

Using seeding methods for efficient and accurate simulation of folding dynamics

Hongbin Wan

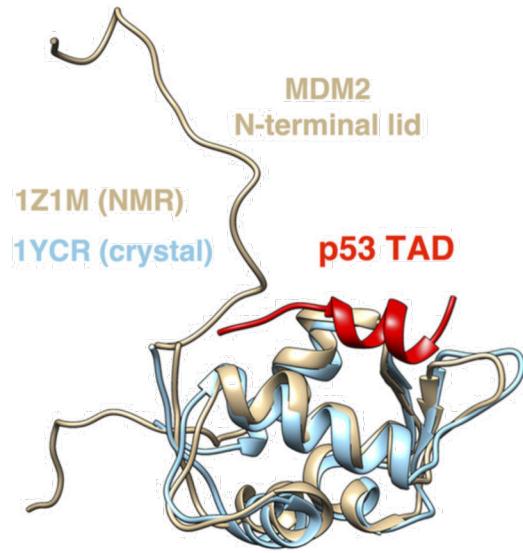
Vincent Voelz lab

Temple University

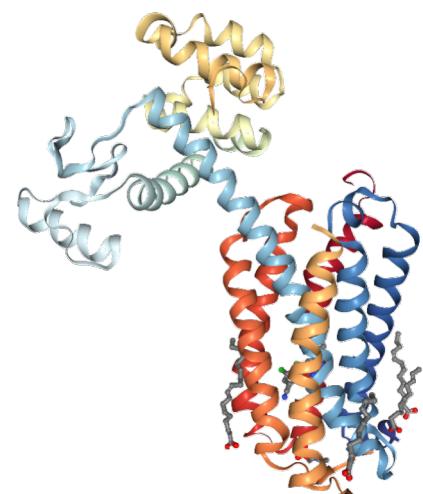
Protein Folding Consortium Workshop, 2018

Protein conformational dynamics are directly linked to protein function in many important biological processes

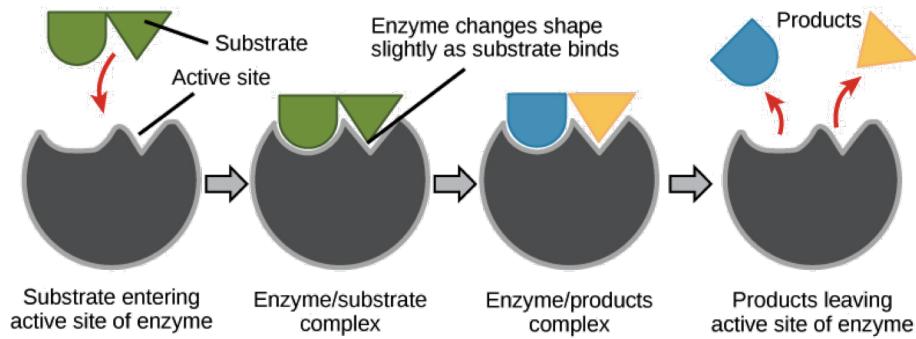
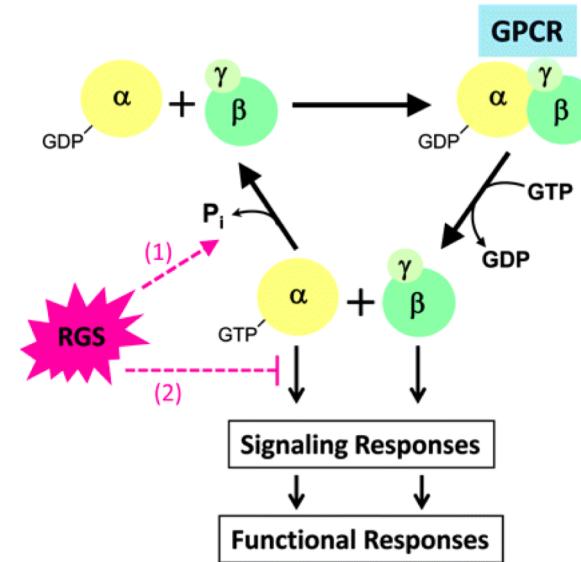
- protein-protein interaction



- allosteric regulation

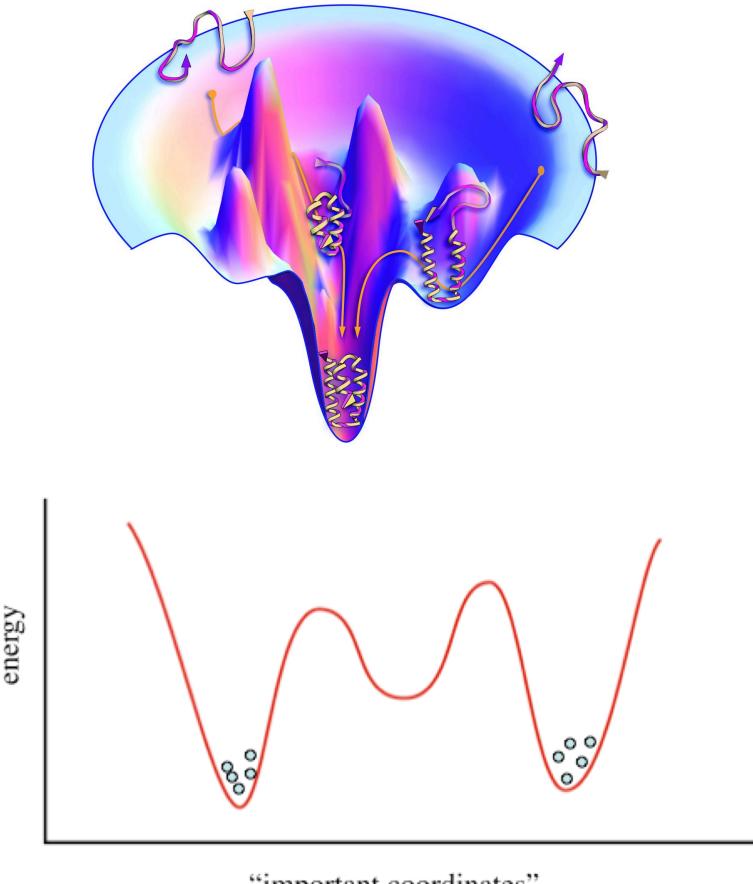


- Cell signaling



- Enzyme mechanism

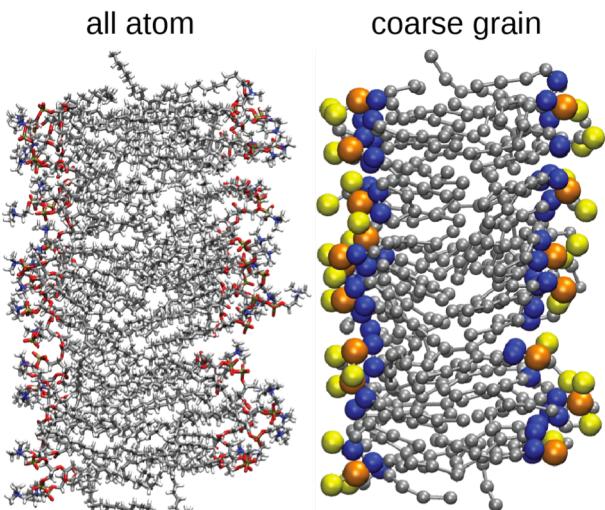
It is challenging to simulate conformational dynamics using MD



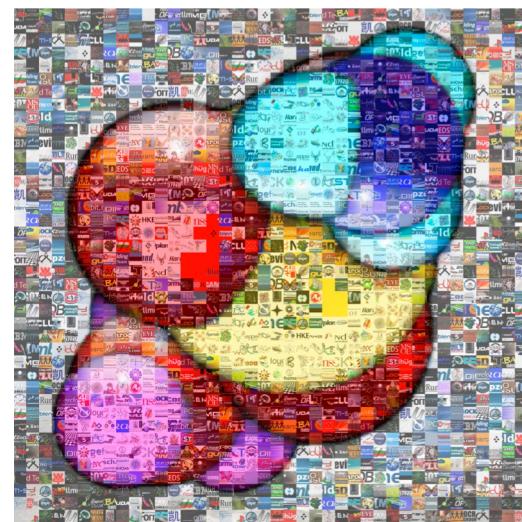
- Long timescale is needed to capture biologically important processes.
- Energy landscapes are rugged and high-dimensional.
- Important degrees of freedom are hard to define ahead of time.
- Needs lots of samplings (time, computational resources) to explore all reaction coordinates.

Q: how could we deal with these challenges?

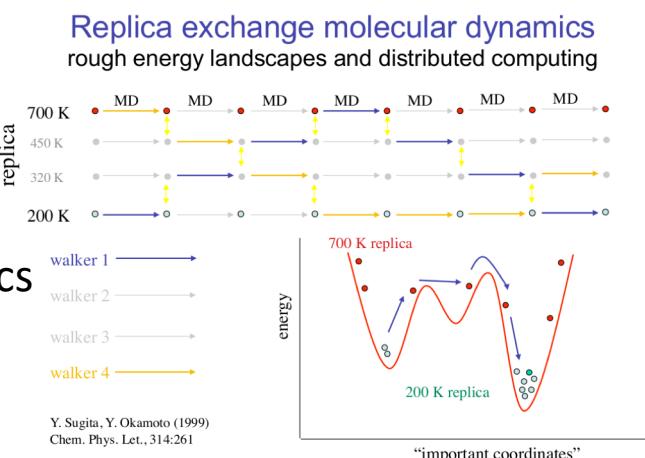
Many methods exist to address these sampling problems



Coarse-grain model



Special purpose hardware

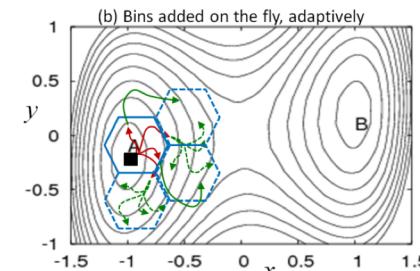
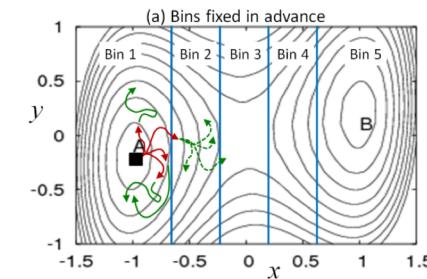


Replica-exchange

- Mainly thermodynamics

Weighted Ensemble

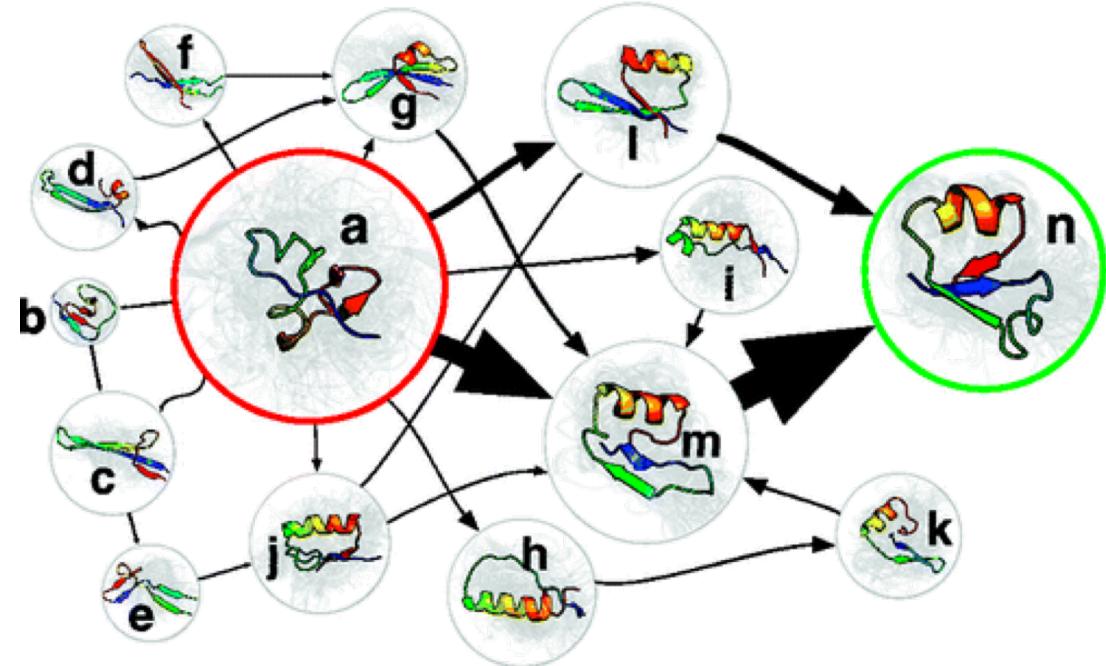
- Mainly kinetics



Enhanced sampling methods

Markov State models are a good way to capture both thermodynamics and kinetics

- Top 10 folding pathways of MSM of N-terminal domain of the crystal structure of ribosomal protein L9(NTL9)
- **MSMs estimate long time dynamics from large numbers of short trajectories (local equilibrium).**



Suppose we already have constructed an MSM.

Can we leverage this information to build new
MSMs?

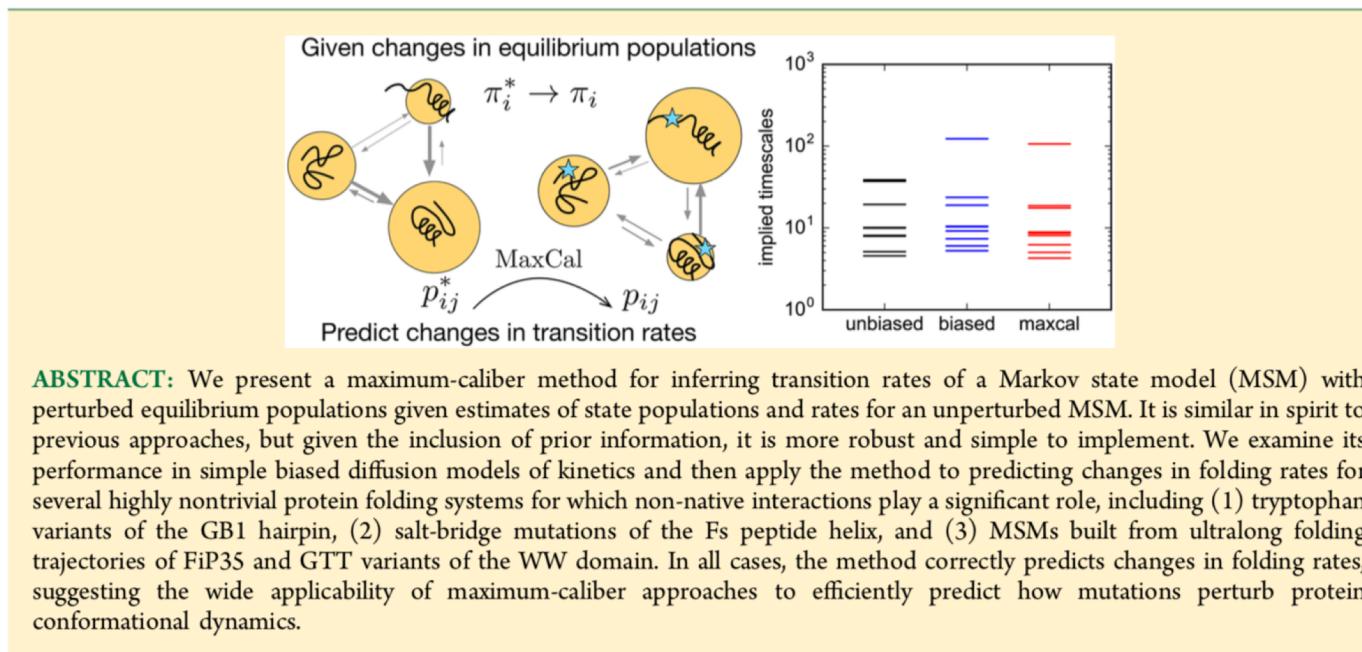
One approach: Using maximum caliber (MaxCal) methods to infer perturbation effects

A Maximum-Caliber Approach to Predicting Perturbed Folding Kinetics Due to Mutations

Hongbin Wan, Guangfeng Zhou, and Vincent A. Voelz*[†]

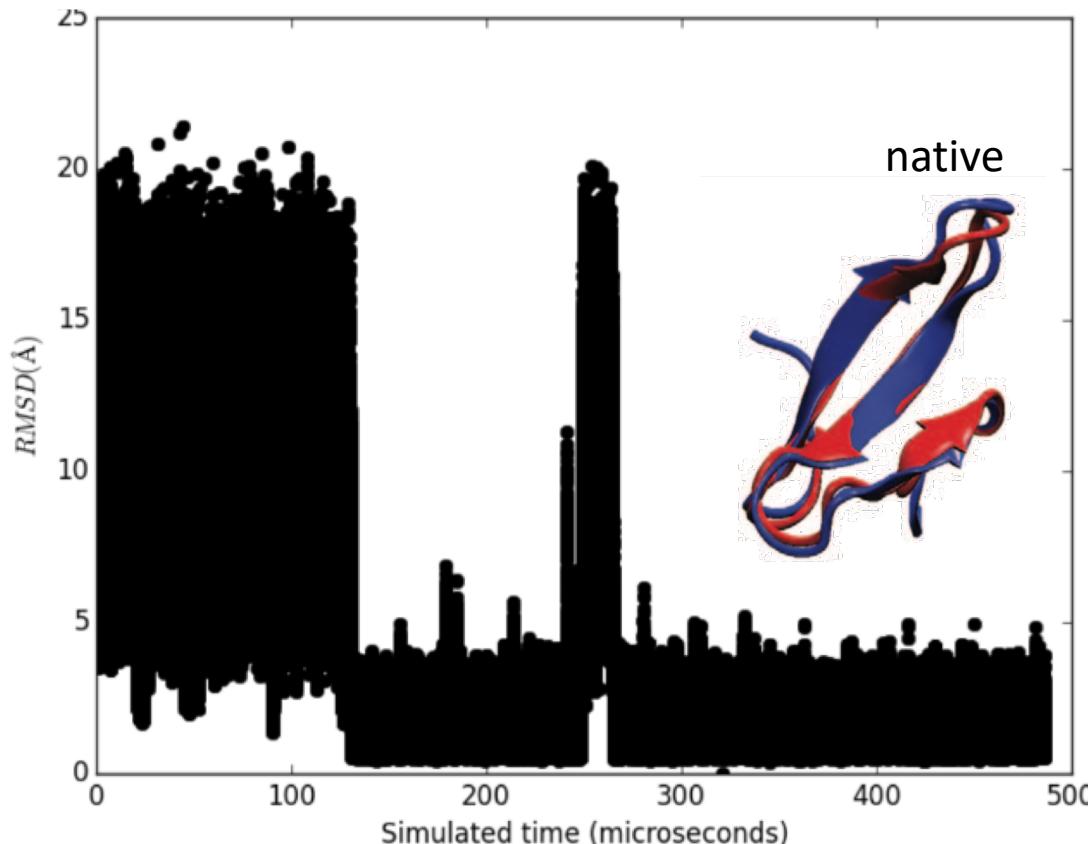
Department of Chemistry, Temple University, Philadelphia, Pennsylvania 19122, United States

 Supporting Information

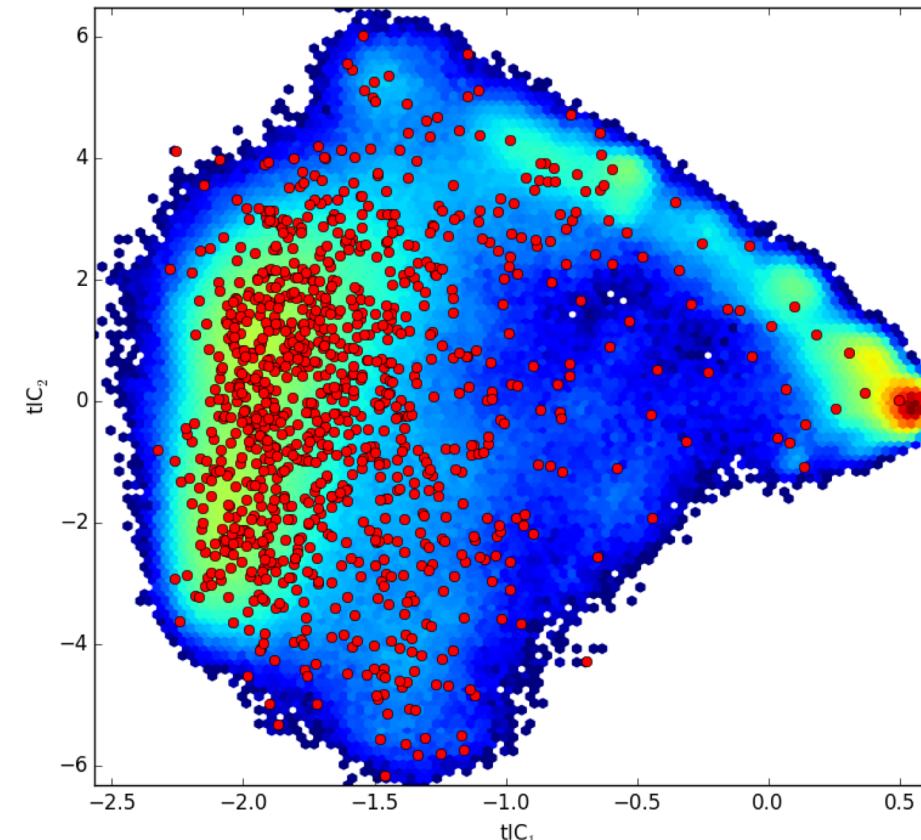


Another approach: adaptive seeding

Ultra-long trajectory of GTT ww-domain folding



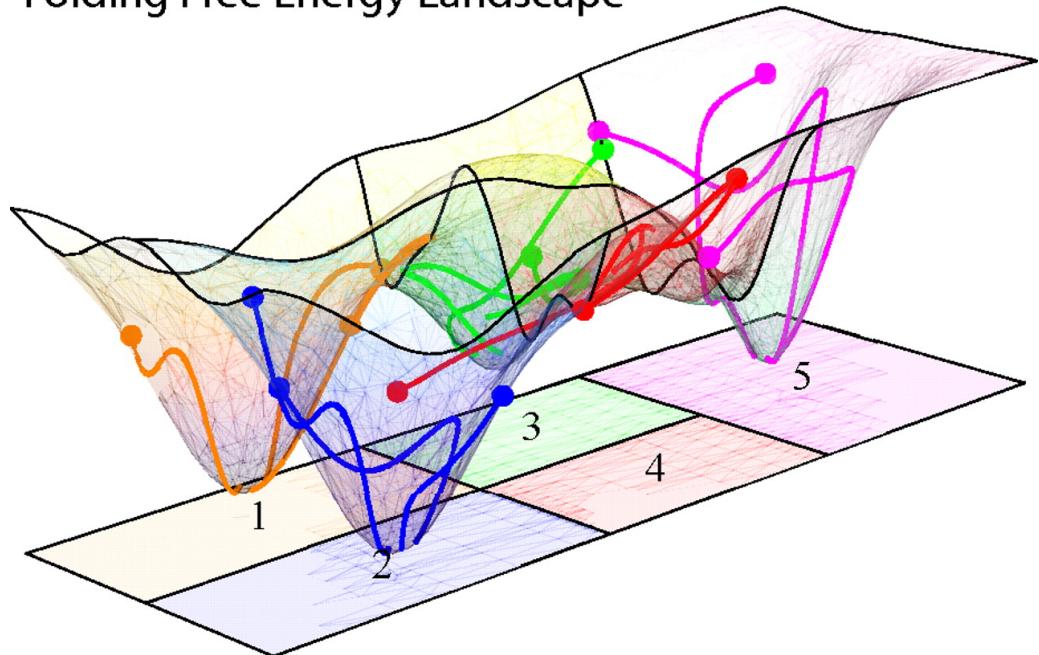
MSM built from the trajectory



- Long trajectories spend 100s of microseconds in the folded state.
- Can we seed ensembles of short-trajectories from each MSM microstate to sample transitions?

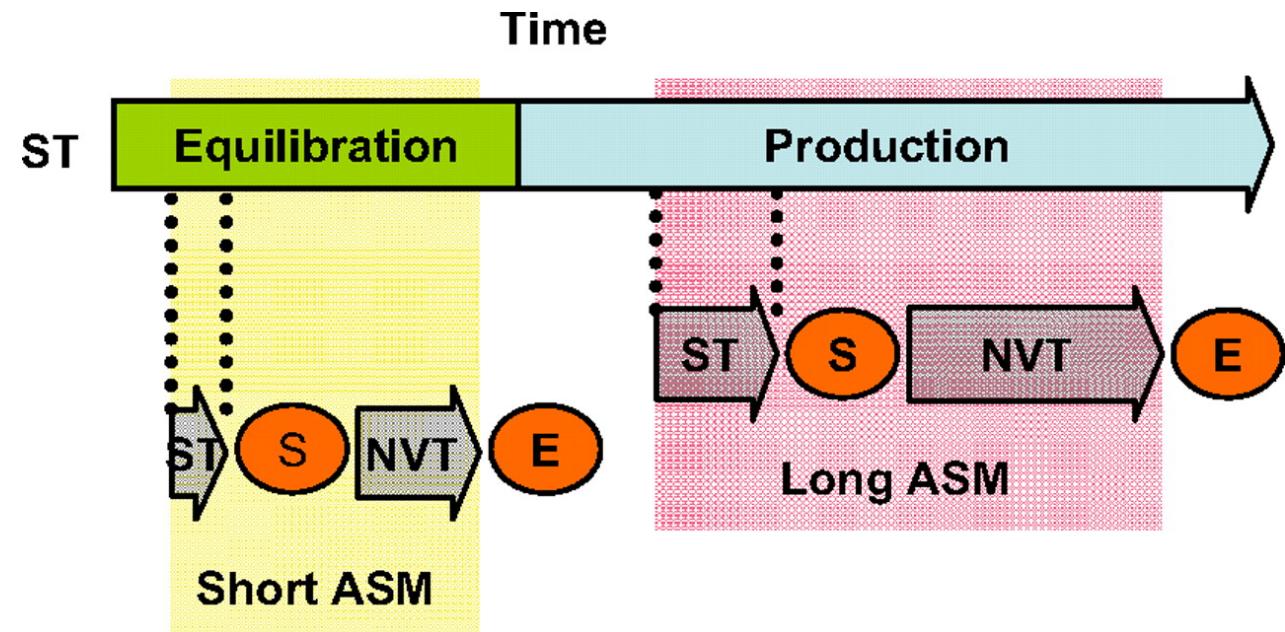
Adaptive seeding method (ASM)

Folding Free Energy Landscape



Metastable States
from Markov State Models

Schematic of the adaptive seeding scheme

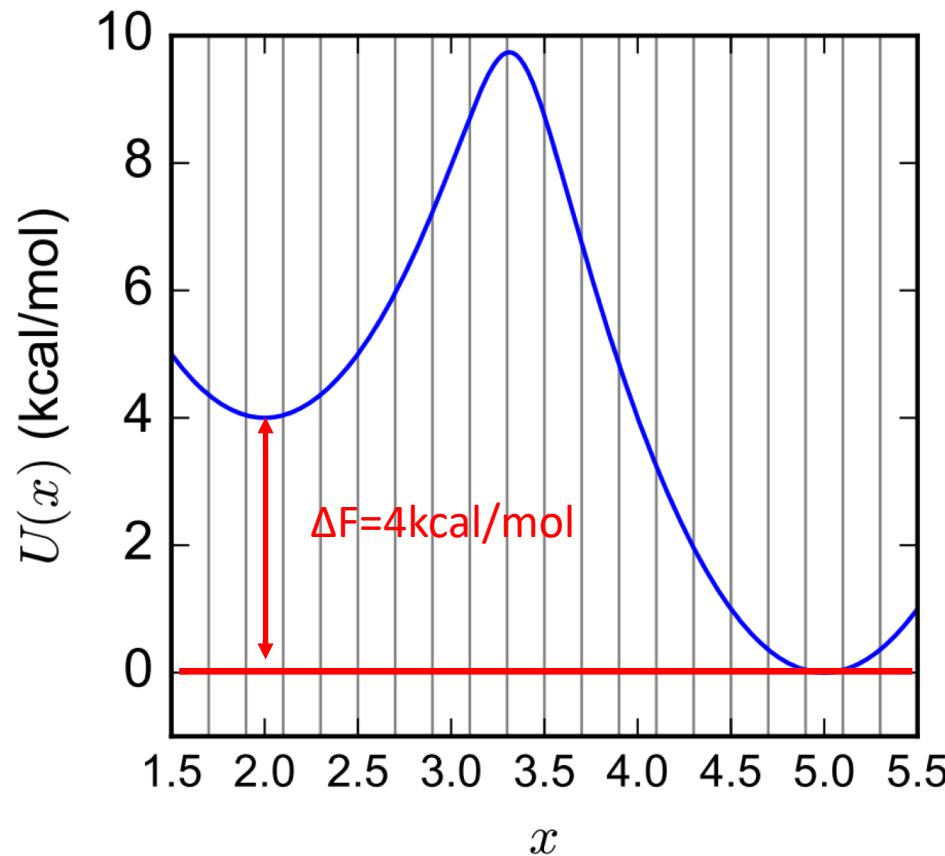


Adaptive seeding of a 1-D two-state potential

$$U(x) = -\frac{2k_B T}{0.596} \ln[e^{-2(x-2)^2/2} + e^{-2(x-5)^2/2}]$$

- Monte Carlo random-walk with step-size=0.05 simulations.
- Seeds: 20bins(microstates) * 20trajs/bin.
- Slowest relaxation time $\sim 5 \times 10^6$ steps

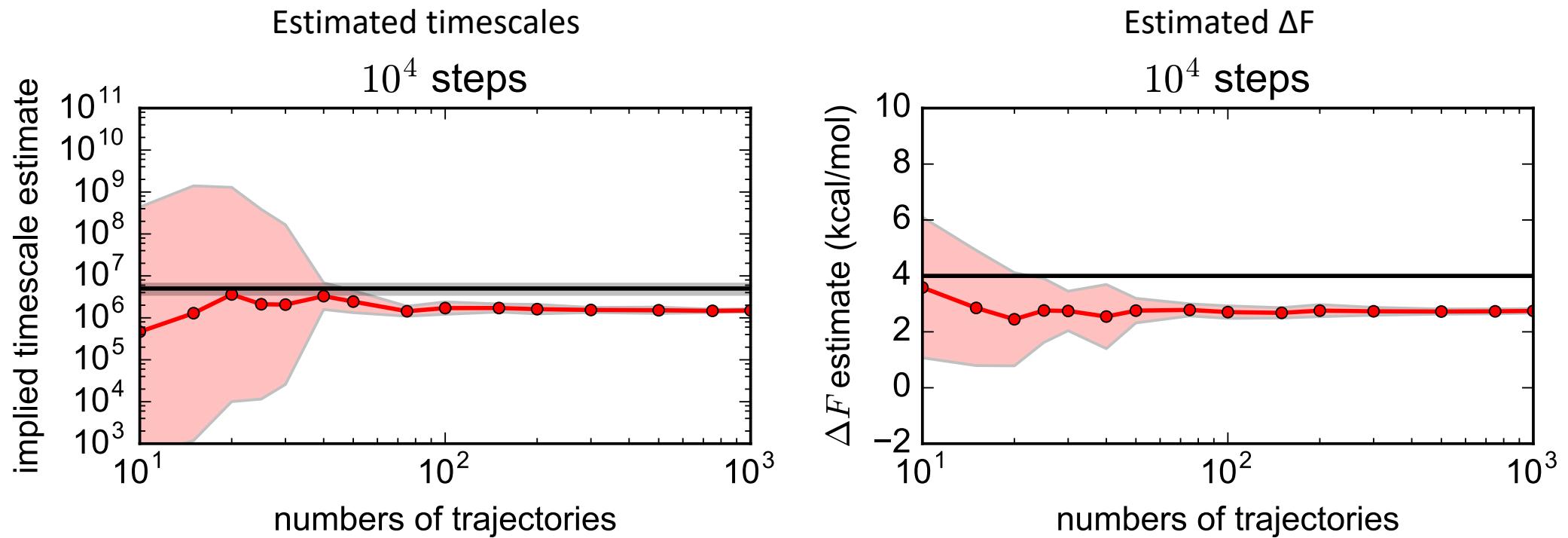
Shoot 20 trajectories from each bin



Problem

- The samples are not drawn from equilibrium
- But most MSM estimators assume that samples **ARE** drawn from equilibrium

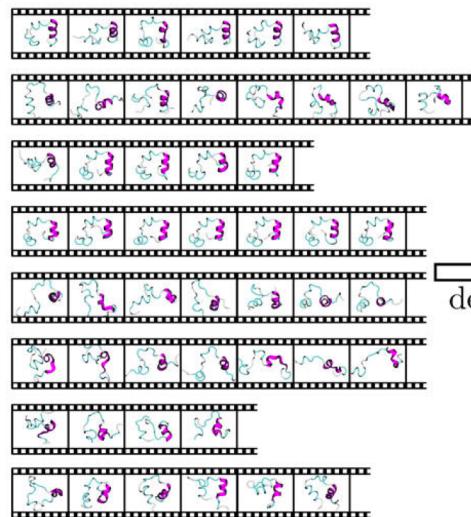
Standard MSM methods under-estimate timescales and free energies



Are there better estimators we can use?

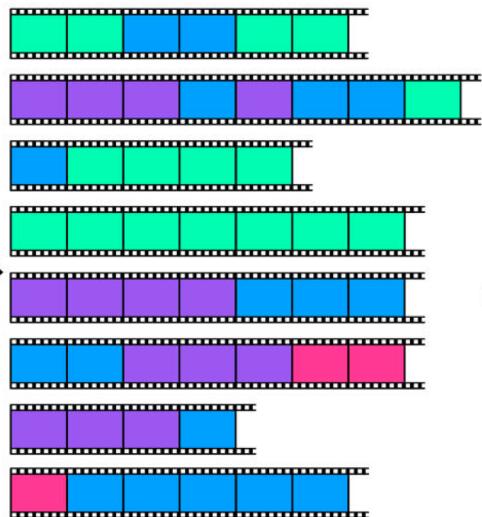
MSMs are built by estimating transition probabilities from observed counts

(a) simulation data



state
decomposition

(b) state assignments



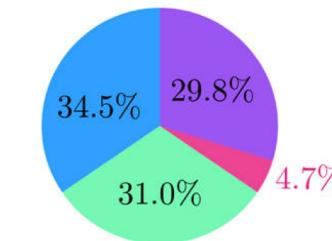
(c) observed transition counts matrix

11	1	0	0
3	9	2	0
0	4	9	1
0	1	0	1

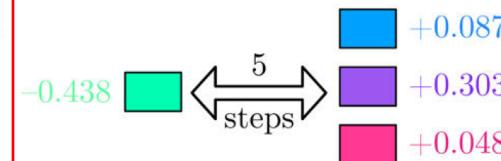
(d) reversible transition probability matrix

0.85	0.15	0	0
0.14	0.62	0.21	0.03
0	0.24	0.72	0.04
0	0.25	0.25	0.5

(e) adjusted populations



(f) slowest process eigenfluxes



Maximum-likelihood estimator (MLE)

- c_{ij} is the number of counts from state i to j (in lag time τ)
 - p_{ij} are transition probabilities from state i to j
 - π_i are the equilibrium populations of state i

$$p_{ij} = \frac{(c_{ij} + c_{ji})\pi_j}{N_j\pi_i + N_i\pi_j}$$

Maximum-likelihood estimator (MLE) if populations π_i are known

- c_{ij} is the number of counts from state i to j (in lag time τ)
 - p_{ij} are transition probabilities from state i to j
 - π_i are the equilibrium populations of state i

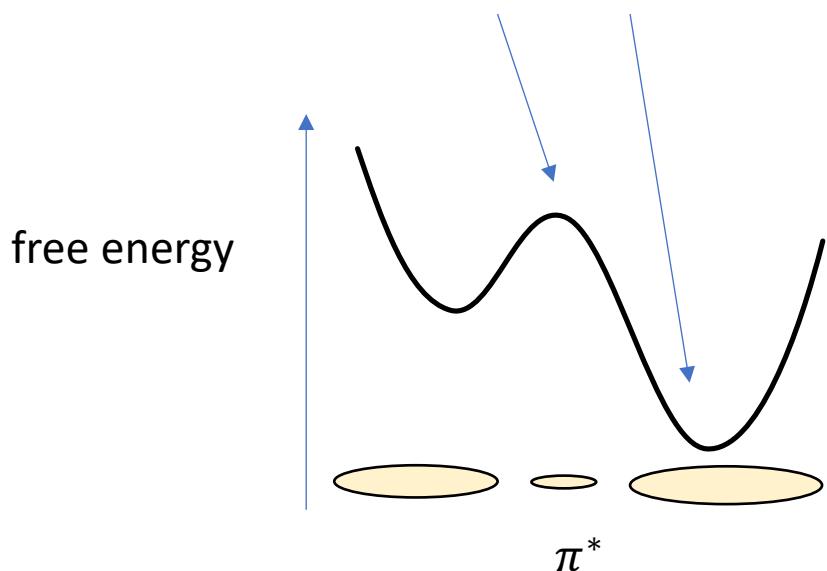
Self-consistent iteration to solve for
the Lagrange multipliers λ_i

Then the transitions rates are:

$$p_{ij} = \frac{(c_{ij} + c_{ji})\pi_j}{\lambda_j\pi_i + \lambda_i\pi_j}$$

Maximum-likelihood estimator (MLE) with π_i^* -weighted trajectory counts

Suppose 100 trajectories
are seeded at each state



some a priori estimate of
the state populations at
equilibrium



IDEA: weight the counts
from each trajectory to
better resemble counts
collected at equilibrium

$$c'_{ij} = \sum_k \pi^{*(k)} c_{ji}^{(k)}$$

Then use the MLE estimator
with detailed balance
constraint from before

$$L(p) = \prod_{i,j}^k p_{ij}^{c'_{ij}}$$

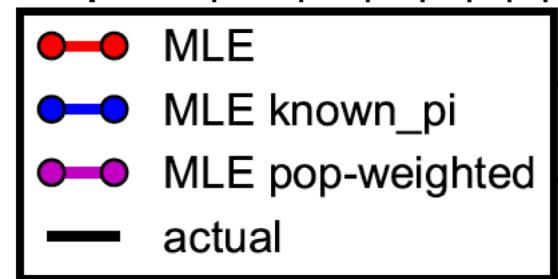
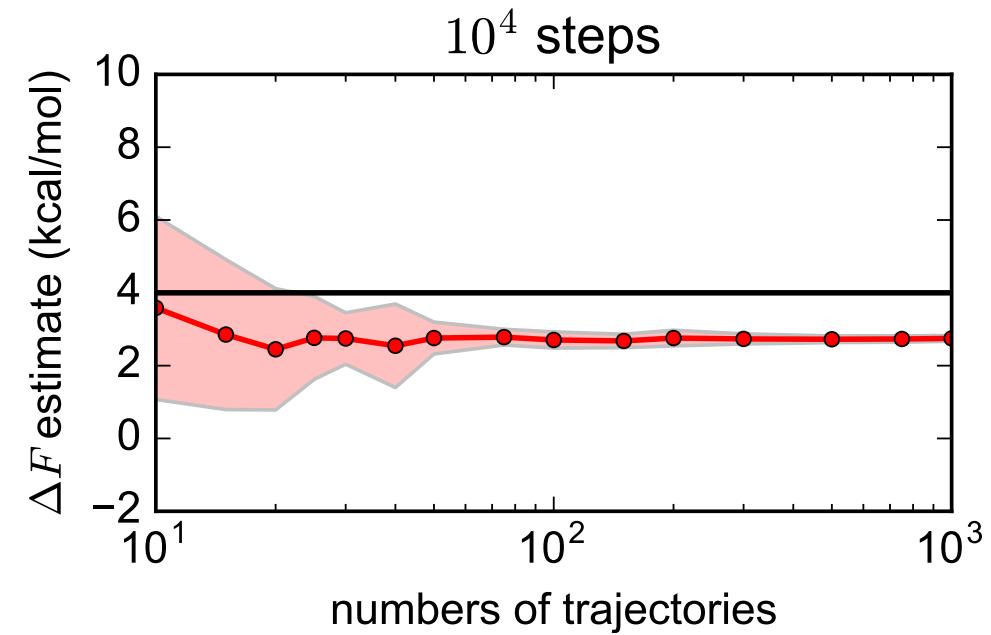
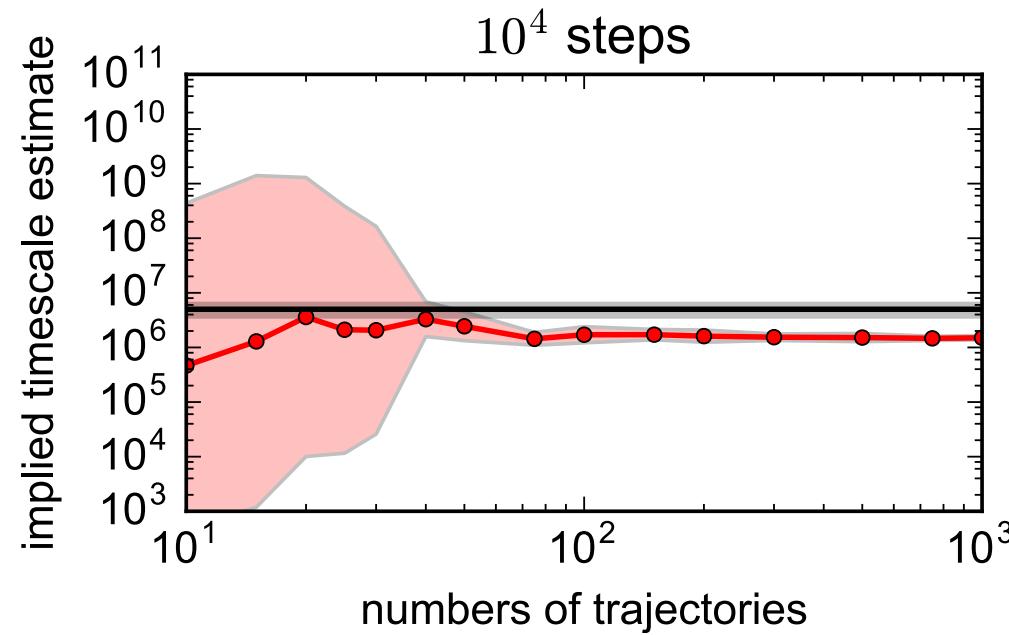
- $c_{ji}^{(k)}$ are counts from trajectory k .
- $\pi^{*(k)}$ is the weight of trajectory k

Row-normalized counts estimator

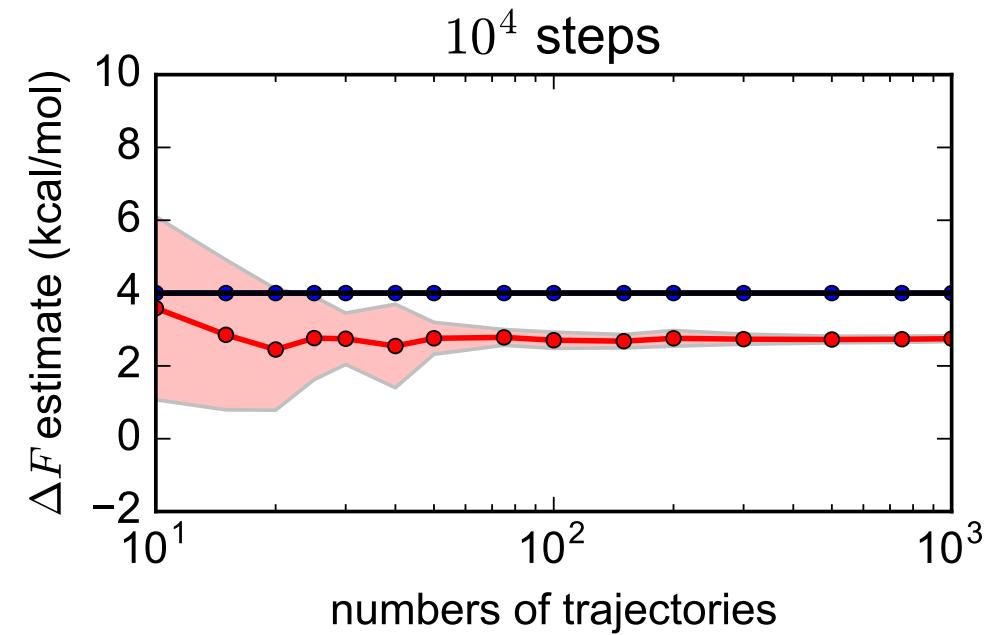
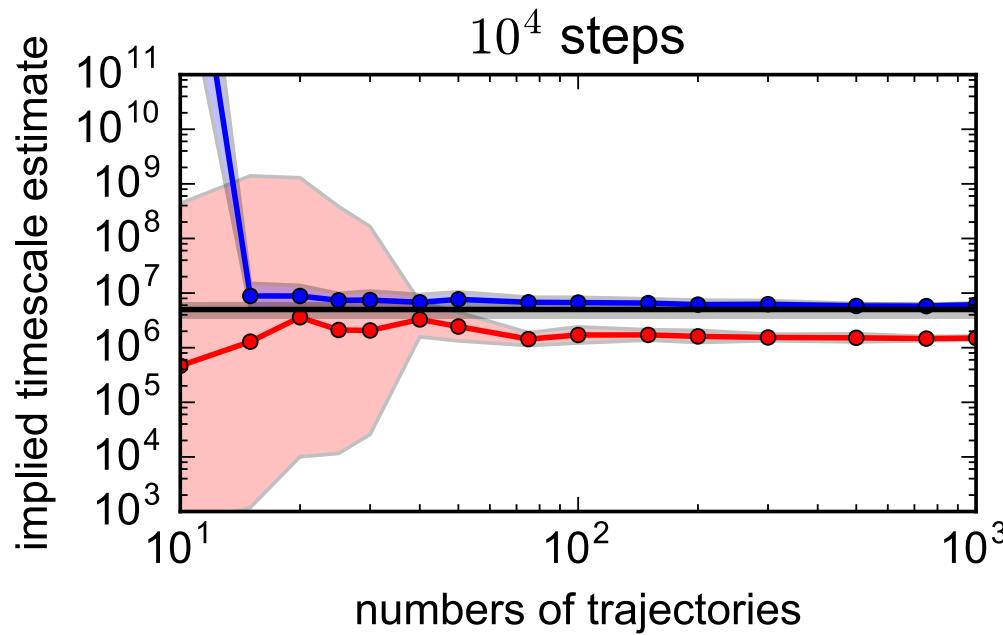
$$p_{ij} = \frac{c_{ij}}{\sum_j c_{ij}}$$

- Ignores detailed balance (!)
- Works only if $\sum_j c_{ij} \neq 0$.
- Hard for short trajectories with a larger lag time.

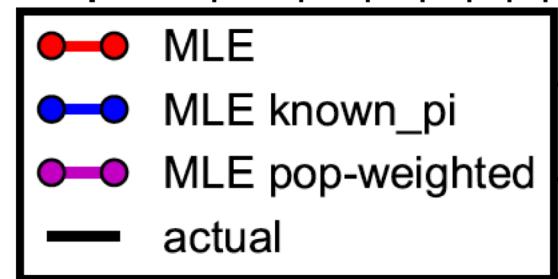
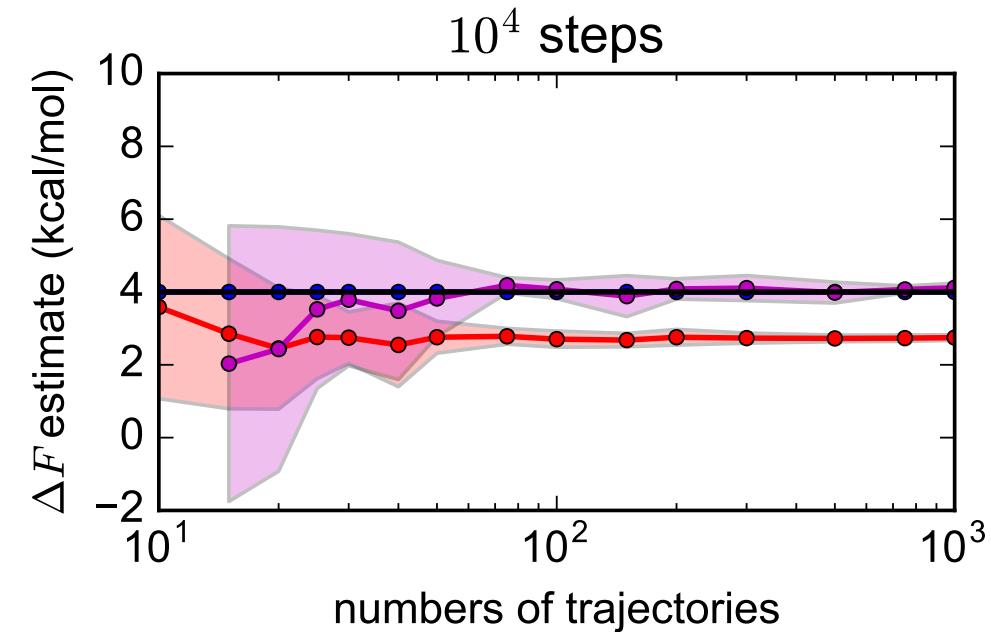
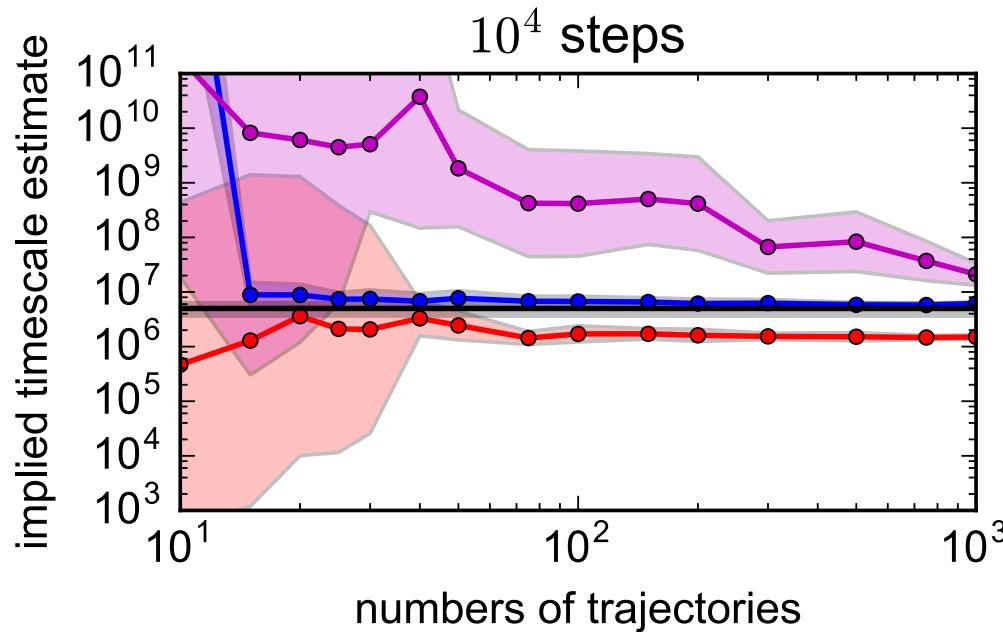
MLE(red) under-estimates timescales and free energies



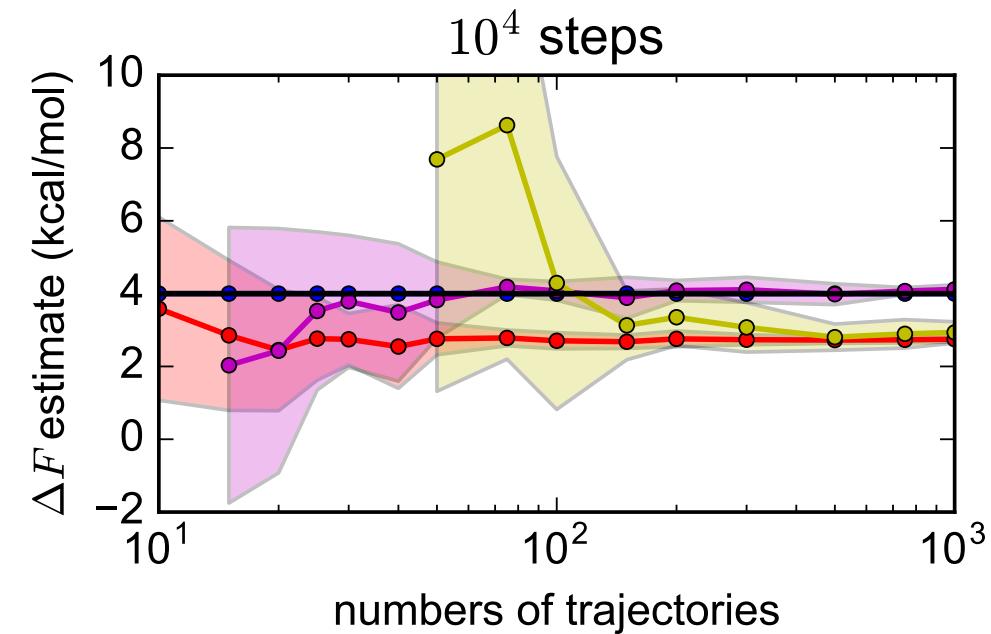
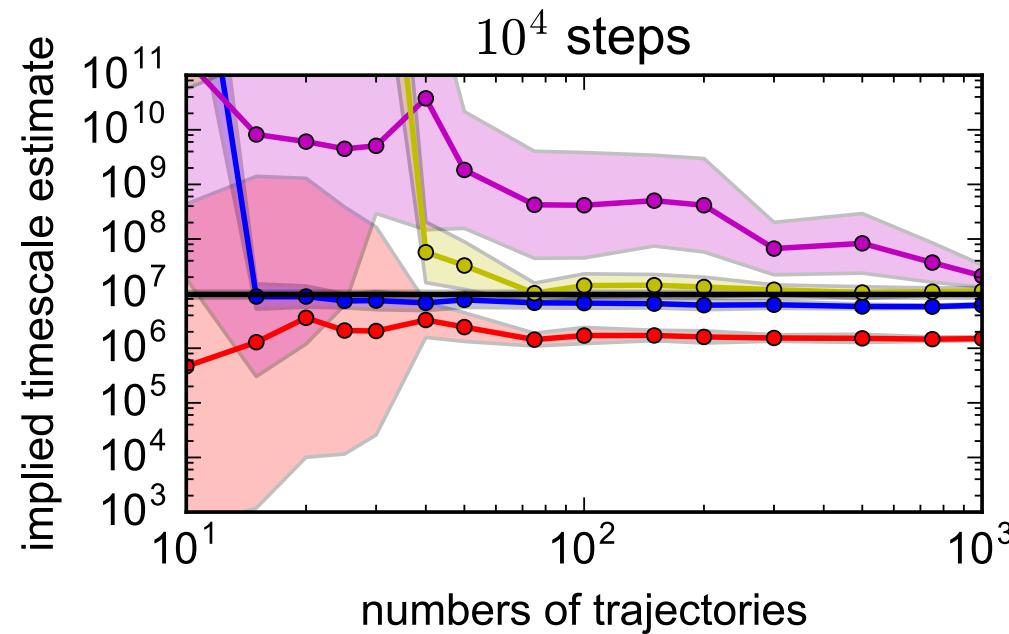
MLE-known_population (blue) converged quickly
on both rate and free-energy estimations.



MLE-pop-reweight(magenta) converges as more sampling is collected

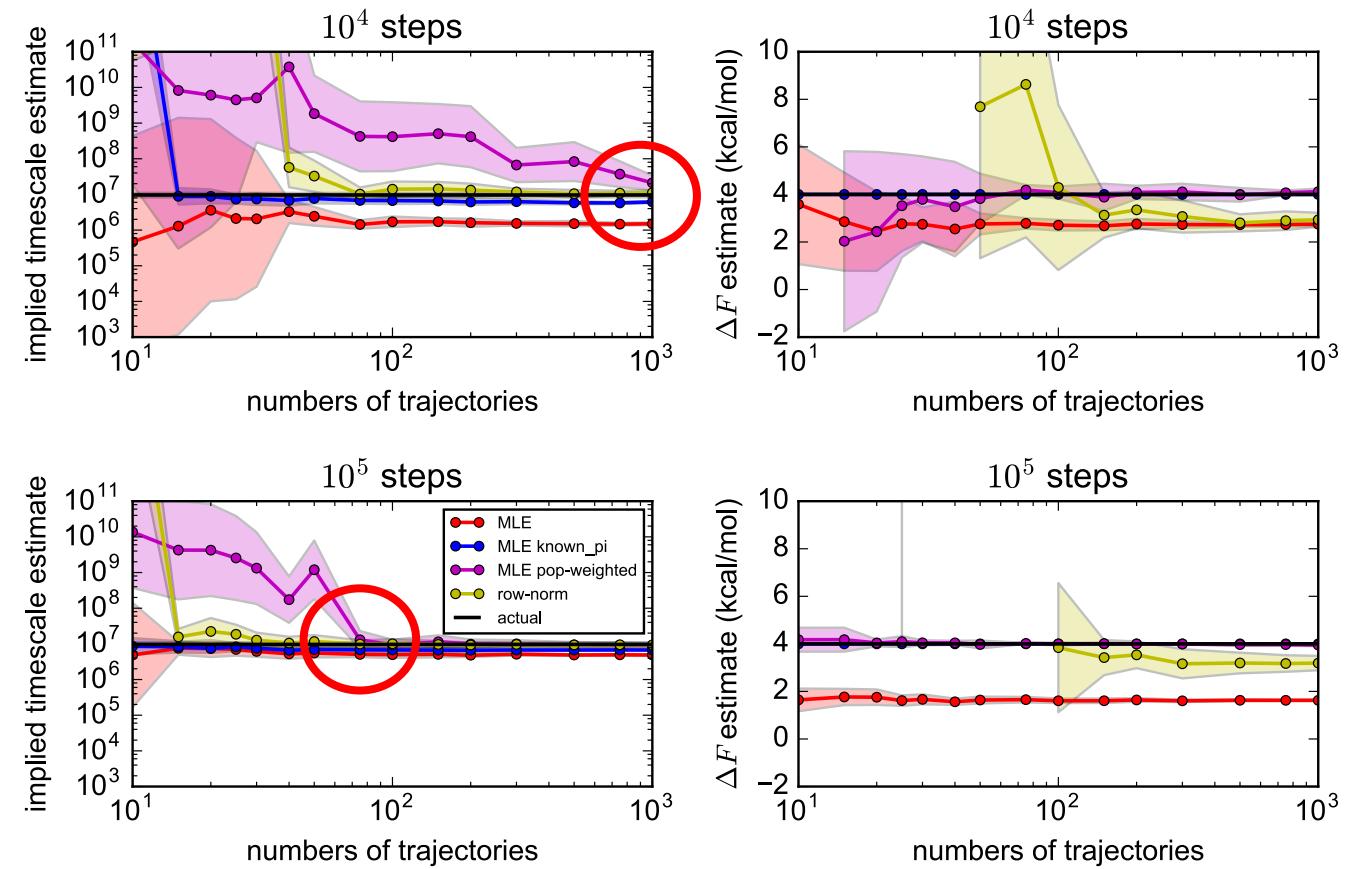


Row-norm needs a lot of trajectories to work well.

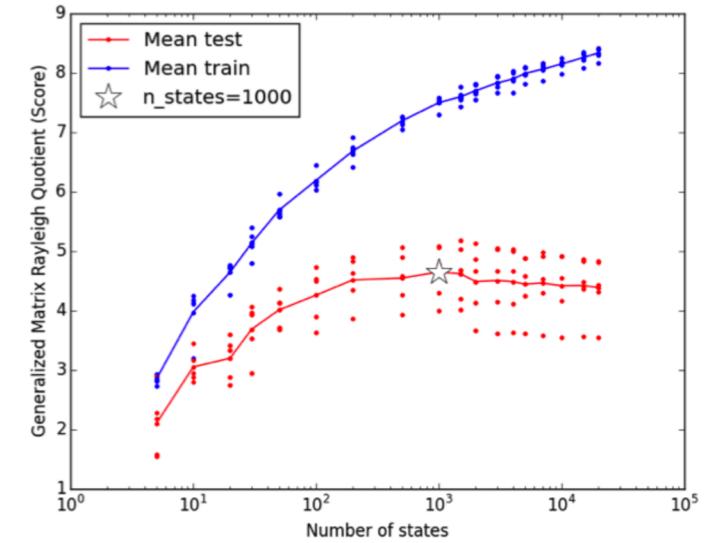
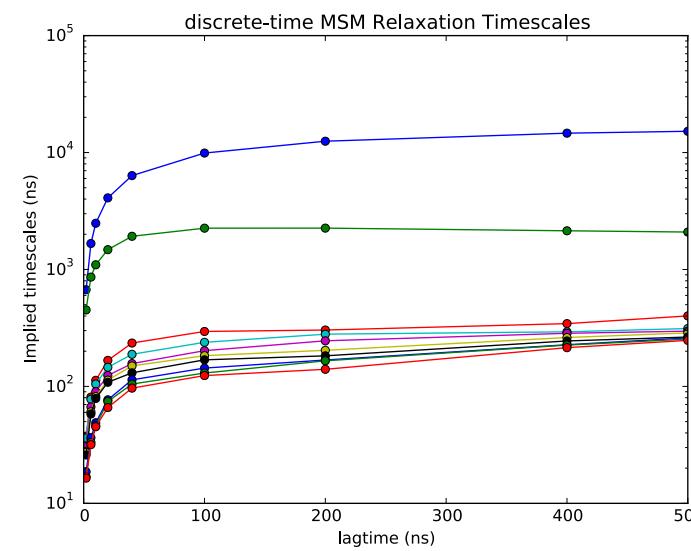
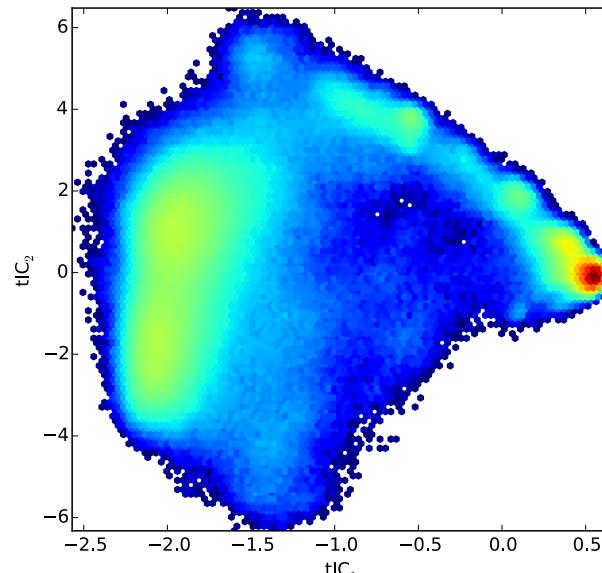


MLE-pop-reweight(magenta) converges if enough sampling is collected

- $N_{\text{bins}} * 10^7$ (slowest motion) steps(sampling) are needed for MLE-pop-reweights to get converged.

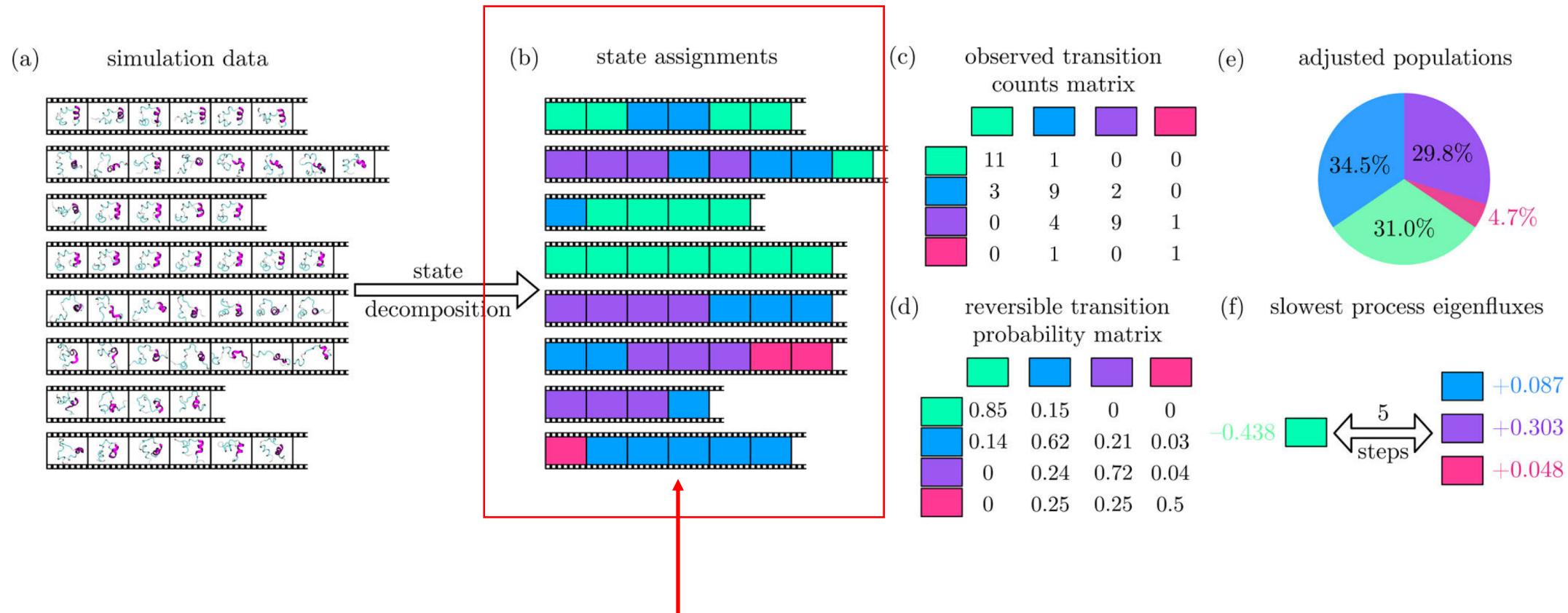


GTT at 360K with charmm22-star forcefield



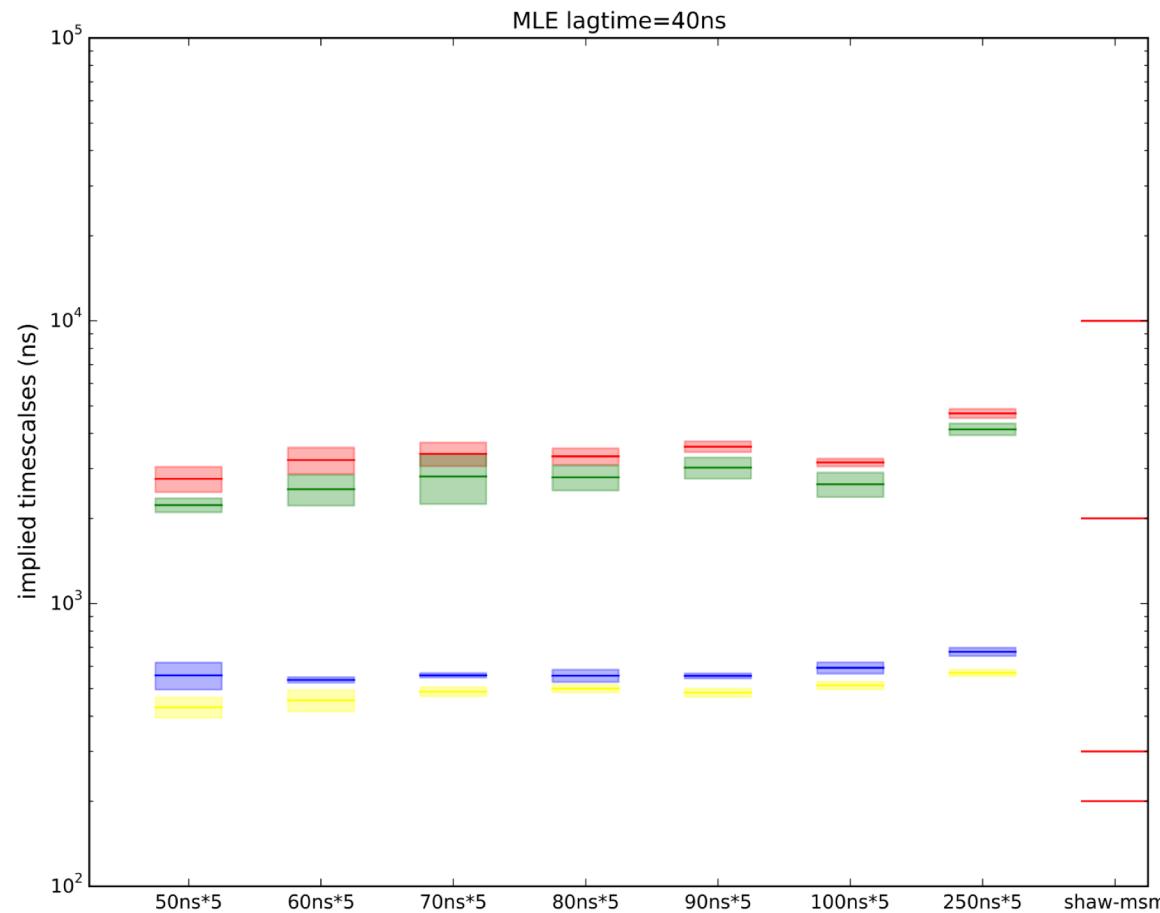
- Two independent simulation trajs with lengths of **651us** and **486us**.
- tICA with 10ns lagtime + k-center clustering algorithm.
- $C_\alpha + C_\beta$ pairwise distances was used as the input coordinates for tICA.
- GMRQ was used to determine that 1000 micro-states and 8 tICAs were optimal to capture the dynamics.
- MSM was built with a lag-time of 100ns.

Making-up short trajectories from an existing dataset

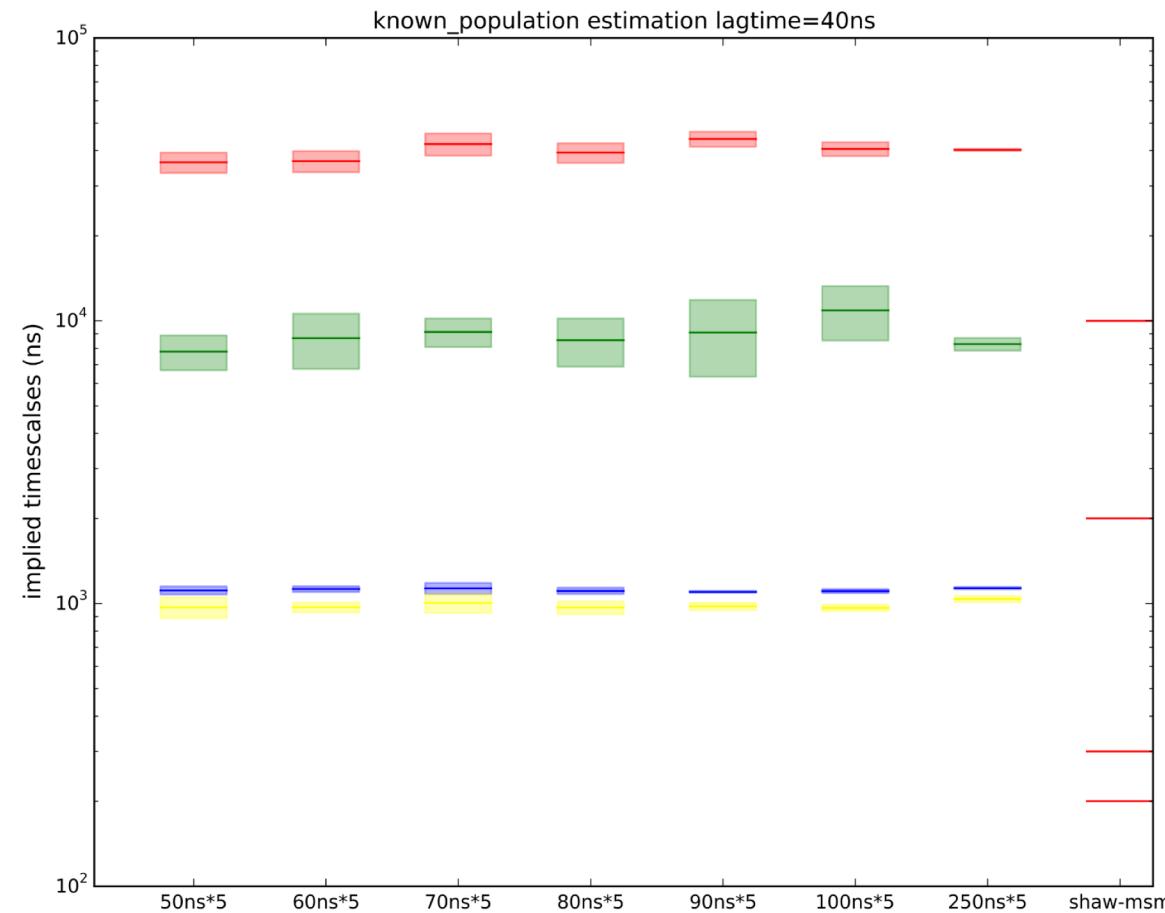


Expectation: it should re-produce both kinetic and thermodynamic properties

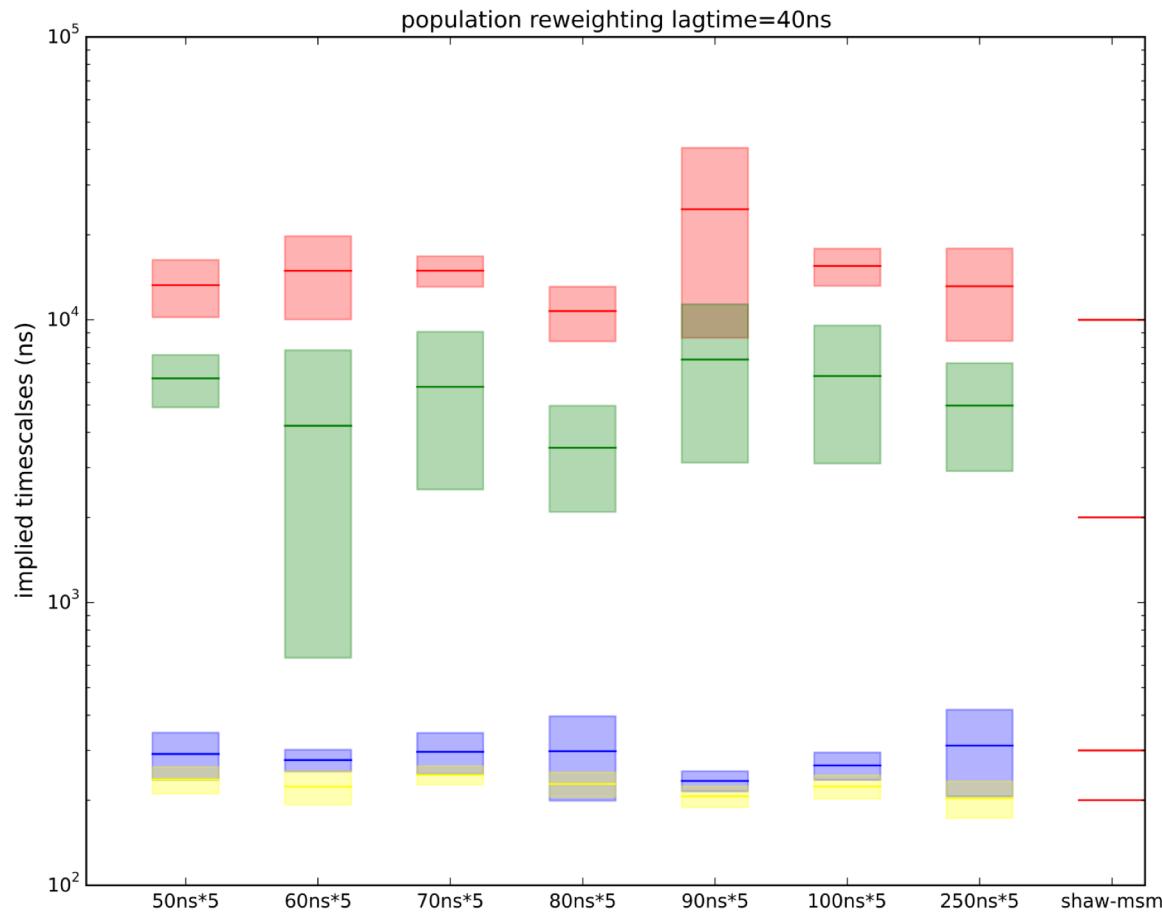
MLE underestimates folding rate



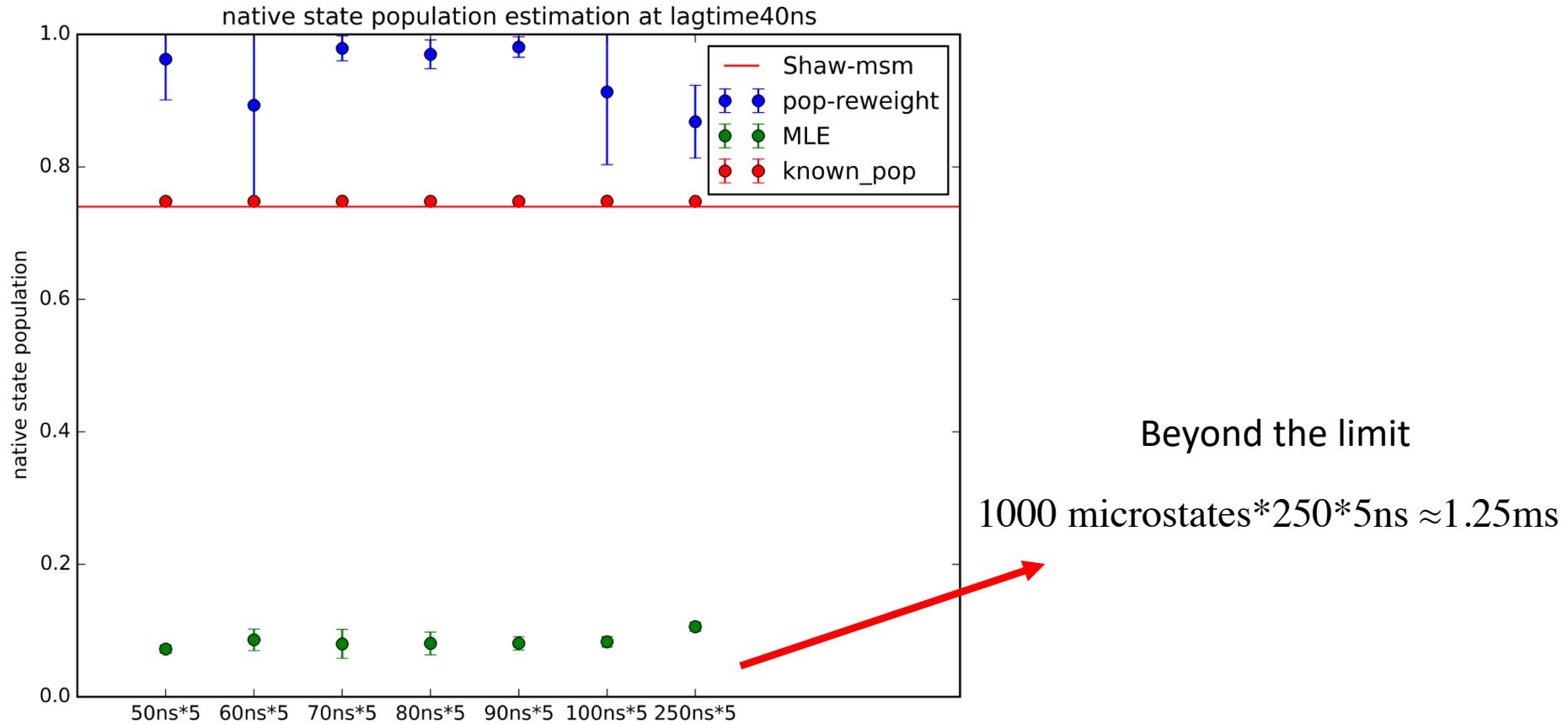
MLE-known_population overestimates folding time



MLE-pop-reweight converges in this large dataset

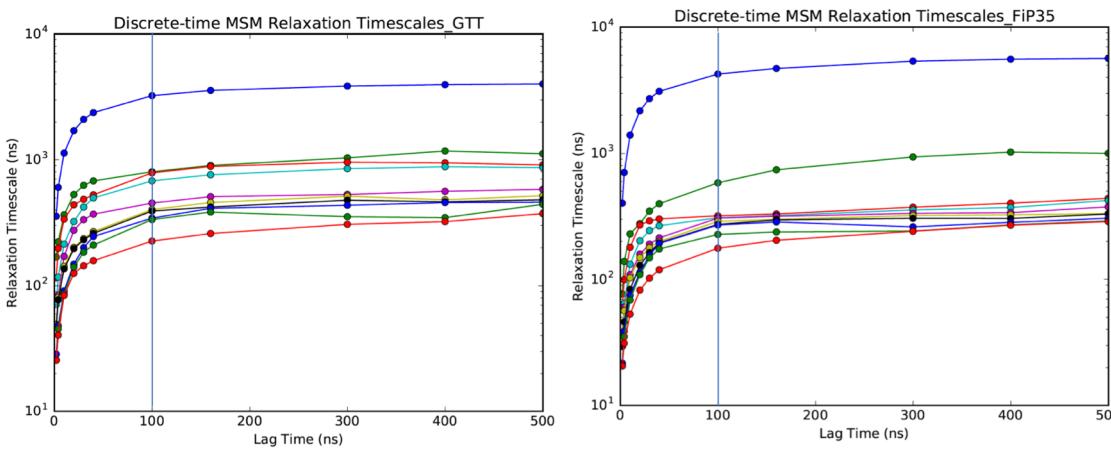
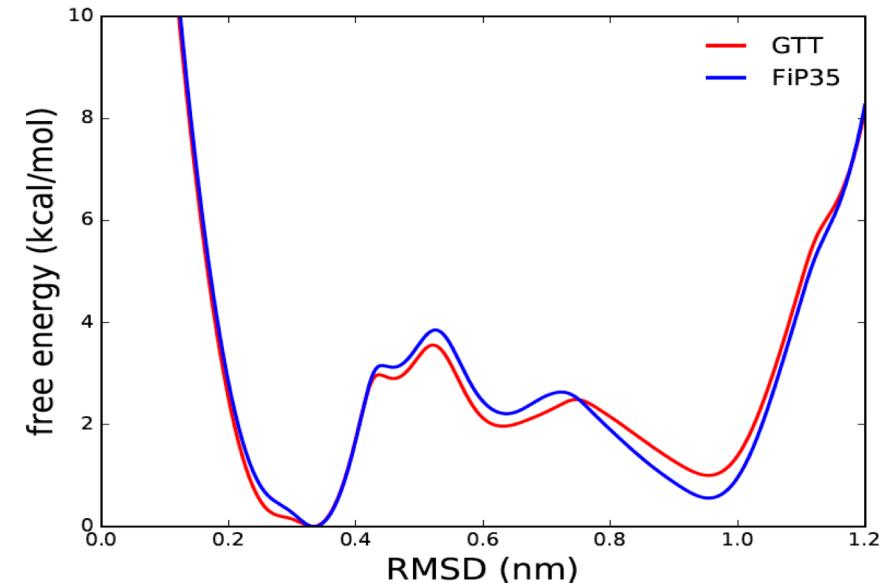
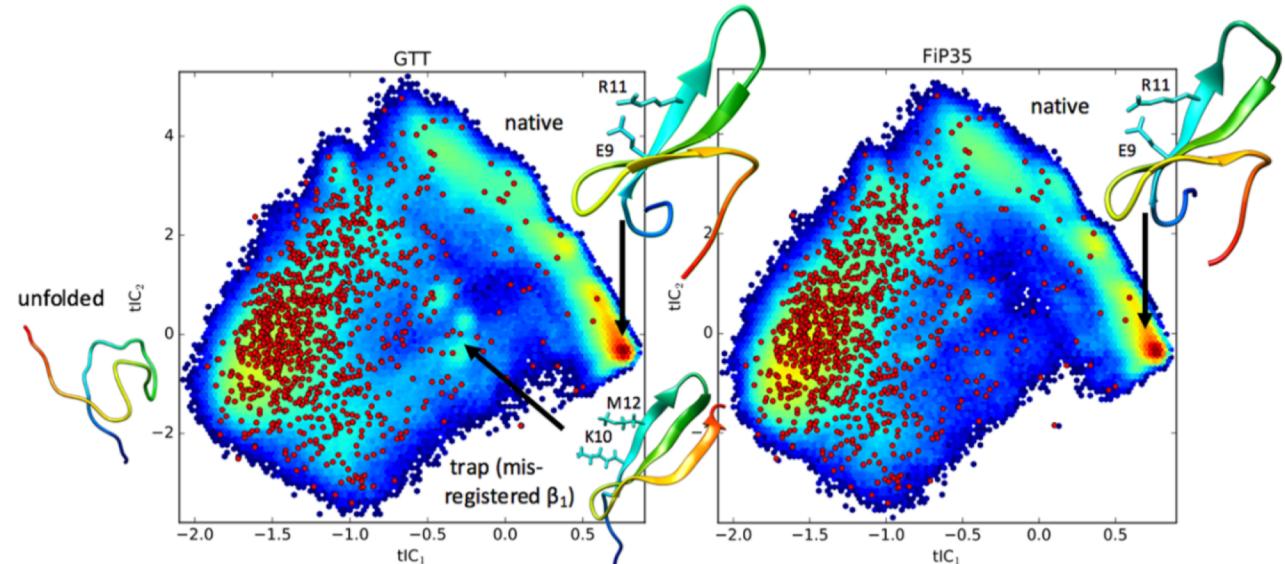


The native population estimation from MLE-pop-reweight(blue) gets better as more sampling collected



- $N_{\text{bins}} * \text{slowest motion} \approx 1000 \text{ microstates} * 10 \mu\text{s} = 10 \text{ ms}$

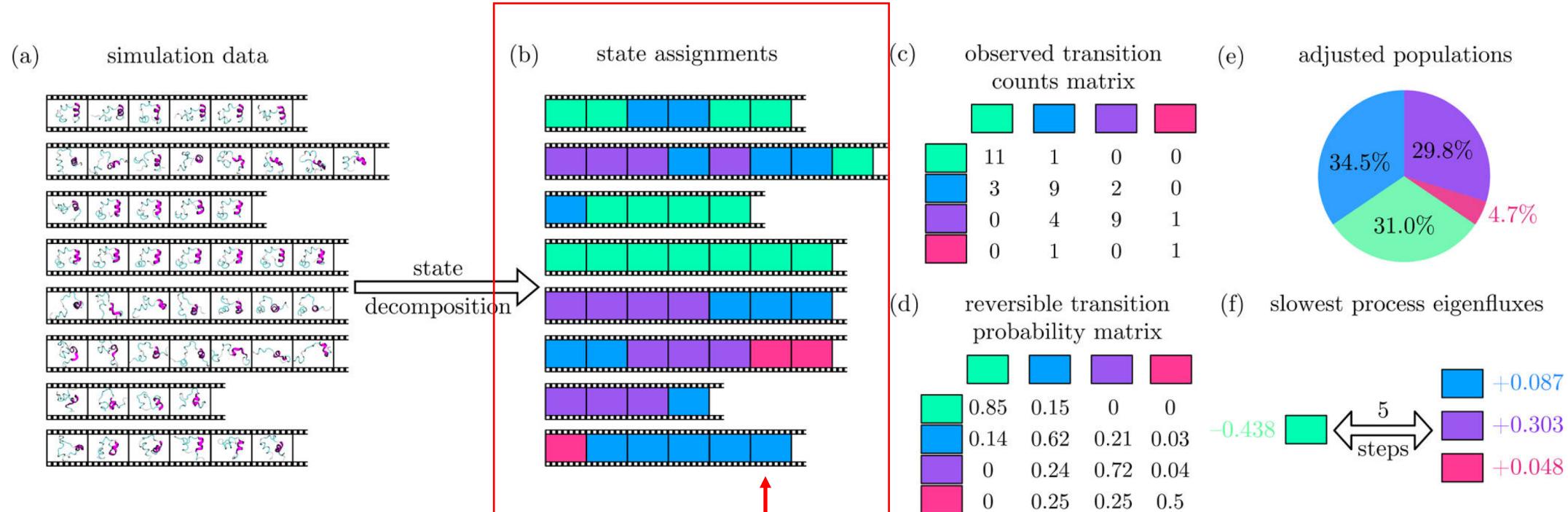
WW domain FiP35 and GTT



- GTT-mutation folds faster than FiP35 and stabilizes the native state.
- These two simulations are at same temperature and used same forcefield.

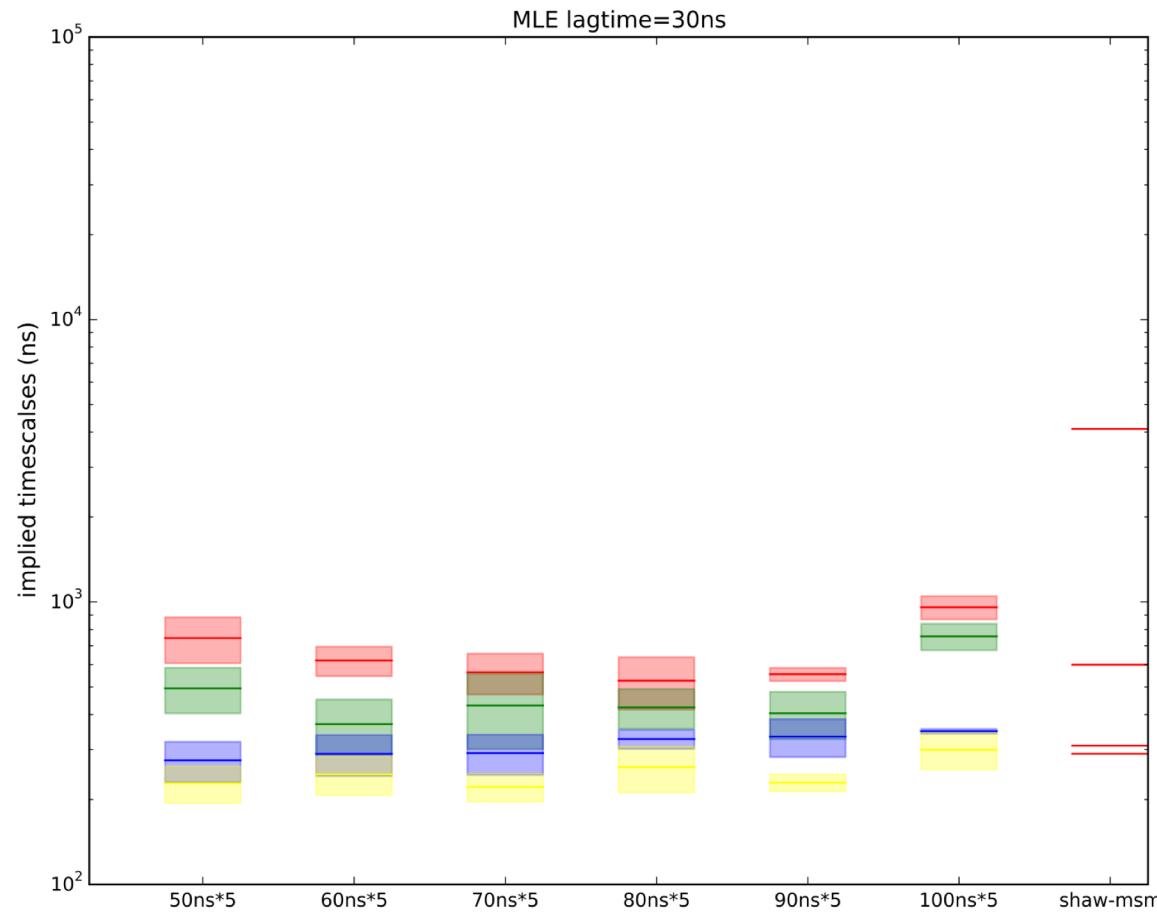
Q: to see if we can accurately predict FiP35 folding behaviors from seeding simulations of GTT(ref)?

Making-up short trajectories from an existing dataset

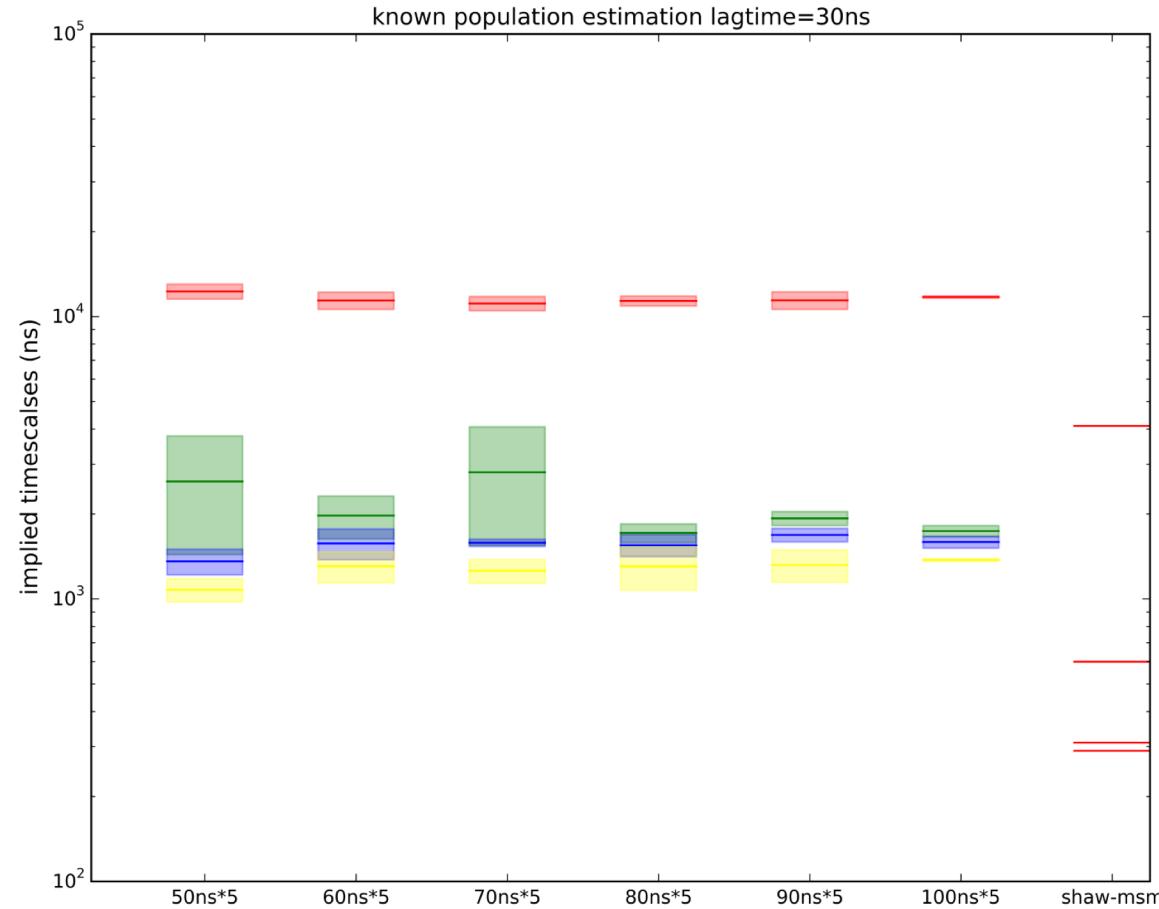


Expectation: seeding from GTT(mutation_ref) can predict both thermodynamic and kinetic properties of FiP35(WT)

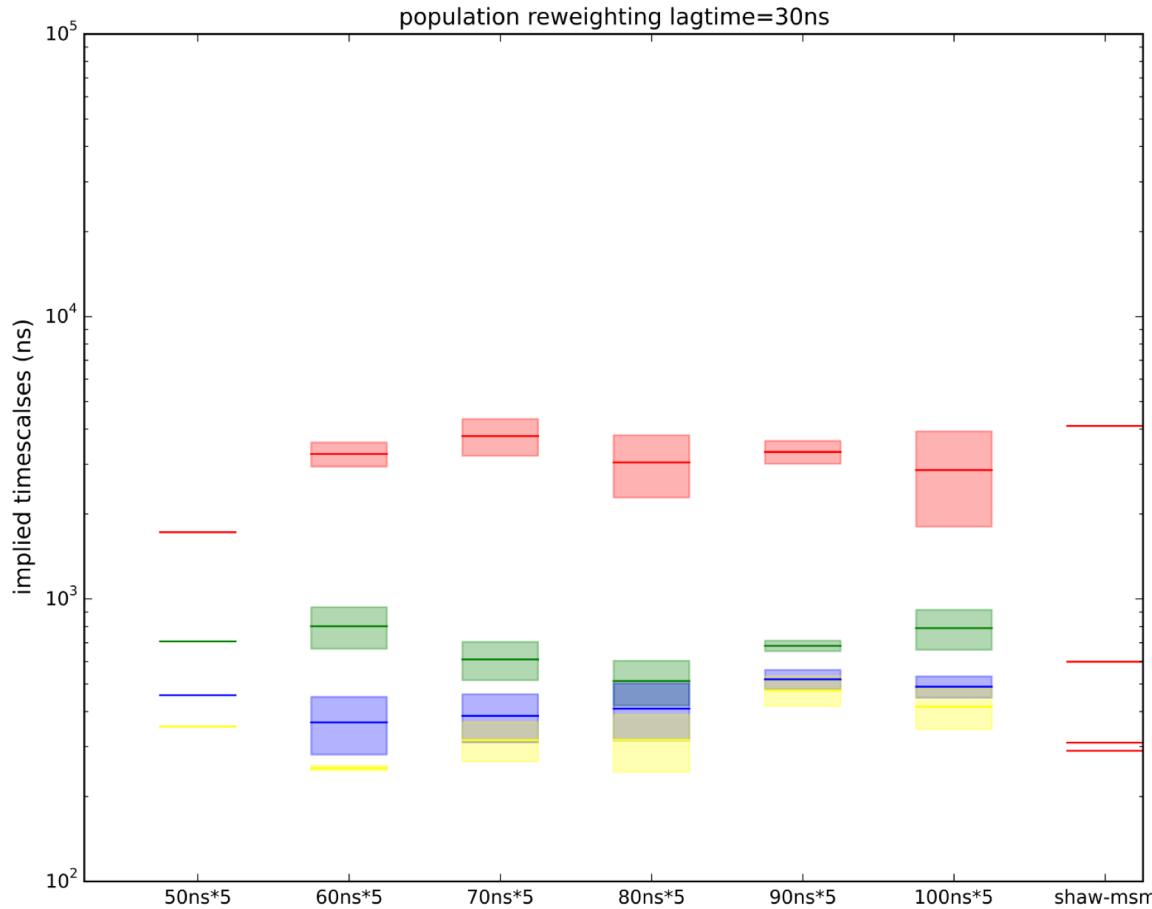
MLE underestimates folding time



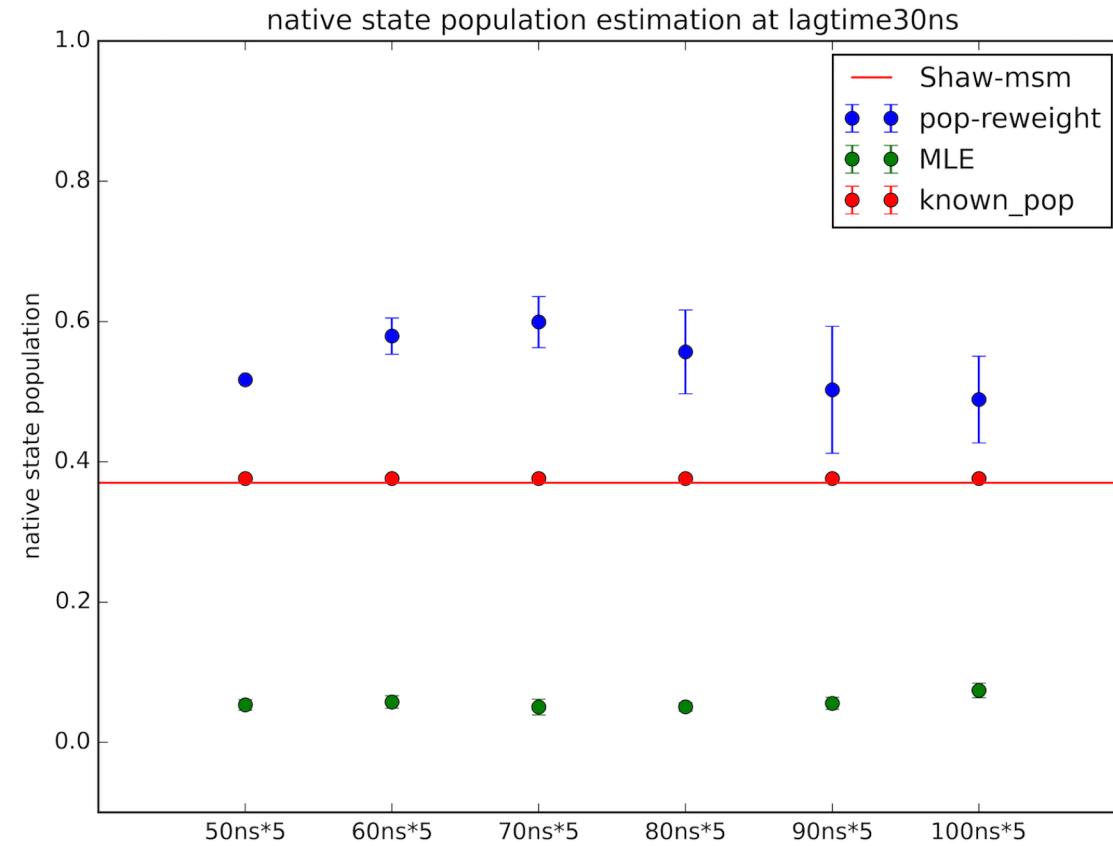
MLE-known-population overestimates folding time



MLE-pop-reweight converges as more sampling is collected



The native population estimation from MLE-pop-reweight(blue) gets better as more sampling is collected



Conclusions:

- Seeding from a reference MSM
 - MLE estimates a faster folding time and a lower native-state population.
 - MLE_known-pop can quickly get the right free energy profile but over-estimates the folding time.
 - MLE_population-reweight can properly estimate the folding time, and the free energy profile converges with more sampling.

Thanks Voelz lab

