal{B}}^{*}(\mathbf{R}_{k})\right) \ln \left[1 - p_{\mathcal{B}}(\mathbf{R}_{k}, \mathbf{w})\right] + \operatorname{Reg}(\mathbf{w})$$
(Method 1)
$$E(\mathbf{w}) = -\sum_{k=1}^{N} p_{\mathcal{B}}^{*}(\mathbf{R}_{k}) \ln \pi_{\mathcal{B}}(r(\mathbf{R}_{k}, \mathbf{w})) - \sum_{k=1}^{N} \left(1 - p_{\mathcal{B}}^{*}(\mathbf{R}_{k})\right) \ln \left[1 - \pi_{\mathcal{B}}(r(\mathbf{R}_{k}, \mathbf{w}))\right] + \operatorname{Reg}(\mathbf{w})$$
(Method 2)

• <u>Standard regression model:</u>

$$E(\mathbf{w}) = \frac{1}{2N} \sum_{k=1}^{N} \left| p_{\mathcal{B}}(\mathbf{R}, \mathbf{w}) - p_{\mathcal{B}}^{*}(\mathbf{R}_{k}) \right|^{2} + \operatorname{Reg}(\mathbf{w})$$
 (Method 1)

$$E(\mathbf{w}) = \frac{1}{2N} \sum_{k=1}^{N} \left| \pi_{\mathcal{B}}(r_k(\mathbf{R}, \mathbf{w})) - p_{\mathcal{B}}^*(\mathbf{R}_k) \right|^2 + \operatorname{Re} g(\mathbf{w})$$
 (Method 2)

 $\operatorname{Reg}(\mathbf{w}) = \begin{cases} \lambda \|\mathbf{w}\| & \text{Lasso (LASR)} \\ \lambda \rho \|\mathbf{w}\| + \beta (1 - \rho) \|\mathbf{w}\|^2 & \text{Elastic Net (ELAN)} \end{cases}$





List of standard regressors:

Lasso (LASR) Elastic Net (ELAN) Lasso Least Angle (LLAR) Huber (HUBR) Linear (LINR) Ridge (RIDR) Bayeian (BAYR) Orthogonal Mathing Pursuit (ORMP) Passive-Aggressive (PAGR)

Neighbors/trees:

Decision Tree (DECT) Random Forest (RAFR) Extra Trees (EXTR)

Splitting rule based on variance reduction:

$$\sigma^{2} = \frac{\sum_{k=1}^{n} \left(p_{\mathcal{B}}^{*}(\mathbf{R}_{k}) - \overline{p}_{\mathcal{B}} \right)^{2}}{n}$$

Boosters:

AdaBoost (ADAB) Gradient Boost (GRBR)

Light Gradient Boost (LGBM)

Boosting:

```
\pi_{B}^{(0)}(\mathbf{R}) = \arg\min_{\mathbf{w}} E(\mathbf{w})
for m = 1 to M
\overline{p}_{B,k} = p_{B,k} - \pi_{B}^{(m-1)}(\mathbf{R}_{k})
(\mathbf{w}_{m}, \beta_{m}) = \arg\min_{\beta, \mathbf{w}} \sum_{k=1}^{N} (\overline{p}_{B,k} - \pi_{B,k}^{(m-1)} - \beta h(\mathbf{R}_{k}; \mathbf{w}))^{2}
\pi_{B}^{(m)}(\mathbf{R}) = \pi_{B}^{(m-1)}(\mathbf{R}) + l_{r}\beta_{m}h(\mathbf{R}; \mathbf{w}_{m})
endfor
```

Reaction coordinate of C7eq -> C7ax conformational transition in Ala dipeptide in vacuum



Alanine dipeptide





Committor values sampled using aimless shooting algorithm.





Reaction coordinate of C7eq -> C7ax conformational transition in Ala dipeptide



Alanine dipeptide



Enhanced sampling of (φ, ψ)



Method 1

0.50

 R_2 of P_B

0.25

0.75

1.00

0.00

0.05

ADAB

BAYR

CREM -

DECT

ELAN

EXTR ·

GRBR

HUBR ·

KNNR -

LLAR

LASR

LGBM

ORMP

PAGR

RAFR -

RIDR

0.00

LINR

Method 2



Reaction coordinate of C7eq -> C7ax conformational transition in Ala dipeptide



LGBM

CREM

3

3







Alanine dipeptide



$$r(\mathbf{q}, \mathbf{w}) p_{\mathcal{B}}(\mathbf{q}, \mathbf{w})$$
 = $b + \sum_{i=1}^{45} w_{i-1} \cos q_i + \sum_{i=46}^{90} w_{i-1} \sin q_i$

SHAP (SHapley Additive ePlanation – L. S. Shapley in *Contributions to the Theory of Games* (2016)).

$$\varphi_i = \sum_{S \subseteq F \setminus \{i\}} \binom{|F|}{|S|}^{-1} \left[f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S) \right]$$

Each dot in the SHAP plot is an observation in our data set. "High" and "low" are high and low values of the features. Positive and negative values describe whether high or low values of the feature push the predictions toward one basin or the other (0 or 1 commitor probability values).





Feature extraction forC7eq -> C7ax conformational transition in Ala dipeptide







$r(\mathbf{q},\mathbf{w}) \Big _{-b}$	$\sum_{k=1}^{45} w \cos \alpha \pm$	$\sum_{n=1}^{90} w \sin a$
$p_{\mathcal{B}}(\mathbf{q},\mathbf{w})\int^{-D+}$	$\sum_{i=1}^{N} w_{i-1} \cos q_i +$	$\sum_{i=46}^{W_{i-1}} \operatorname{SIII} q_i$

Dihedral Index	Abbriv.	1 2	3	4	Dihedral Index	Abbriv.	1	2	3	4
1	β_1	12	5	6	24	γ_8	15	9	11	13
2	β_2	$1 \ 2$	5	7	25	γ_9	15	9	11	14
3	β_3	32	5	6	26	ψ_1	7	9	15	16
4	β_4	32	5	7	27	ψ_2	7	9	15	17
5	β_5	4 2	5	6	28	ψ_3	10	9	15	16
6	β_6	4 2	5	7	29	ψ_4	10	9	15	17
7	θ_1	25	7	8	30	ψ_5	11	9	15	16
8	θ_2	25	7	9	31	ψ_6	11	9	15	17
9	θ_3	65	7	8	32	δ_1	9	15	17	18
10	θ_4	65	7	9	33	δ_2	9	15	17	19
11	ϕ_1	57	9	10	34	δ_3	16	15	17	18
12	ϕ_2	57	9	11	35	δ_4	16	15	17	19
13	ϕ_3	57	9	15	36	ϵ_1	15	17	19	20
14	ϕ_4	87	9	10	37	ϵ_2	15	17	19	21
15	ϕ_5	87	9	11	38	€3	15	17	19	22
16	ϕ_6	87	9	15	39	ϵ_4	18	17	19	20
17	γ_1	79	11	12	40	ϵ_5	18	17	19	21
18	γ_2	79	11	13	41	ϵ_6	18	17	19	22
19	γ_3	79	11	14	42	ζ_1	2	$\overline{7}$	5	6
20	γ_4	$10 \ 9$	11	12	43	η_1	5	9	$\overline{7}$	8
21	γ_5	$10 \ 9$	11	13	44	κ_1	9	17	15	16
22	γ_6	109	11	14	45	λ_1	15	19	17	18
23	γ_7	$15 \ 9$	11	12						







- X-22-25-Υ ψ X-12
 - X-12-15-Y → ψ(-1)







Reaction coordinate of trans->cis conformational transition in disarcosine in GB solvent









Enhanced sampling of



Reaction coordinate of trans->cis conformational transition in disarcosine





200 🔌

250

⁰ ⁵⁰ 100₁₅₀₂₀₀₂₅₀₃₀₀₃₅₀ ³⁵⁰





LGBM





Reaction coordinate of trans->cis conformational transition in disarcosine

 $21 \ 18 \ 17 \ 22$

22 25 27 28

 $22 \ 25 \ 27 \ 32$

23 22 25 26

 $23 \ 22 \ 25 \ 27$

24 22 25 27

 $24 \ 22 \ 25 \ 26$

25 27 32 33

25 27 28 29

25 27 32 35

25 27 28 31

25 27 32 34

25 27 28 30

26 25 27 32

 $26\ 25\ 27\ 28$

28 27 32 35

28 27 32 34

28 27 32 33

29 28 27 32

30 28 27 32

31 28 27 32

 ζ_3

 η_1

 η_2

 ψ_3

 ψ_4

 ψ_5

 ψ_6

 κ_1

 λ_1

 κ_2

 λ_2

 κ_3

 λ_3

 η_3

 η_4

 κ_4

 κ_5

 κ_6

 μ_1

 μ_2

 μ_3



Meth 1



13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

$r(\mathbf{q}, \mathbf{w})$ $p_{_{\mathcal{B}}}(\mathbf{q}, \mathbf{v})$	$\left. \begin{array}{c} \mathbf{v} \\ \mathbf{v} \end{array} \right\} =$	b	+	6 2 <i>i</i> =	6 =1	$W_{i-1} \cos \theta$	$sq_i + \sum_{i=6'}^{132}$	W_{i}	₋₁ S	$in q_i$
-	411 -	_				D-1 1 11	1 411 -	1.0		
Dihedral Index	Abbriv.	1	2	3	4	Dihedral Ii	idex Abbriv.	1 2	3	4
1	α_1	1	5	$\overline{7}$	12	34	ϕ_2	$15 \ 17$	22 2	25
2	α_2	1	5	7	8	35	ϵ_3	$15 \ 17$	18 2	21
3	β_1	2	1	5	7	36	ϕ_3	$15 \ 17$	22 2	24
4	β_2	2	1	5	6	37	ω_3	16 15	17 2	22
5	β_3	3	1	5	7	38	ω_4	16 15	17 1	18
6	β_4	3	1	5	6	39	ψ_1	17 22	25 2	27
7	β_5	4	1	5	7	40	ψ_2	17 22	25 2	26
8	β_6	4	1	5	6	41	ϕ_4	18 17	22 2	25
9	γ_1	5	7	8	10	42	ϕ_5	18 17	22 2	24
10	$\phi(-1)_1$	5	7	12	13	43	ϕ_6	18 17	22 2	23
11	γ_2	5	7	8	9	44	ζ_1	19 18	3 17 2	22
12	$\phi(-1)_2$	5	7	12	15	45	ζ_2	20 18	17.2	22

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

5 7 8 11

 $5\ 7\ 12\ 14$

6 5 7 12

6 5 7 8

7 12 15 16

8 7 12 14

9 8 7 12

10 8 7 12

11 8 7 12

12 15 17 22

12 15 17 18

15 17 18 20

15 17 22 23

15 17 18 19

 $\psi(-1)_1$ 7 12 15 17

 $\phi(-1)_4 = 8 - 7 - 12 - 15$

 $\phi(-1)_6 \quad 8 \quad 7 \quad 12 \quad 13$

 $\psi(-1)_3$ 13 12 15 17

 $\psi(-1)_4$ 13 12 15 16

 $\psi(-1)_5$ 14 12 15 17

 $\psi(-1)_6$ 14 12 15 16

 γ_3

 $\phi(-1)_{3}$

 α_3

 α_4

 $\psi(-1)_2$

 $\phi(-1)_{5}$

 δ_1

 δ_2

 δ_3

 ω_1

 ω_2

 ϵ_1

 ϕ_1

 ϵ_2



Meth 2



Mechanism of trans->cis conformational transition in disarcosine











CVs include sines and cosines of 17 dihedral angles plus 152 symmetry functions of the form:

$$G_2^i = \sum_j e^{-\eta (R_{ij}-R_s)^2} \cdot f_c(R_{ij}),$$

Try out schemes for automatically selecting descriptors: G. Imbalzano *et al. JCP* (2018).





- Interoperability of d-AFED/TAMD, metadynamics, and umbrella sampling leads to an algorithm ideally suited for use with machine learning method for regressing high-dimensional FESs, as it sweeps over the full landscape, producing an initially sparse sampling (the density of which grows with simulation length) that regression methods can easily interpolate.
- From a large set of redundant CVs, committor learning is possible via either a direct mapping of the CVs to the committor or via a model of the committor in terms of a reaction coordinate expressed in terms of the CVs.
- For the systems studied, light gradient boosting proved to be the best performing of the machine-learning models investigated.
- SHAP analysis allows the most important features to be extracted, which reduces the number of degrees of freedom needed to express the committor or reaction coordinate.