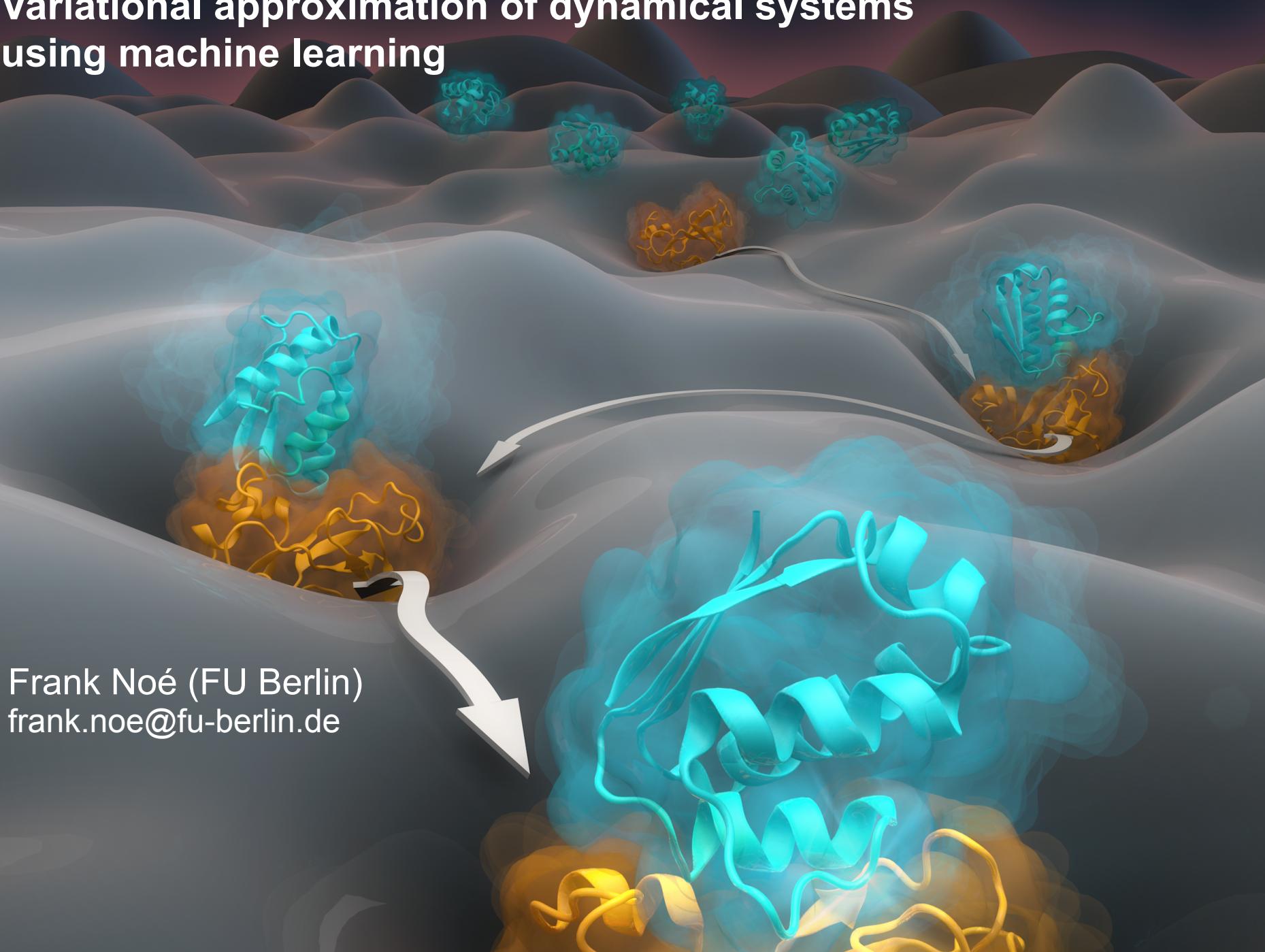


Variational approximation of dynamical systems using machine learning

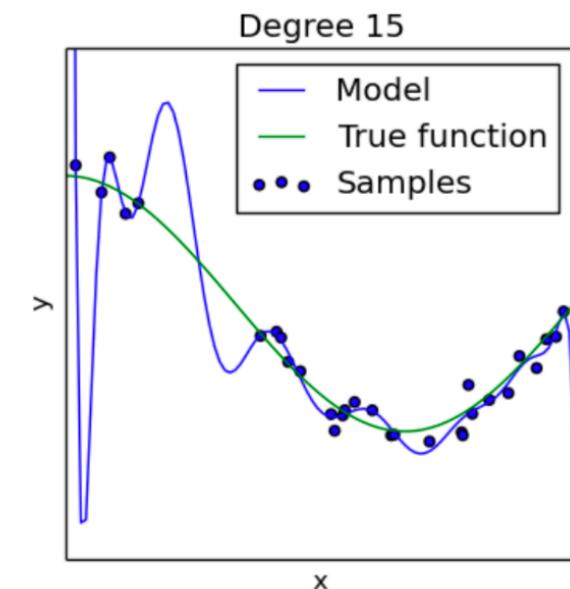
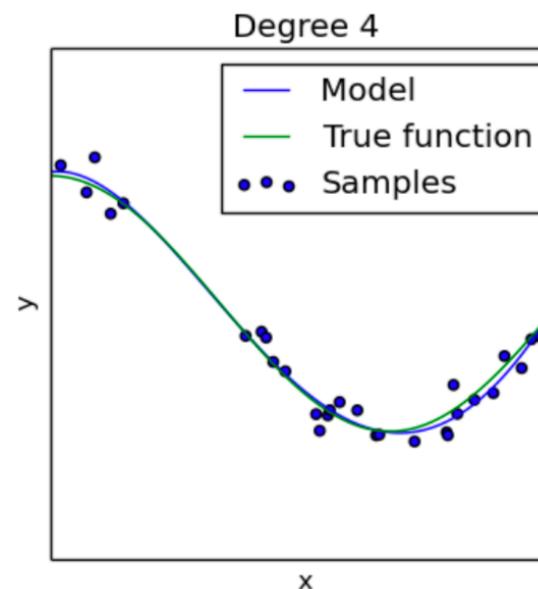
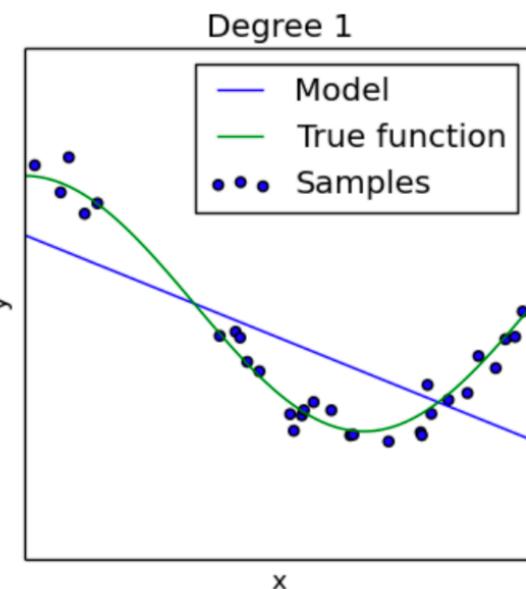


Frank Noé (FU Berlin)
frank.noe@fu-berlin.de

Aim: learn model of dynamical system from observations $x_t, x_{t+\tau}$

How do we choose?

- Method / Representation (DMD, EDMD, TICA, VAMPnets, time-Autoencoder, ...)
- Type of Basis Set (polynomial, cos/sin, characteristic functions, ...)
- Number of Basis Functions
- Type of Kernel in a Kernel approach
- Regularization Hyperparameters
- Type of Neurons and Architecture of Neural Net



Cost or loss function C quantifies performance of network with parameters θ to predict observations \mathbf{X} .

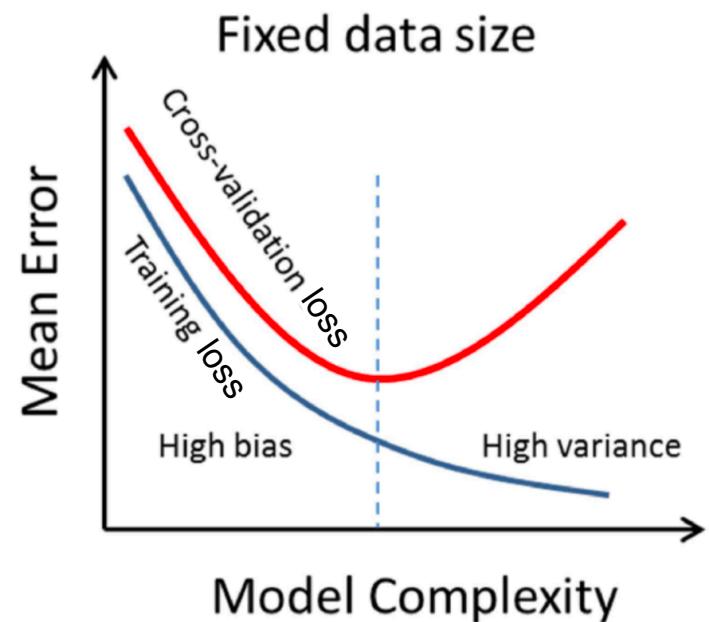
Learning network weights

$$\hat{\theta} = \arg \min_{\theta} C(\mathbf{X}, \mathbf{Y}, \theta)$$

Minimizing cost $C \equiv$ maximizing score $-C$. Often we just write $C(\theta)$.

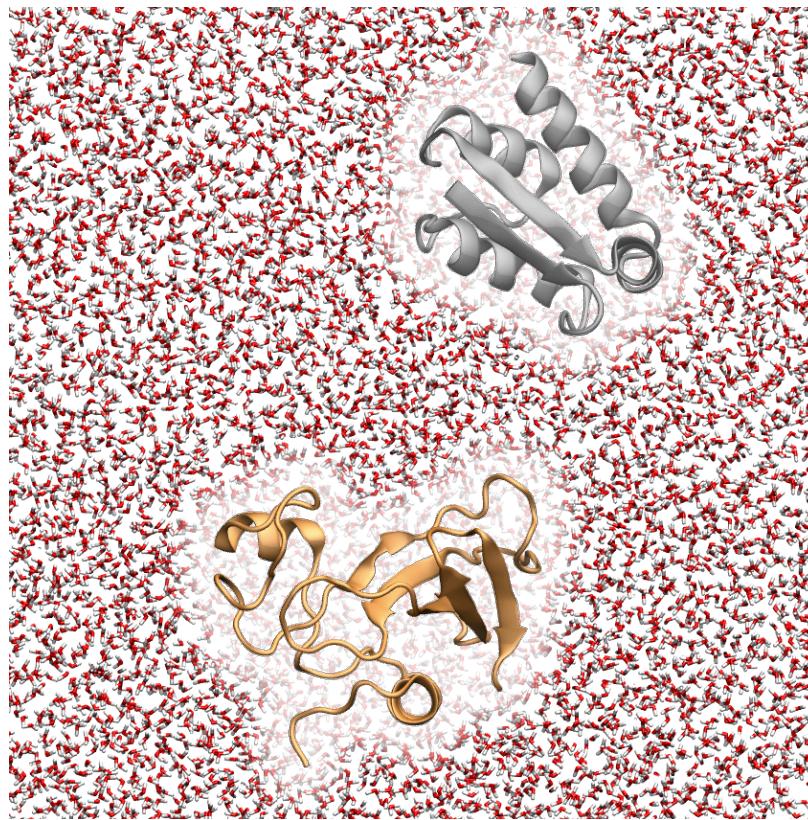
Standard approach (general and flexible):

- Parameter optimization: Training data
- Hyperparameter optimization: Validation data
- Quantifying model performance: Test data
- Details: **Statistical Estimator Theory** (e.g., Vapnik and others)



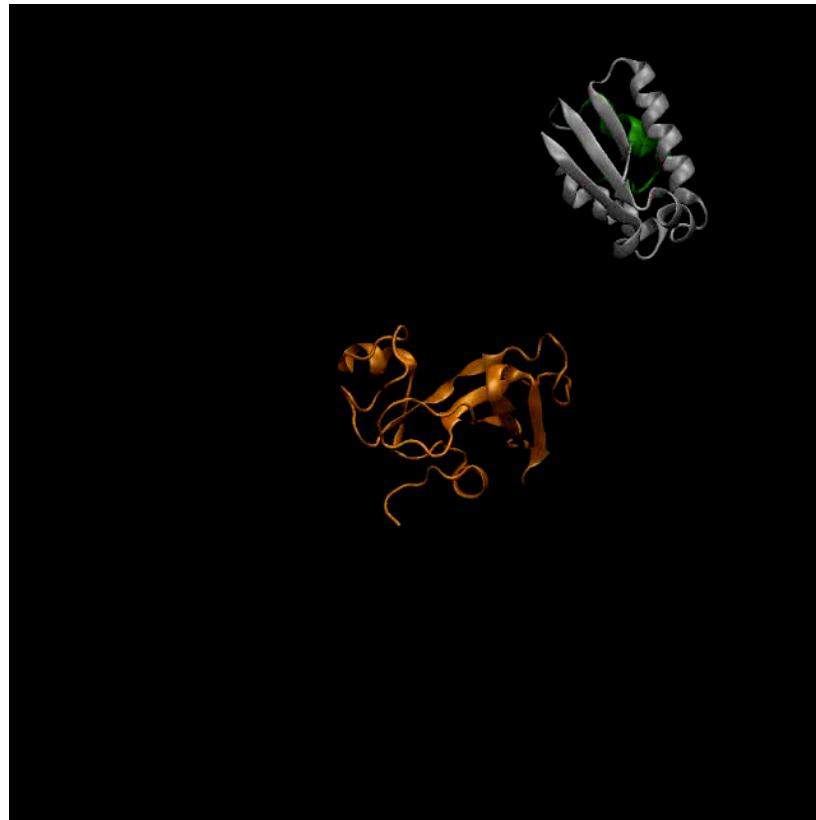
Simulating biological timescales at atomic resolution

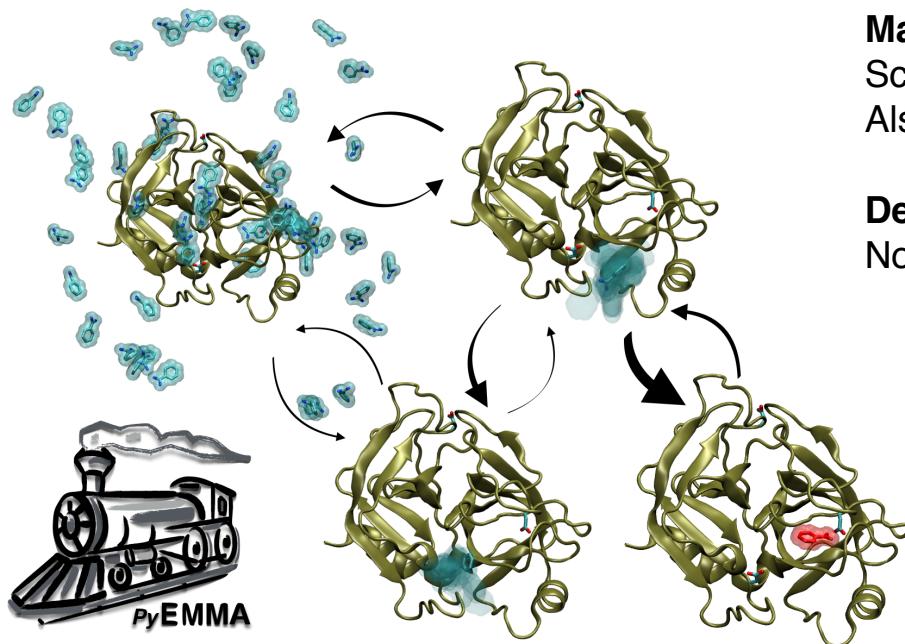
**Microsecond
MD Trajectories**



Simulating biological timescales at atomic resolution

**Microsecond
MD Trajectories**





Mathematical theory:

Schütte et al, **J Comp Phys** 1999,

Also: Weber, Deuflhard, Friesecke, Dellnitz ...

Developments for high-throughput molecular dynamics:

Noé, Pande, Swope, Hummer (mid 2000's)

Propagator

$$\rho_{t+\tau} = \mathcal{P}_\tau \rho_t = \int p_\tau(y \mid x) \rho_t(x) dx$$

Transfer operator / Perron-Frobenius operator

= propagator for densities $u(x) = \frac{\rho(x)}{\pi(x)}$ with stationary density $\pi(x)$.

$$u_{t+\tau} = \mathcal{T}_\tau u_t = \int \frac{\pi(x)}{\pi(y)} p_\tau(y \mid x) \rho_t(x) dx$$

Koopman operator

Adjoint to \mathcal{P} , adjoint to \mathcal{T} with respect to π

$$f_{t+\tau} = \mathcal{K}_\tau f_t = \int p_\tau(y \mid x) f_t(y) dy = \mathbb{E}[f_{t+\tau}(x)]$$

with detailed balance: $\pi(x)p_\tau(y \mid x) = \pi(y)p_\tau(x \mid y)$ we have $\mathcal{K}_\tau \equiv \mathcal{T}_\tau$

See: Mesic **Nonlinear Dyn.** 41, 309 (2005).

Slow processes (unique equilibrium distribution, detailed balance)

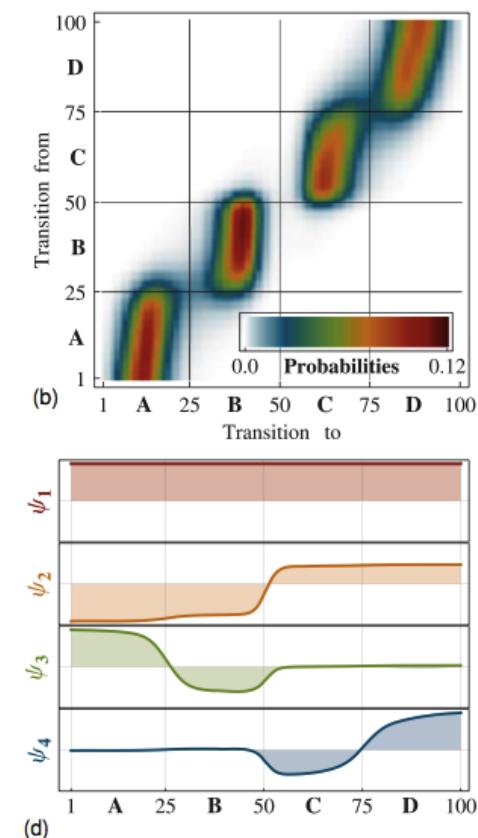
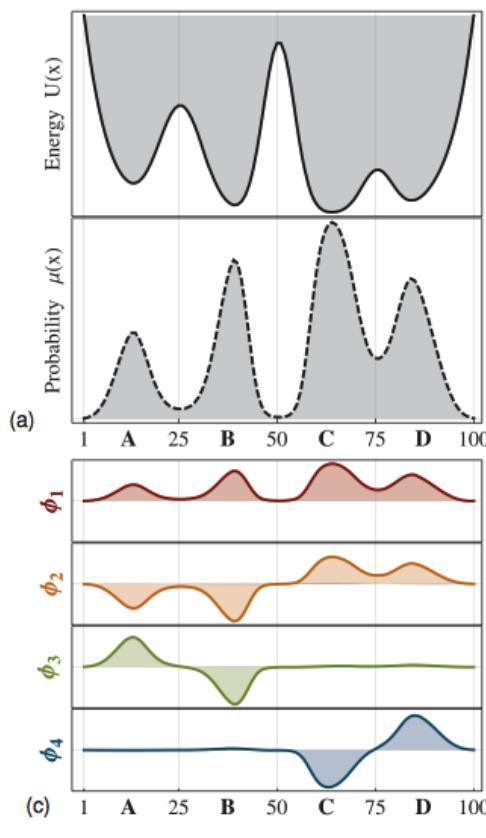
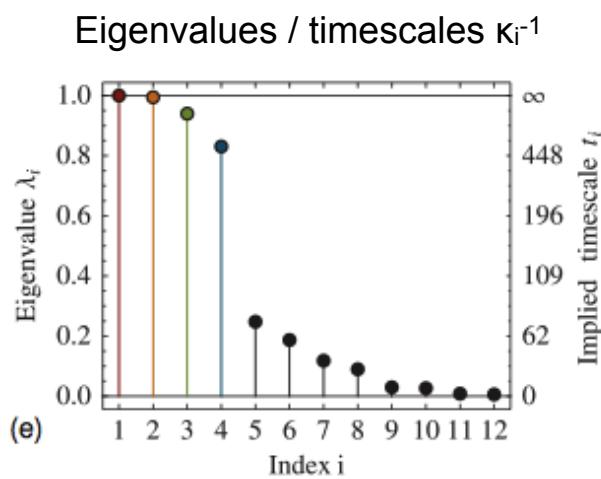
Backward propagator

$$\rho_\tau = \mathcal{T}(\tau)\rho_0$$

Spectral decomposition

$$\rho_\tau = \sum_{i=1}^{\infty} e^{-\tau \kappa_i} \langle \psi_i | \rho_0 \rangle \psi_i + \text{fast part}$$

Processes:



Schütte et al: *J. Comput. Phys.* (1999), Prinz et al.: *J. Chem. Phys.* 134, p174105 (2011)

Variational approach for Markov processes

Data-based version of: Fan, **PNAS** 35, 652-655 (1949)

The first m eigenfunctions ψ_1, \dots, ψ_m are the solution to the problem

$$\begin{aligned} & \max_{f_1, \dots, f_m} \sum_{i=1}^m \mathbb{E} [f_i(\mathbf{x}_t) f_i(\mathbf{x}_{t+\tau})] \\ \text{s.t. } & \mathbb{E} [f_i(\mathbf{x}_t)^2] = 1 \\ & \mathbb{E} [f_i(\mathbf{x}_t) f_j(\mathbf{x}_{t+\tau})] = 0, \text{ for } i \neq j \end{aligned} \tag{1}$$

and the maximum value is the sum of $\lambda_1, \dots, \lambda_m$

Properties:

- ψ_i and ψ_j are uncorrelated for $i \neq j$.
- ψ_i are the directions of slow kinetics with maximal autocorrelations $\mathbb{E}_\mu [\psi_i(\mathbf{x}_t) \psi_i(\mathbf{x}_{t+\tau})] = \lambda_i(\tau)$.
- Population changes along ψ_i coordinates decay with $\lambda_i(\tau) = e^{-\frac{\tau}{t_i}}$.
- For every other set of functions, the eigenvalues will be underestimated $\hat{\lambda}_i(\tau) \leq \lambda_i(\tau)$.

Noé and Nüske, **MMS** 11, 635-655 (2013)
Nüske et al, **JCTC** 10, 1739-1752 (2014)

Method of linear variation

Ansatz: Define Basis set $\chi = [\chi_1(\mathbf{x}), \dots, \chi_n(\mathbf{x})]^\top$ and seek the linear expansions:

$$\hat{\psi}_i(\mathbf{x}) = \sum_j r_{ij} \chi_j(\mathbf{x})$$

Noé and Nüske, **MMS** 11, 635-655 (2013)
Nüske et al, **JCTC** 10, 1739-1752 (2014)

Variational approach for reversible Markov processes: Estimator

1. Define

$$\mathbf{X}_0 = \begin{bmatrix} \chi_1(\mathbf{x}_0) & \cdots & \chi_n(\mathbf{x}_0) \\ \vdots & & \vdots \\ \chi_1(\mathbf{x}_{T-\tau}) & \cdots & \chi_n(\mathbf{x}_{T-\tau}) \end{bmatrix} \quad \mathbf{X}_\tau = \begin{bmatrix} \chi_1(\mathbf{x}_\tau) & \cdots & \chi_n(\mathbf{x}_\tau) \\ \vdots & & \vdots \\ \chi_1(\mathbf{x}_T) & \cdots & \chi_n(\mathbf{x}_T) \end{bmatrix}$$

2. Empirical covariance matrices: \mathbf{C}^0 and \mathbf{C}^τ with:

$$\mathbf{C}^0 = \mathbf{X}_0^\top \mathbf{X}_0$$

$$\mathbf{C}^\tau = \mathbf{X}_\tau^\top \mathbf{X}_\tau$$

3. Solve

$$\mathbf{C}^\tau \mathbf{r}_i = \mathbf{C}^0 \hat{\lambda}_i \mathbf{r}_i$$

4. The projections

$$\Psi = \mathbf{X} \mathbf{R}$$

approximate the transfer operator eigenfunctions on the sampled configurations x_t .

Noé and Nüske, **MMS** 11, 635-655 (2013)

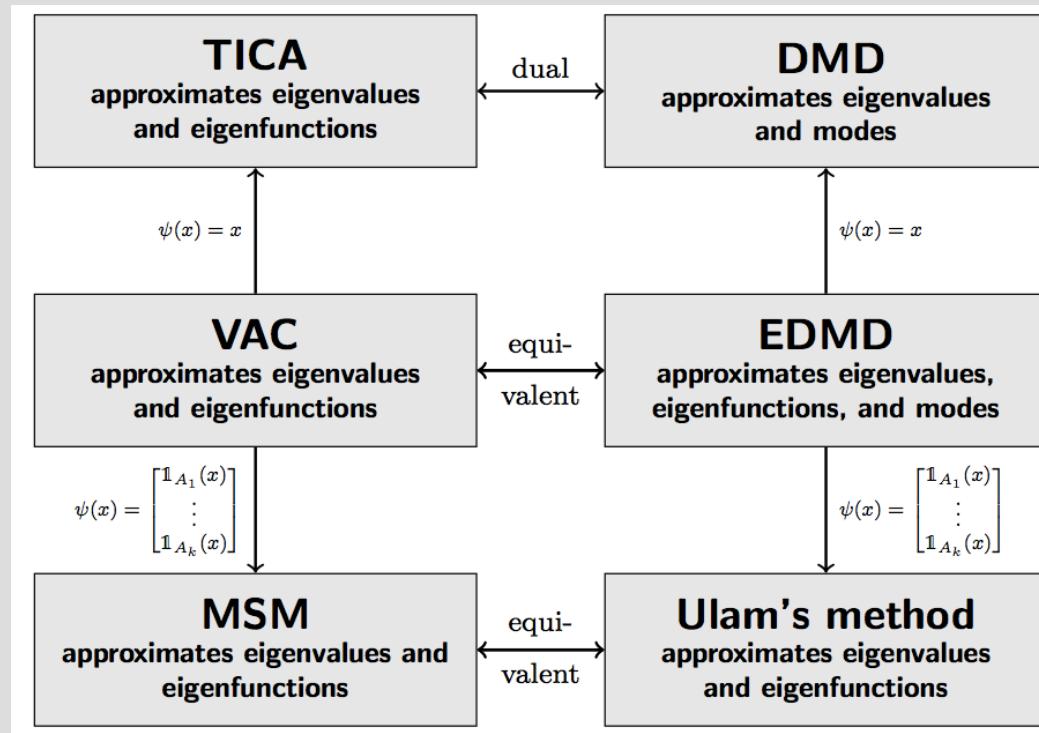
Nüske et al, **JCTC** 10, 1739-1752 (2014)



Comparison between methods

Molgedey and Schuster, **PRL** 72 3634-3637 (1994)
Pérez-Hernández et al, **JCP** 139, 015102 (2013)

Schmidt, Sesterhenn,
Ann. Meet. APS Div. Fluid Mech. (2008)



Schütte et al: **J. Comput. Phys.** (1999)
also: Noé, Pande, Hummer, Weber, Swope, ...

Klus, Nüske, Koltai, Wu, Krevrekidis, Schütte, Noé: Data-driven model reduction and transfer operator approximation (**J. Nonlin. Sci.** 2018 / **arXiv:1703.10112**)

Variational approach for Markov processes (VAMP)

Koopman operator

$$\begin{aligned}\mathcal{K}_\tau f(x) &= \mathbb{E}[f(x_{t+\tau}) | x_t = x] \\ &= \int p_\tau(x, y) f(y) dy\end{aligned}$$

Wu and Noé, arXiv:1707.04659 (2017)



Variational approach for Markov processes (VAMP)

Koopman operator

$$\begin{aligned}\mathcal{K}_\tau f(x) &= \mathbb{E}[f(x_{t+\tau}) | x_t = x] \\ &= \int p_\tau(x, y) f(y) dy\end{aligned}$$

Singular value decomposition:

$$\mathcal{K}_\tau f = \sum_i \sigma_i \langle \phi_i, f \rangle_{\rho_1} \psi_i$$

Wu and Noé, arXiv:1707.04659 (2017)



Variational approach for Markov processes (VAMP)

Koopman operator

$$\begin{aligned}\mathcal{K}_\tau f(x) &= \mathbb{E}[f(x_{t+\tau}) | x_t = x] \\ &= \int p_\tau(x, y) f(y) dy\end{aligned}$$

Singular value decomposition:

$$\mathcal{K}_\tau f = \sum_i \sigma_i \langle \phi_i, f \rangle_{\rho_1} \psi_i$$

- ρ_0, ρ_1 : empirical distribution of $x_t, x_{t+\tau}$
- If data are in equilibrium: stationary distribution $\mu = \rho_0 = \rho_1$
- $\{\phi_i\}$ and $\{\psi_i\}$ are both orthonormal bases with respect to $\langle \cdot, \cdot \rangle_{\rho_1}$ and $\langle \cdot, \cdot \rangle_{\rho_0}$,
- σ_i denotes the i th largest singular value.

Wu and Noé, arXiv:1707.04659 (2017)



Variational approach for Markov processes (VAMP)

Theorem **VAMP variational principle.** *The k dominant singular components of a Koopman operator are the solution of the following maximization problem:*

$$\begin{aligned} \sum_{i=1}^k \sigma_i^r &= \max_{\mathbf{f}, \mathbf{g}} \mathcal{R}_r [\mathbf{f}, \mathbf{g}], \\ \text{s.t. } \langle f_i, f_j \rangle_{\rho_0} &= 1_{i=j}, \\ \langle g_i, g_j \rangle_{\rho_1} &= 1_{i=j}, \end{aligned} \tag{10}$$

where $r \geq 1$ can be any positive integer. The maximal value is achieved by the singular functions $f_i = \psi_i$ and $g_i = \phi_i$ and

$$\mathcal{R}_r [\mathbf{f}, \mathbf{g}] = \sum_{i=1}^k \langle f_i, \mathcal{K}_\tau g_i \rangle_{\rho_0}^r \tag{11}$$

is called the VAMP- r score of \mathbf{f} and \mathbf{g} .

Wu and Noé, arXiv:1707.04659 (2017)



Implementation: time-lagged canonical covariance analysis (TCCA)

1. Compute

$$\mathbf{C}_{00} = \frac{1}{T-\tau} \mathbf{X}^{\top} \mathbf{X}$$

$$\mathbf{C}_{01} = \frac{1}{T-\tau} \mathbf{X}^{\top} \mathbf{Y}$$

$$\mathbf{C}_{11} = \frac{1}{T-\tau} \mathbf{Y}^{\top} \mathbf{Y}$$

with

$$\mathbf{X} = (\chi_0(\mathbf{x}_1), \chi_0(\mathbf{x}_2), \dots, \chi_0(\mathbf{x}_{T-\tau}))^{\top}$$

$$\mathbf{Y} = (\chi_1(\mathbf{x}_{1+\tau}), \chi_1(\mathbf{x}_{2+\tau}), \dots, \chi_1(\mathbf{x}_T))^{\top}$$

Wu and Noé, arXiv:1707.04659 (2017)

Implementation: time-lagged canonical covariance analysis (TCCA)

1. Compute

$$\begin{aligned}\mathbf{C}_{00} &= \frac{1}{T-\tau} \mathbf{X}^{\top} \mathbf{X} \\ \mathbf{C}_{01} &= \frac{1}{T-\tau} \mathbf{X}^{\top} \mathbf{Y} \\ \mathbf{C}_{11} &= \frac{1}{T-\tau} \mathbf{Y}^{\top} \mathbf{Y}\end{aligned}$$

with

$$\begin{aligned}\mathbf{X} &= (\chi_0(\mathbf{x}_1), \chi_0(\mathbf{x}_2), \dots, \chi_0(\mathbf{x}_{T-\tau}))^{\top} \\ \mathbf{Y} &= (\chi_1(\mathbf{x}_{1+\tau}), \chi_1(\mathbf{x}_{2+\tau}), \dots, \chi_1(\mathbf{x}_T))^{\top}\end{aligned}$$

2. Perform the truncated SVD

$$\mathbf{C}_{00}^{-\frac{1}{2}} \mathbf{C}_{01} \mathbf{C}_{11}^{-\frac{1}{2}} \approx \mathbf{U}_k \hat{\Sigma}_k \mathbf{V}_k^{\top}$$

3. Output $\hat{\Sigma}_k$, $\psi = \mathbf{U}_k^{\top} \mathbf{C}_{00}^{-\frac{1}{2}} \chi_0$ and $\phi = \mathbf{V}_k^{\top} \mathbf{C}_{11}^{-\frac{1}{2}} \chi_1$

Wu and Noé, arXiv:1707.04659 (2017)



Implementation: time-lagged canonical covariance analysis (TCCA)

1. Compute

$$\begin{aligned}\mathbf{C}_{00} &= \frac{1}{T-\tau} \mathbf{X}^\top \mathbf{X} \\ \mathbf{C}_{01} &= \frac{1}{T-\tau} \mathbf{X}^\top \mathbf{Y} \\ \mathbf{C}_{11} &= \frac{1}{T-\tau} \mathbf{Y}^\top \mathbf{Y}\end{aligned}$$

with

$$\begin{aligned}\mathbf{X} &= (\chi_0(\mathbf{x}_1), \chi_0(\mathbf{x}_2), \dots, \chi_0(\mathbf{x}_{T-\tau}))^\top \\ \mathbf{Y} &= (\chi_1(\mathbf{x}_{1+\tau}), \chi_1(\mathbf{x}_{2+\tau}), \dots, \chi_1(\mathbf{x}_T))^\top\end{aligned}$$

2. Perform the truncated SVD

$$\mathbf{C}_{00}^{-\frac{1}{2}} \mathbf{C}_{01} \mathbf{C}_{11}^{-\frac{1}{2}} \approx \mathbf{U}_k \hat{\Sigma}_k \mathbf{V}_k^\top$$

3. Output $\hat{\Sigma}_k$, $\psi = \mathbf{U}_k^\top \mathbf{C}_{00}^{-\frac{1}{2}} \chi_0$ and $\phi = \mathbf{V}_k^\top \mathbf{C}_{11}^{-\frac{1}{2}} \chi_1$

For the choice $\chi_0 = \chi_1$, TCCA is consistent with EDMD:

$$\mathbf{K}_\tau^\top = \mathbf{C}_{01}^\top \mathbf{C}_{00}^{-1}$$

Wu and Noé, arXiv:1707.04659 (2017)

VAMP: Koopman approximation error

Theorem (Koopman approximation error): For an operator $\hat{\mathcal{K}}_\tau$ defined by

$$\hat{\mathcal{K}}_\tau h = \sum_i \hat{\sigma}_i \langle g_i, h \rangle_{\rho_1} f_i$$

we have

$$\left\| \hat{\mathcal{K}}_\tau - \mathcal{K}_\tau \right\|_{\text{HS}}^2 = \text{tr} \left[\hat{\Sigma} \mathbf{C}_{00} \hat{\Sigma} \mathbf{C}_{11} - 2 \hat{\Sigma} \mathbf{C}_{01} \right] + \sum_i \sigma_i^2$$

- R_E (VAMP-E score)

where $\|\cdot\|_{\text{HS}}$ denotes the Hilbert-Schmidt norm, $\hat{\Sigma} = \text{diag}(\hat{\sigma}_1, \hat{\sigma}_2, \dots)$, and

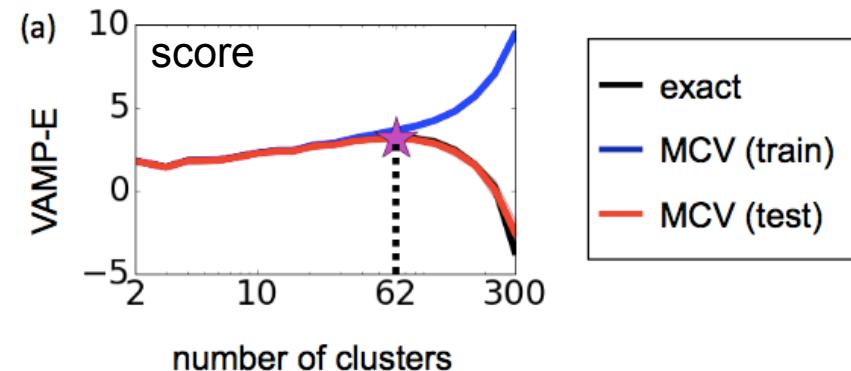
$$\begin{aligned} [\mathbf{C}_{00}]_{ij} &= \mathbb{E}_{\rho_0} [f_i(\mathbf{x}(t)) f_j(\mathbf{x}(t))] \\ [\mathbf{C}_{01}]_{ij} &= \mathbb{E}_{\rho_0} [f_i(\mathbf{x}(t)) g_j(\mathbf{x}(t+\tau))^\top] \\ [\mathbf{C}_{11}]_{ij} &= \mathbb{E}_{\rho_1} [g_i(\mathbf{x}(t)) g_j(\mathbf{x}(t))^\top]. \end{aligned}$$

Wu and Noé, arXiv:1707.04659 (2017)

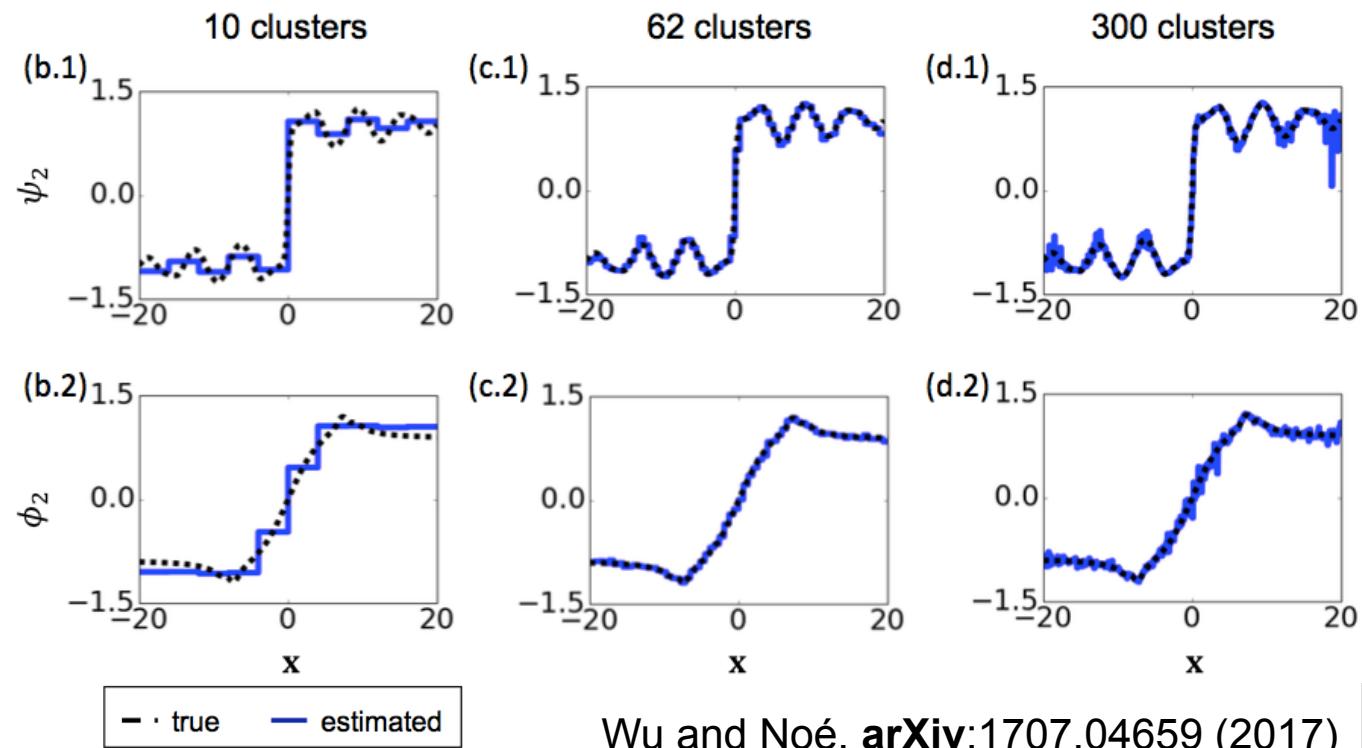
1D-Example

Markov process with Gaussian noise w_t

$$x_{t+1} = \frac{x_t}{2} + \frac{25x_t}{1+x_t^2} + \sqrt{10}(1.1 + \cos(x_t))w_t$$



Eigenfunctions

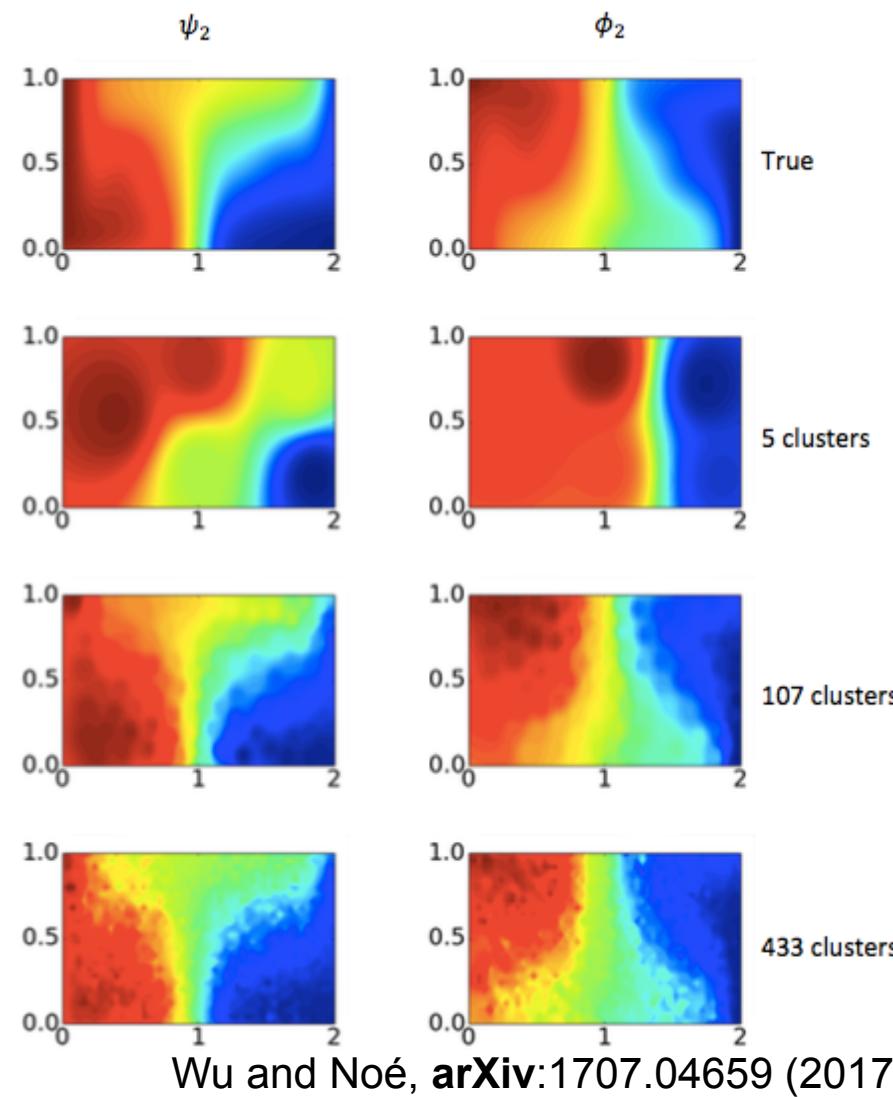
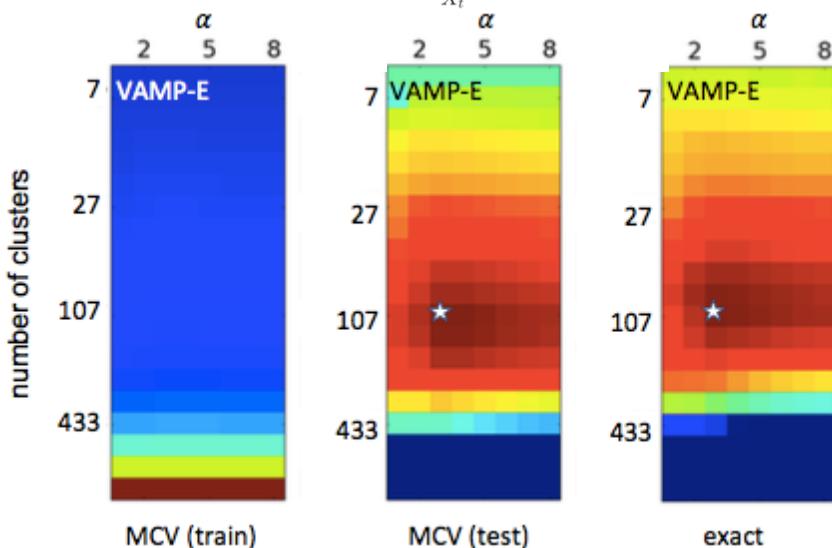
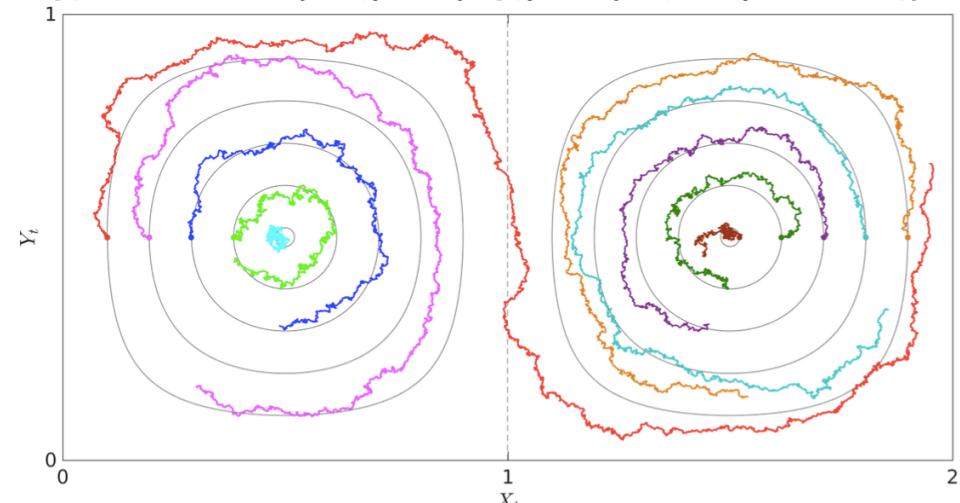


Wu and Noé, arXiv:1707.04659 (2017)

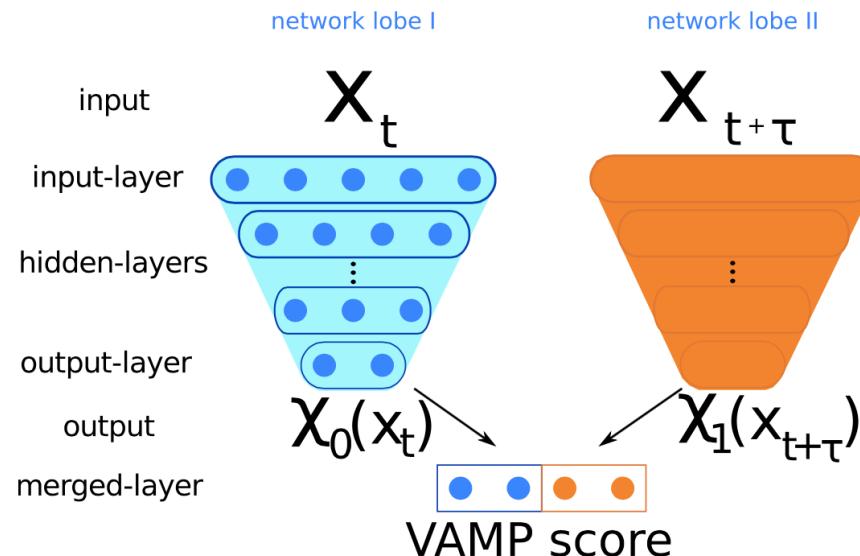
Generalization: VAMP reduces nonequilibrium processes

$$dx_t = -\pi A \sin(\pi x_t) \cos(\pi y_t) - \epsilon(2y_t - 1) + \varepsilon d\mathbf{W}_{t,1},$$

$$dy_t = \pi A \cos(\pi x_t) \sin(\pi y_t) - \epsilon(2x_t - 3) + \varepsilon d\mathbf{W}_{t,2}$$



Wu and Noé, arXiv:1707.04659 (2017)



VAMP variational principle (subspace version)

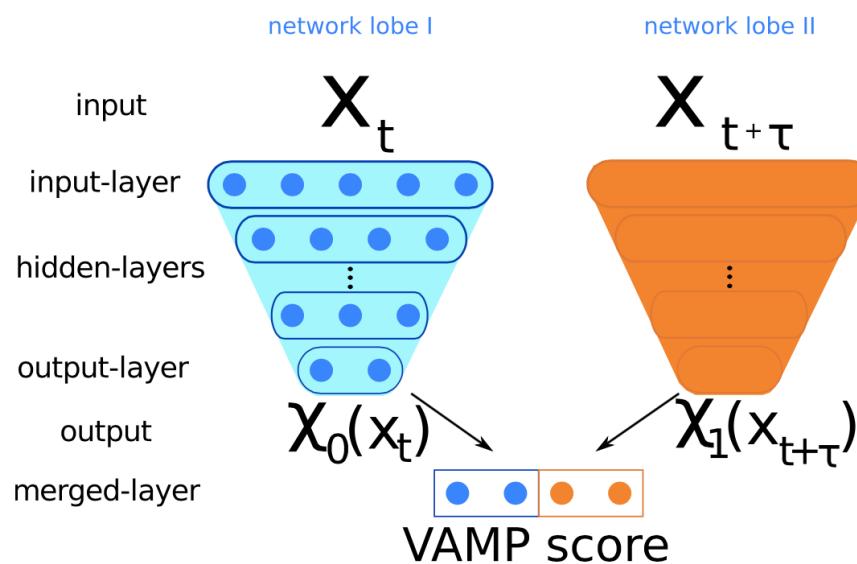
For any two sets of linearly independent functions $\chi_0(\mathbf{x}) = (\chi_{01}(\mathbf{x}), \dots, \chi_{0n}(\mathbf{x}))$ and $\chi_1(\mathbf{x}) = (\chi_{11}(\mathbf{x}), \dots, \chi_{1n}(\mathbf{x}))$, let us call

$$\hat{R}_2[\chi_0, \chi_1] = \left\| \mathbf{C}_{00}^{-\frac{1}{2}} \mathbf{C}_{01} \mathbf{C}_{11}^{-\frac{1}{2}} \right\|_F^2$$

their VAMP-2 score, where \mathbf{C}_{00} , \mathbf{C}_{01} , \mathbf{C}_{11} are the feature correlation matrices as defined earlier and $\|\cdot\|_F$ indicates the Frobenius norm. The maximum value of the VAMP-2 score is achieved when the top n left and right Koopman singular functions belong to $\text{span}(\chi_0)$ and $\text{span}(\chi_1)$, respectively.

Mardt, Pasquali, Wu, Noé **Nat. Commun.** 9, 5 (2018)

VAMPnets



Resulting Markov model:

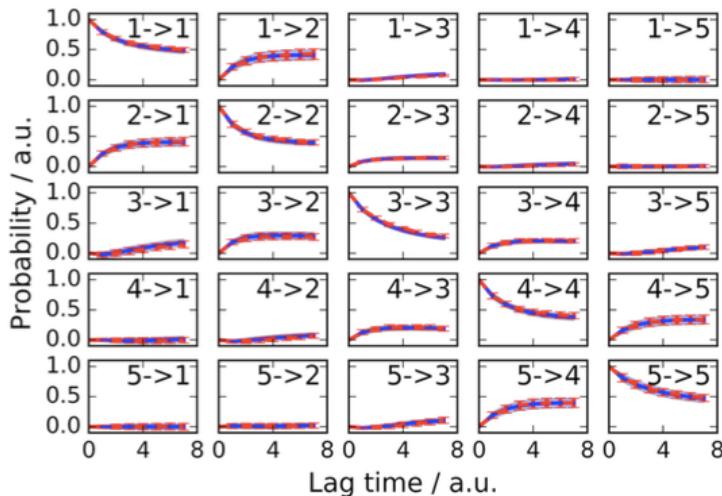
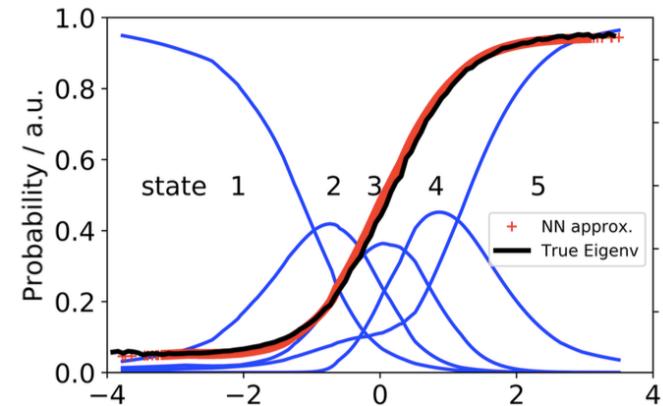
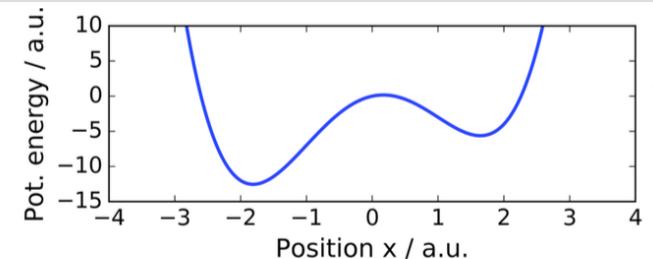
$$\mathbf{K} = \mathbf{C}_{00}^{-1} \mathbf{C}_{01}.$$

Relaxation timescales:

$$t_i(\tau) = -\frac{\tau}{\ln |\lambda_i(\tau)|},$$

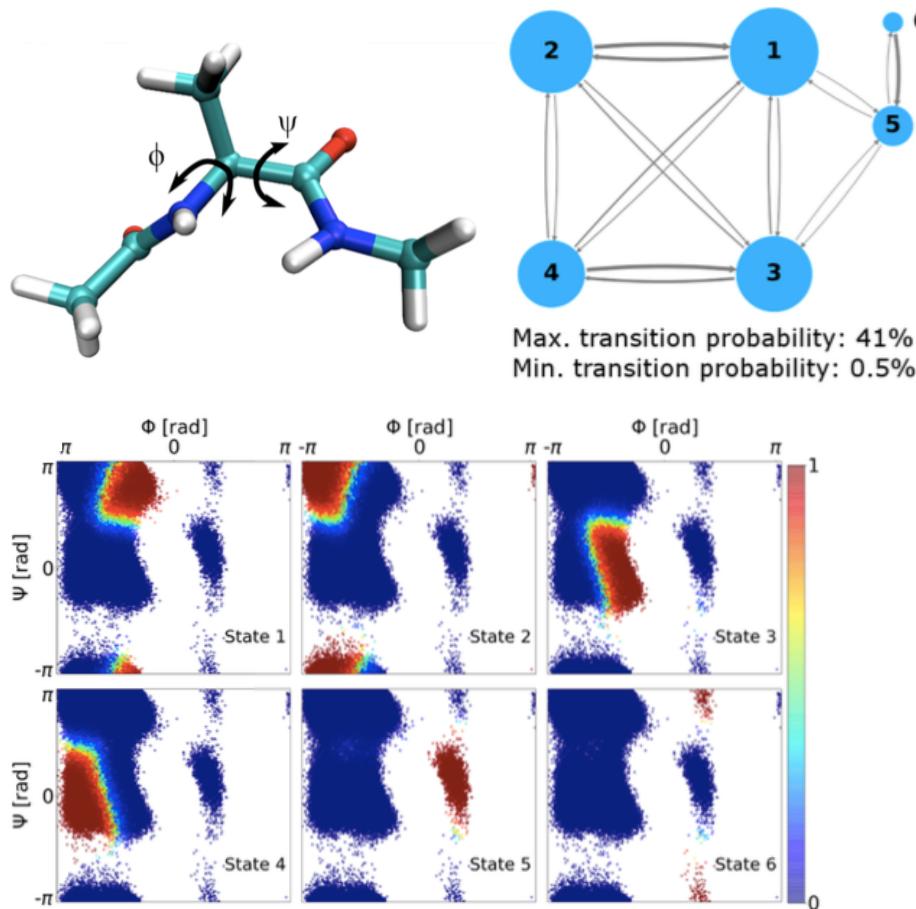
Validate (Chapman-Kolmogorov test):

$$\mathbf{K}(n\tau) = \mathbf{K}^n(\tau),$$

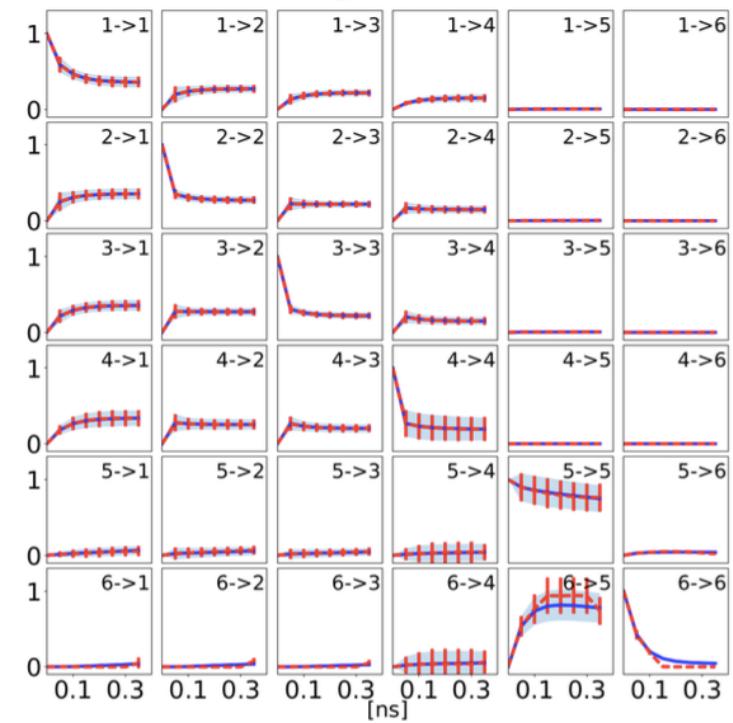
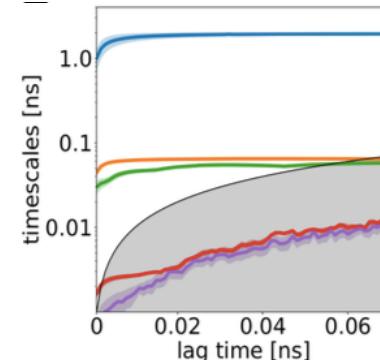


Mardt, Pasquali, Wu, Noé **Nat. Commun.** 9, 5 (2018)

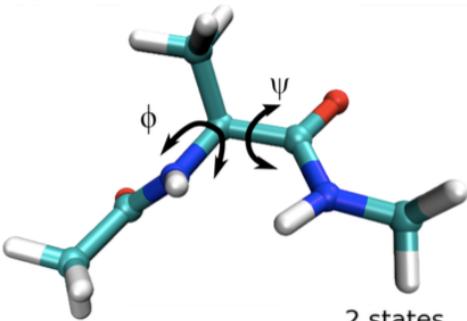
Alanine dipeptide



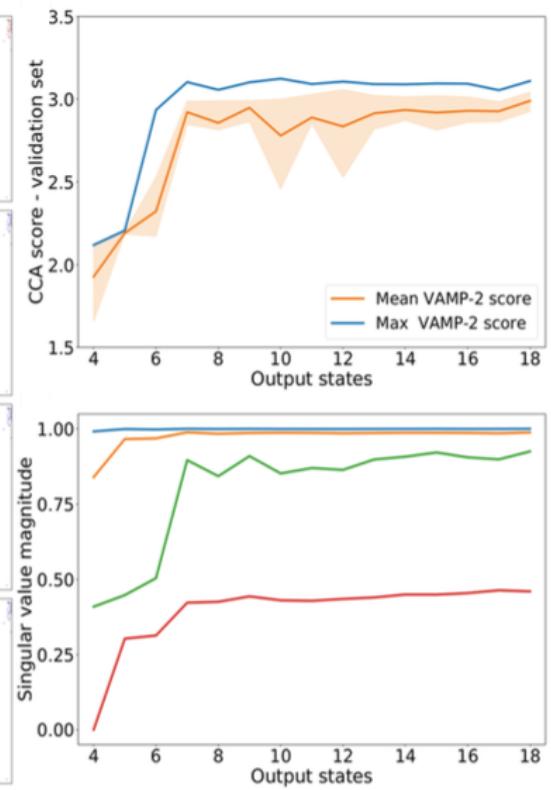
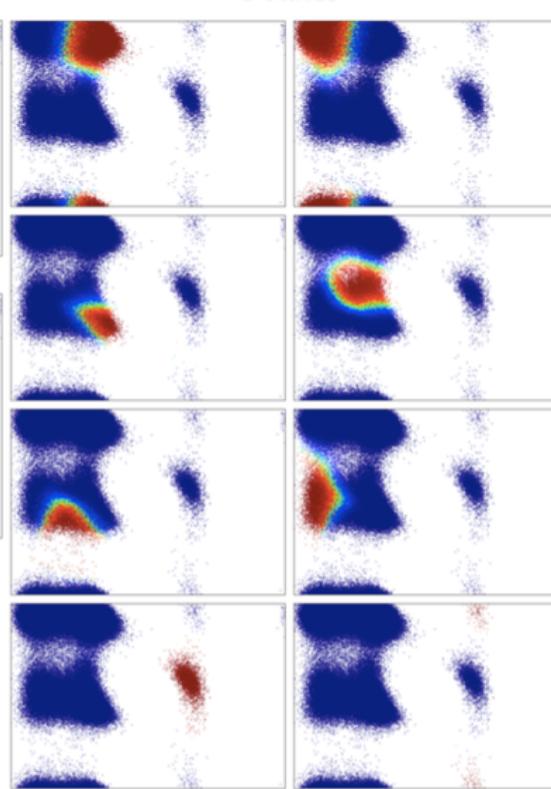
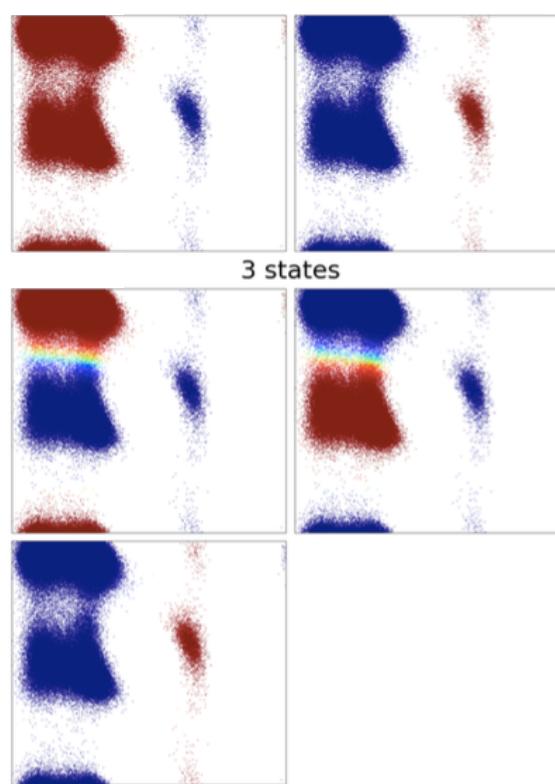
Validation



Mardt, Pasquali, Wu, Noé Nature Communications (2018)

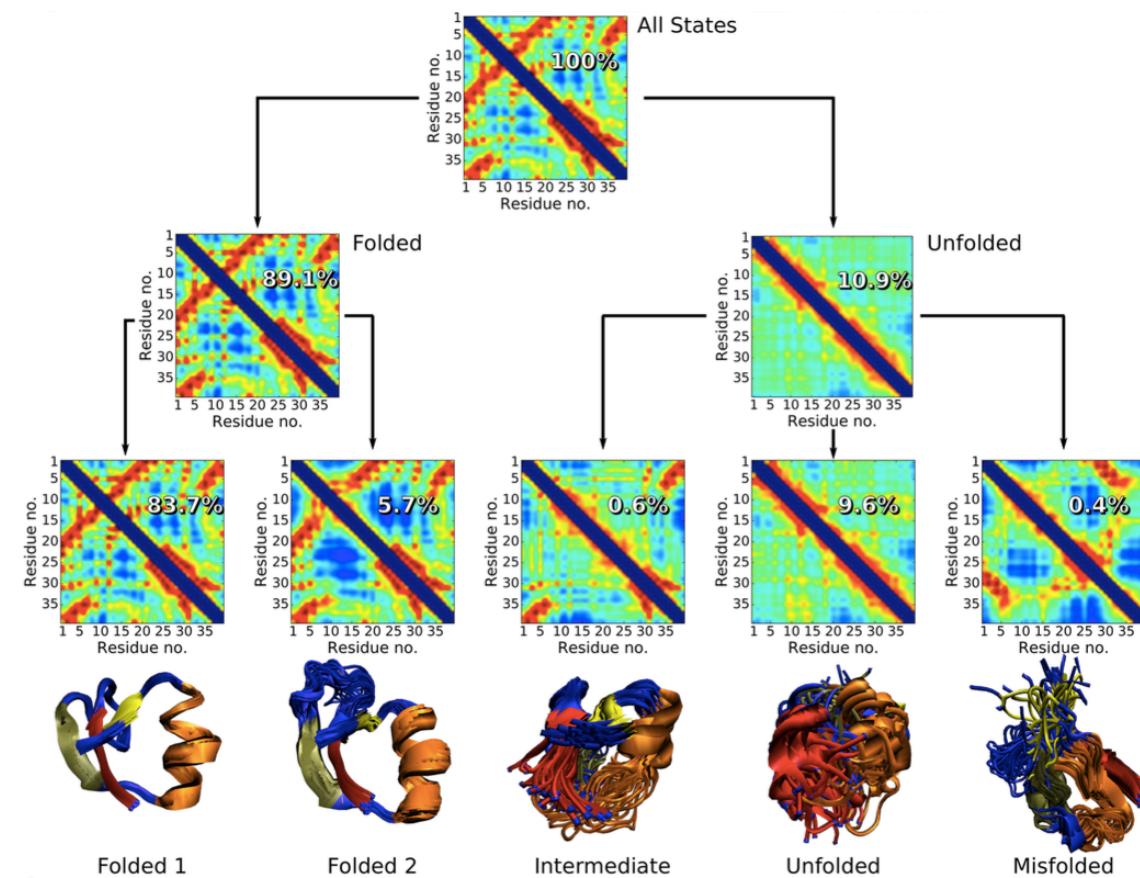


Results as a function of the number of states

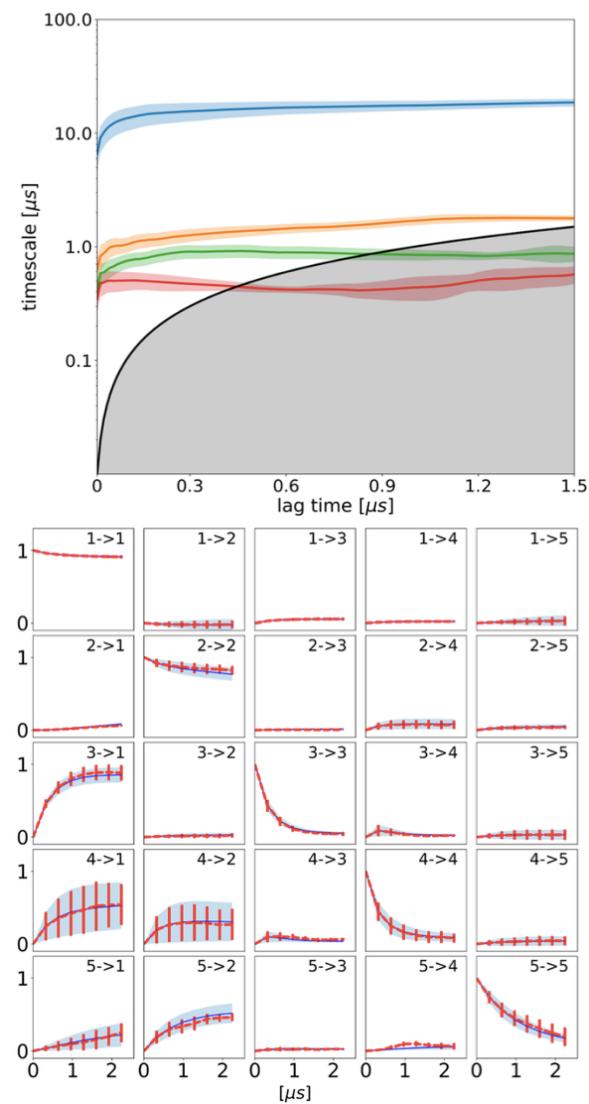


Mardt, Pasquali, Wu, Noé **Nature Communications** (2018)

NTL9 Protein folding

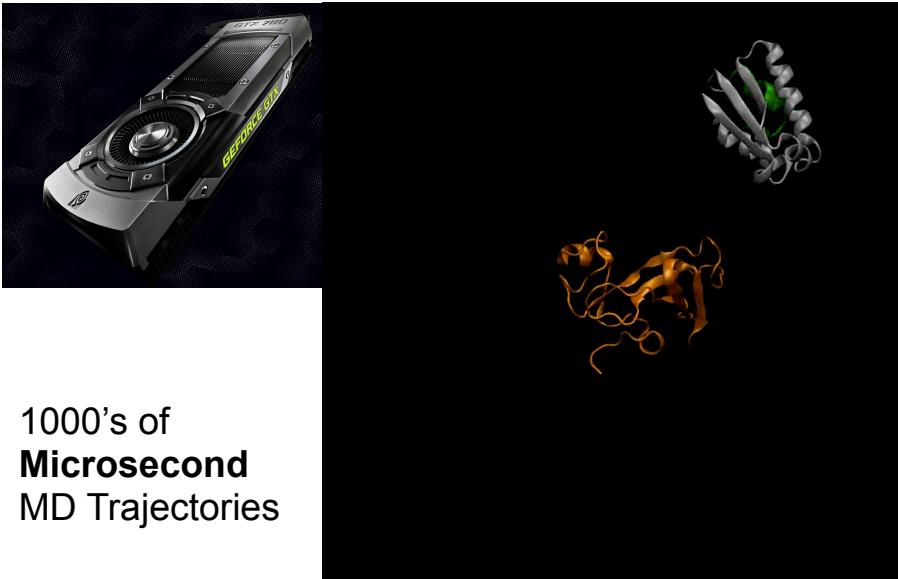


Validation

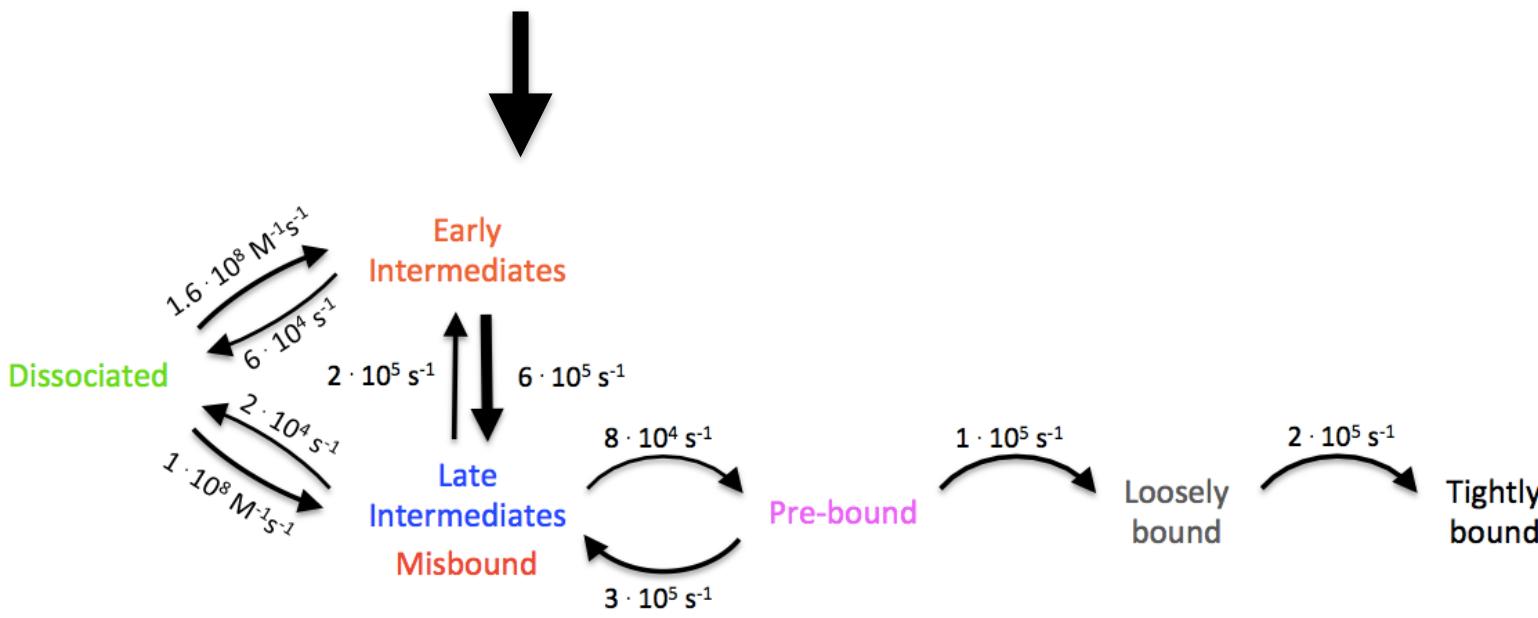


Mardt, Pasquali, Wu, Noé Nature Communications (2018)

Simulating biological timescales at atomic resolution

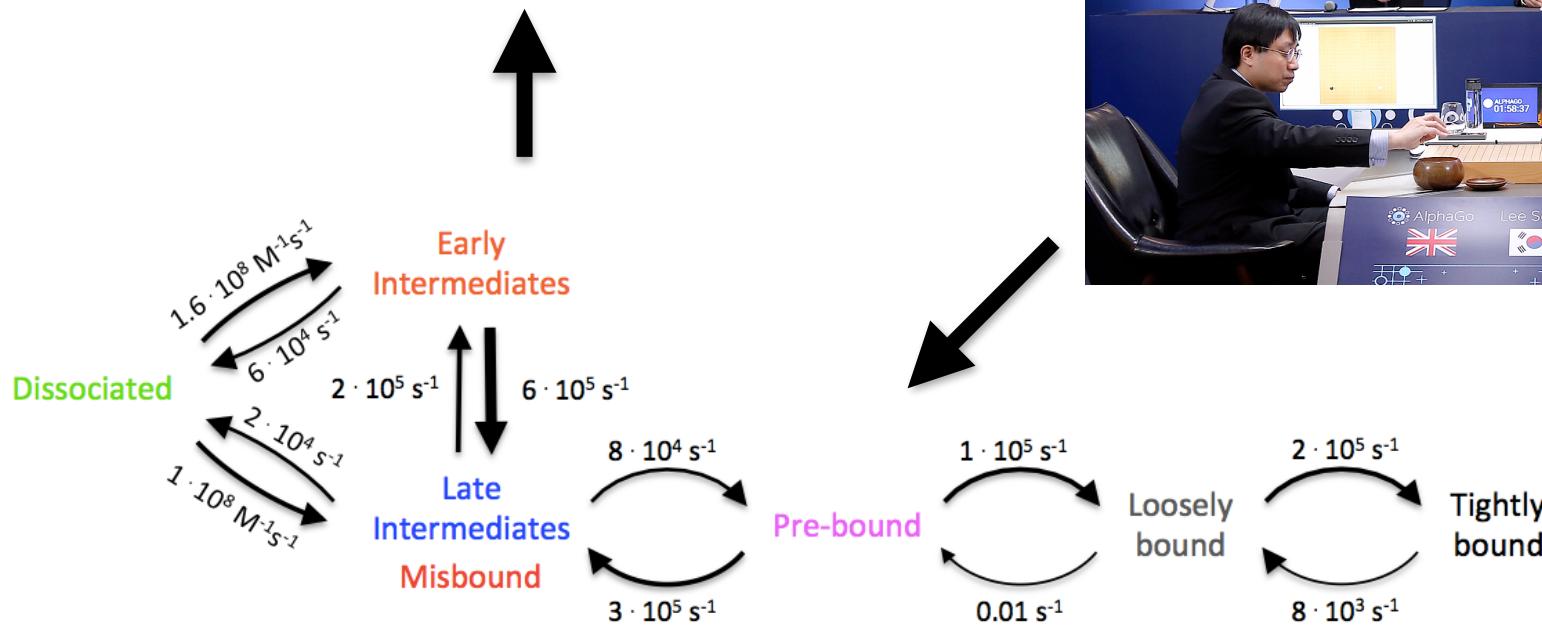


1000's of
Microsecond
MD Trajectories



Markov State Model — **Millisecond kinetics**

Simulating biological timescales at atomic resolution



Adaptive Markov State Model — seconds to hours kinetics



Sampling biological timescales at atomic resolution

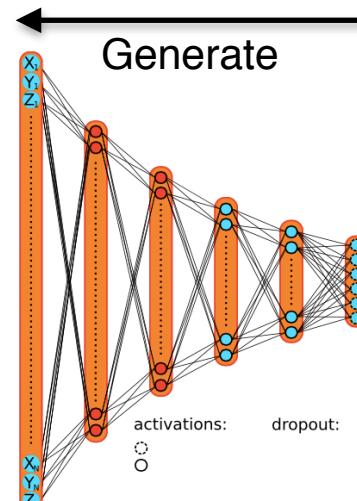
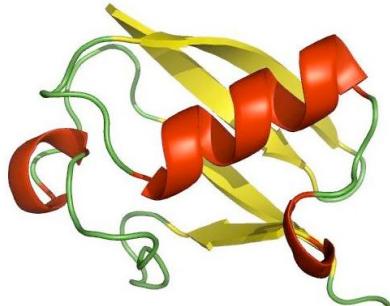


0.1 millisecond binding trajectory

Plattner, Doerr, De Fabritiis, Noé
Nature Chemistry 9, 1005 (2017)

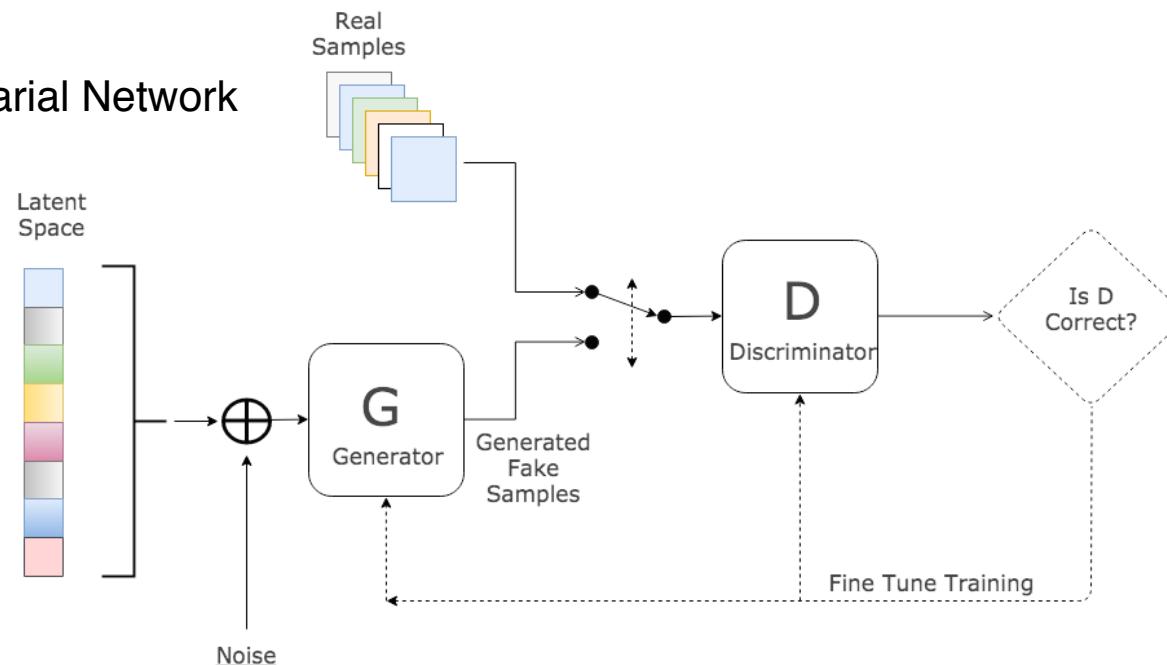
Generative learning and molecular design

Structure and dynamics
from MD Simulation data
Experimental data ...

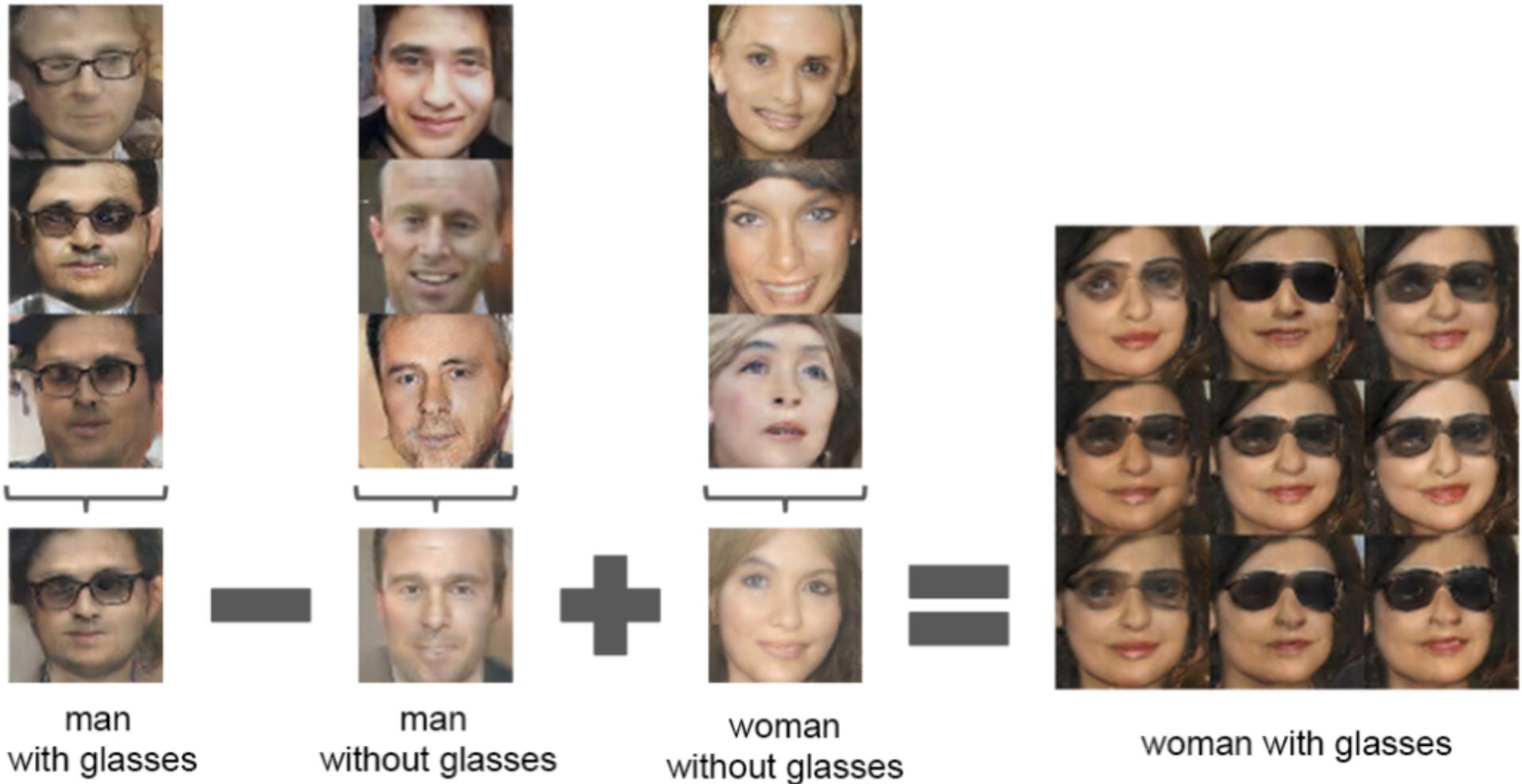


“Latent” variables
(states, kinetics, properties...)

Example:
Generative Adversarial Network

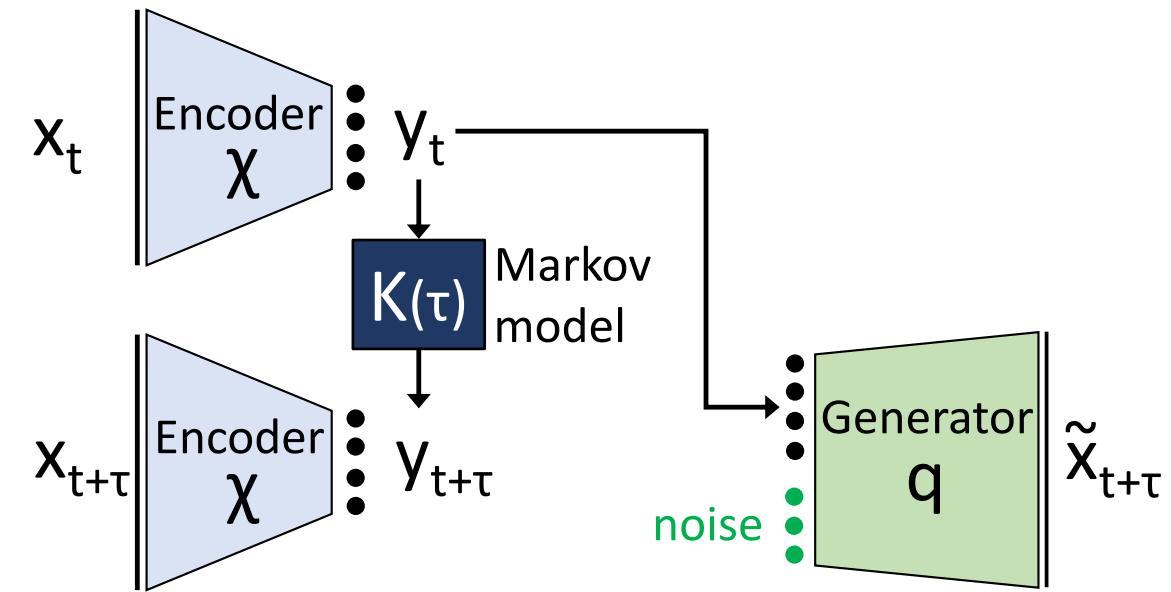


Generative Adversarial Network

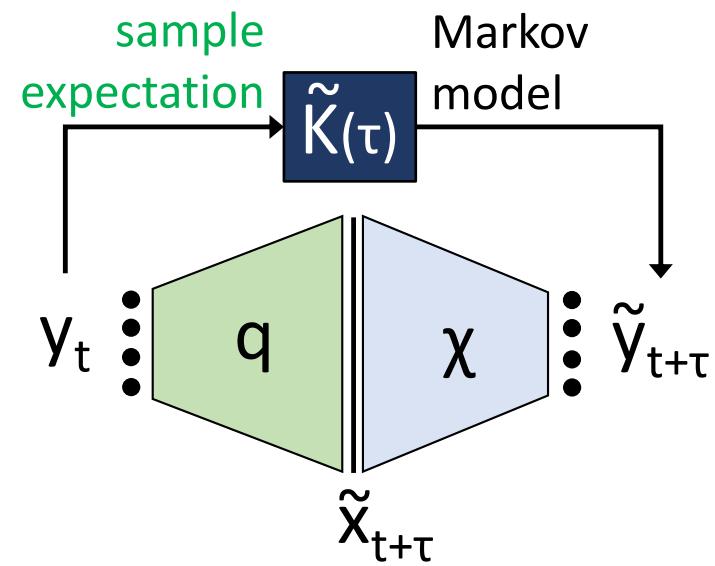


Radford et al: Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks arXiv:1511.06434 (2015)

Deep Generative MSM



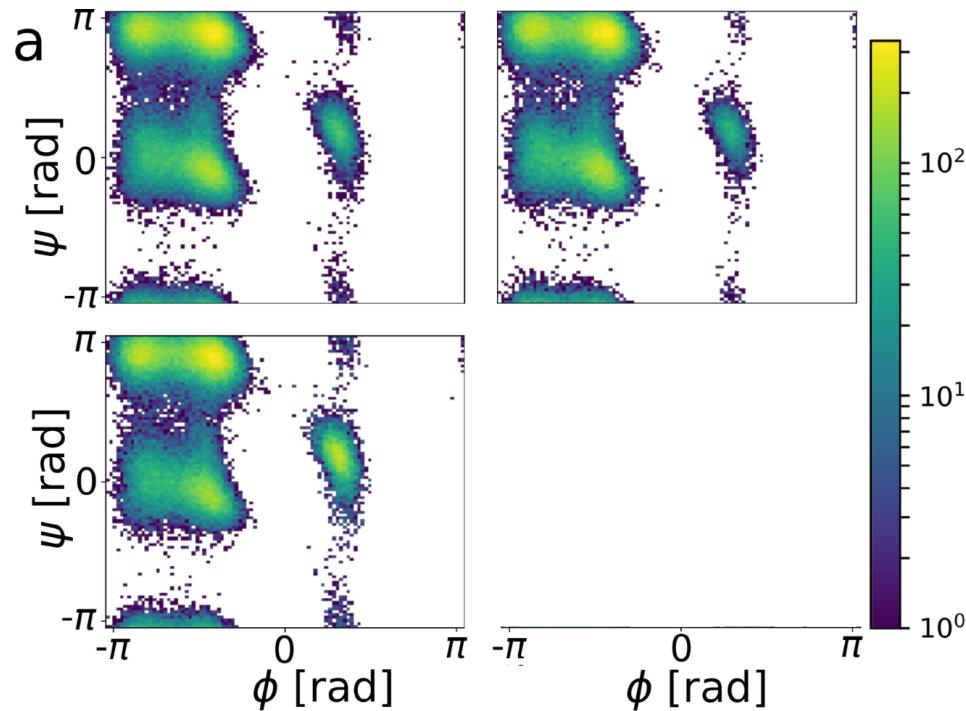
Rewiring Trick



Deep Generative Markov State Models

Data

Deep MSM,
resampled



“good”
classical MSM

Energy distance between distributions:

$$D_E (\mathbb{P}(x), \mathbb{P}(y)) = \mathbb{E} [2 \|x - y\| - \|x - x'\| - \|y - y'\|]$$

with

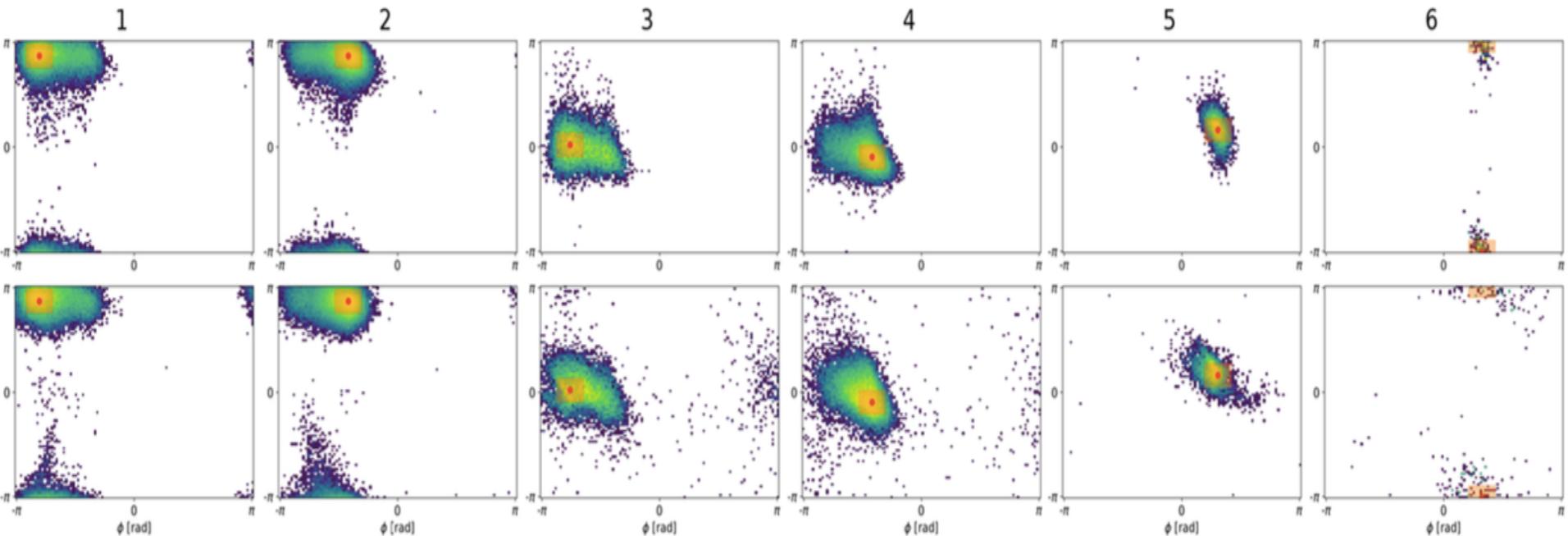
$$x, x' \sim \mathbb{P}(x)$$

$$y, y' \sim \mathbb{P}(y)$$

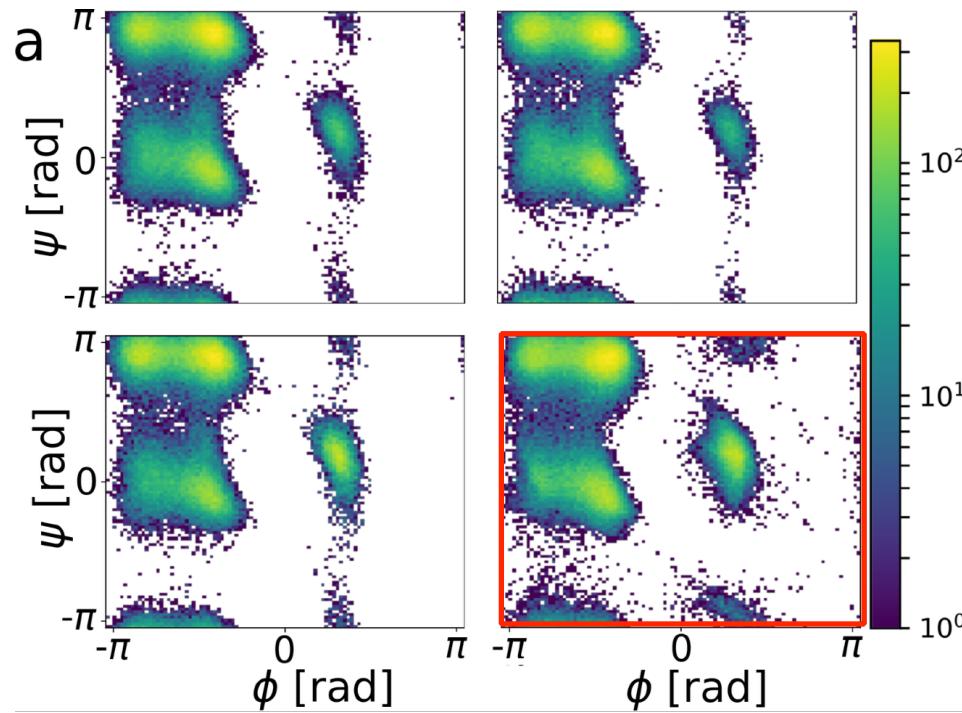
$D_E = 0$ only if distributions are equal

—> Train Generator Network by minimizing Energy Distance

Learning transition densities



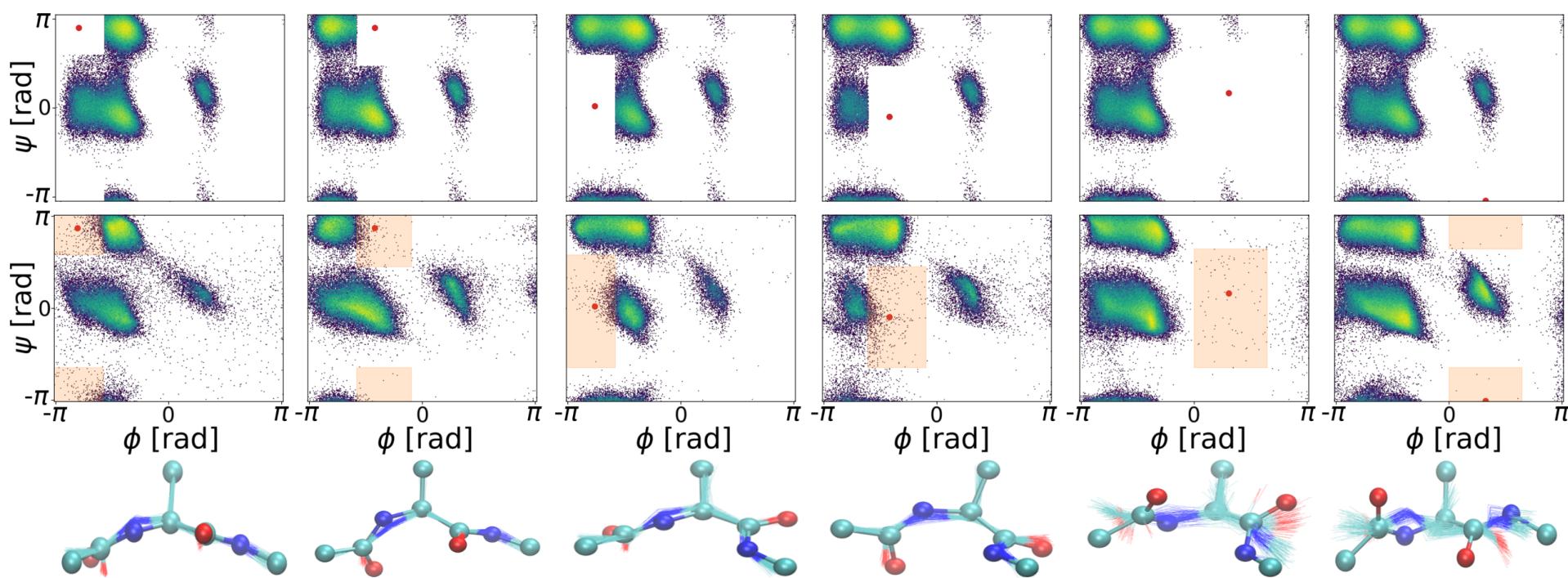
Data



“good”
classical MSM

Deep generative
MSM

Deep Generative Markov State Models



Acknowledgements



Collaborations

Cecilia Clementi (Rice University)
Christof Schütte (FU Berlin)
Eric Vanden-Eijnden (Courant Institut NY)
Thomas Weikl (MPI Potsdam)
Edina Rosta (King's College London)

Vijay Pande (Stanford)
Volker Haucke (FMP Berlin)
Stephan Sigrist (FU Berlin)
Katja, Faelber, Oliver Daumke (MDC)
John Chodera (MSKCC NY)
Gianni de Fabritiis (Barcelona)

Funding



Deutsche
Forschungsgemeinschaft