

Markov-state modeling of biomolecular systems



Toni Giorgino

National Research Council of Italy

toni.giorgino@cnr.it



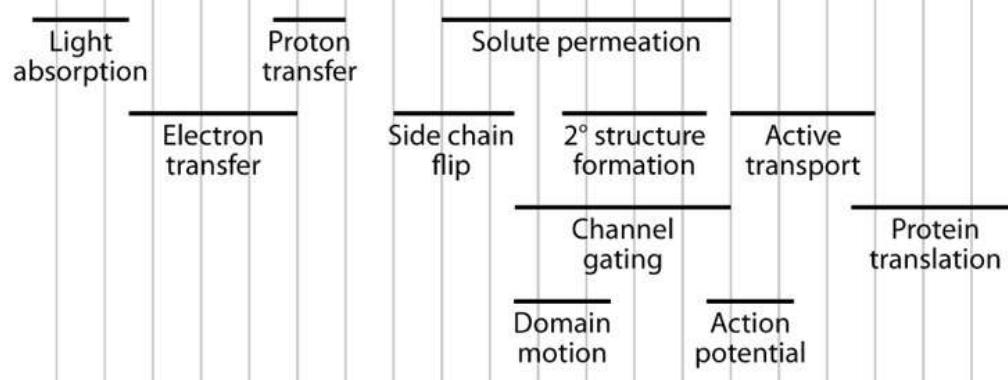
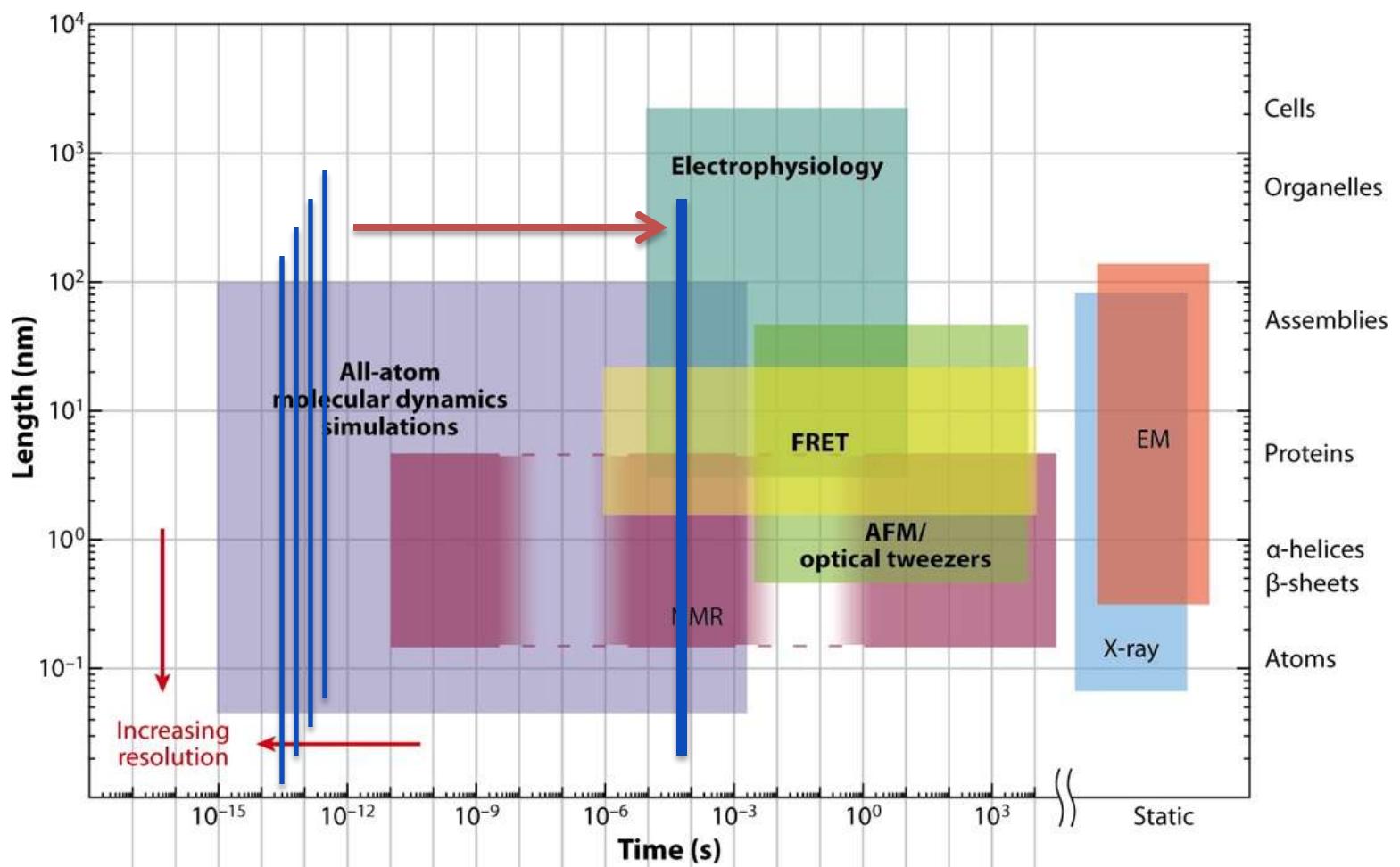
Introduction

- The aim of this class is to provide a practical overview of Markov state models in computational structural biology
- MSM emerging because:
 - reconstruct *kinetic information**^{*}, including state transition networks, from simulated trajectories
 - start from **unbiased simulations** (no *a priori* reaction coordinate hypothesis necessary)
 - microsecond-scale (high-throughput) trajectory data are becoming accessible (e.g., with GPUs)
- Success cases: *ab initio* folding, drug binding, peptide binding, ...

* as well as the corresponding structures

Motivation

- Can we use an *ensemble* of simulations to estimate dynamic behaviours which happen on timescales longer than each of the observed trajectories?
- In other words, can we leverage several “short-sighted” (non-equilibrium) observations to *extrapolate* long-time behaviour?



Dror RO, et al. 2012.

Annu. Rev. Biophys. 41:429–52

Warning

- Still very active field
- For serious work, many more details are in...
 - the theory of Markov state models
 - the discretization, projection, and estimation of models from trajectories
- Suggested further steps: worked out real-world examples distributed with software packages*

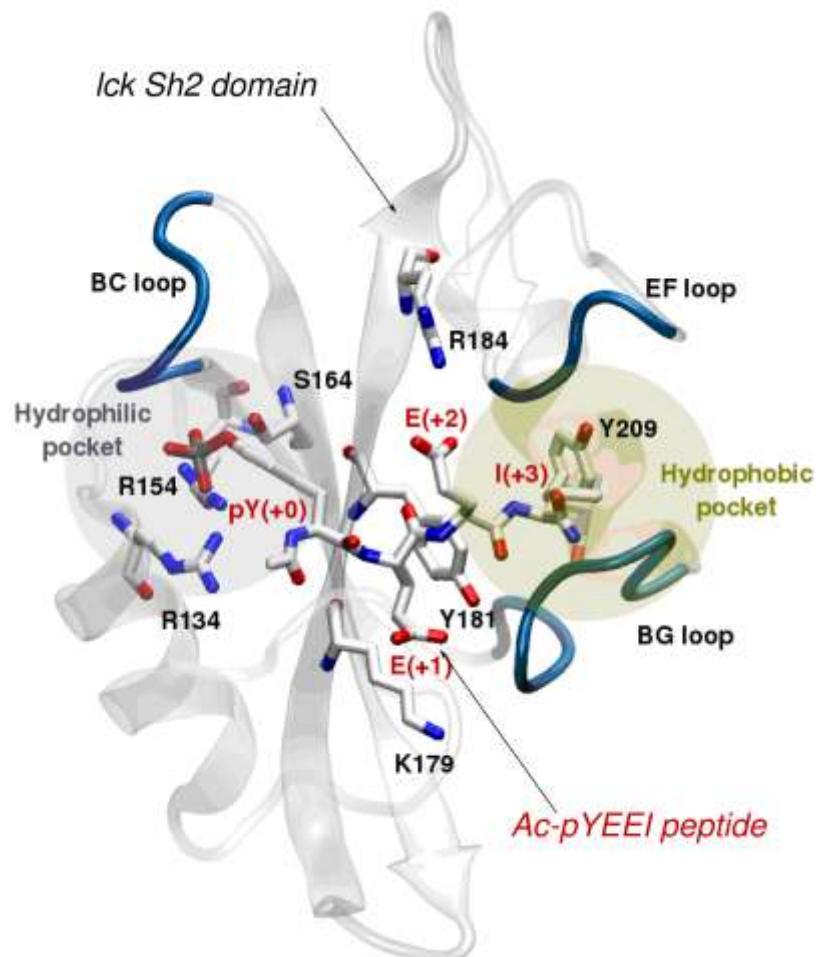
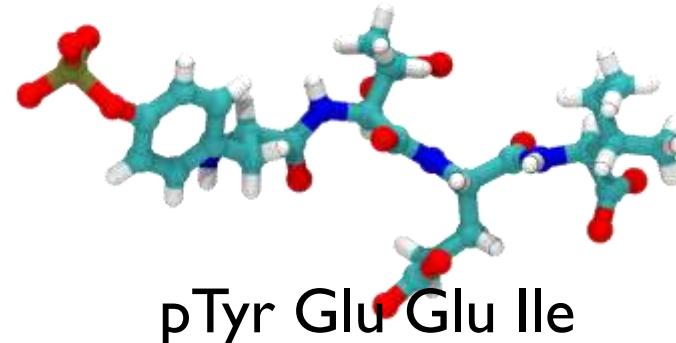
* see the last slides

Binding as an irreversible two-state model

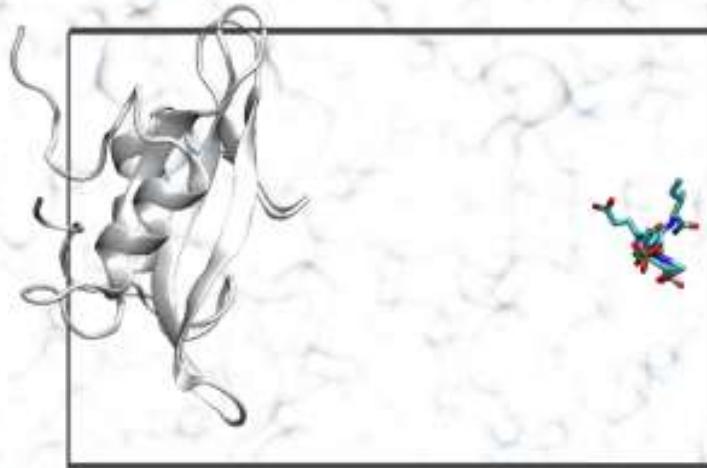
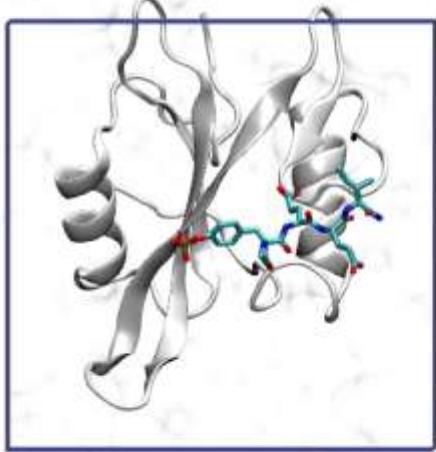
**Unbiased simulations:
SH2 – pYEEI binding example**

SH2 Domains

- Small domain (100 AA), well-conserved
- Recognize short pY* sequences with high specificity
- Prominent role in signaling
- Found in ~110 human proteins
- “Socket with two holes”



In silico set-up

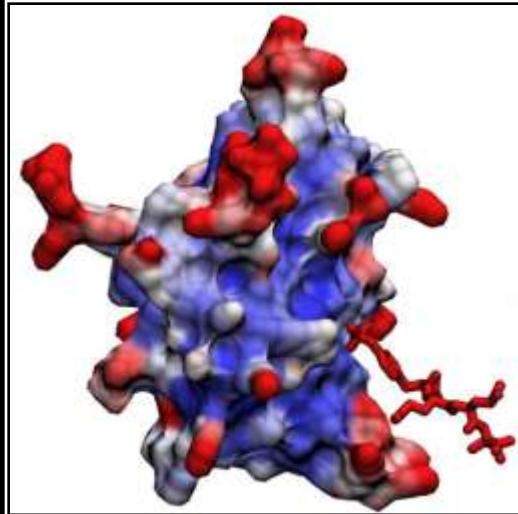


10 Å

7 nm

10 nm

Teaser (1:36)



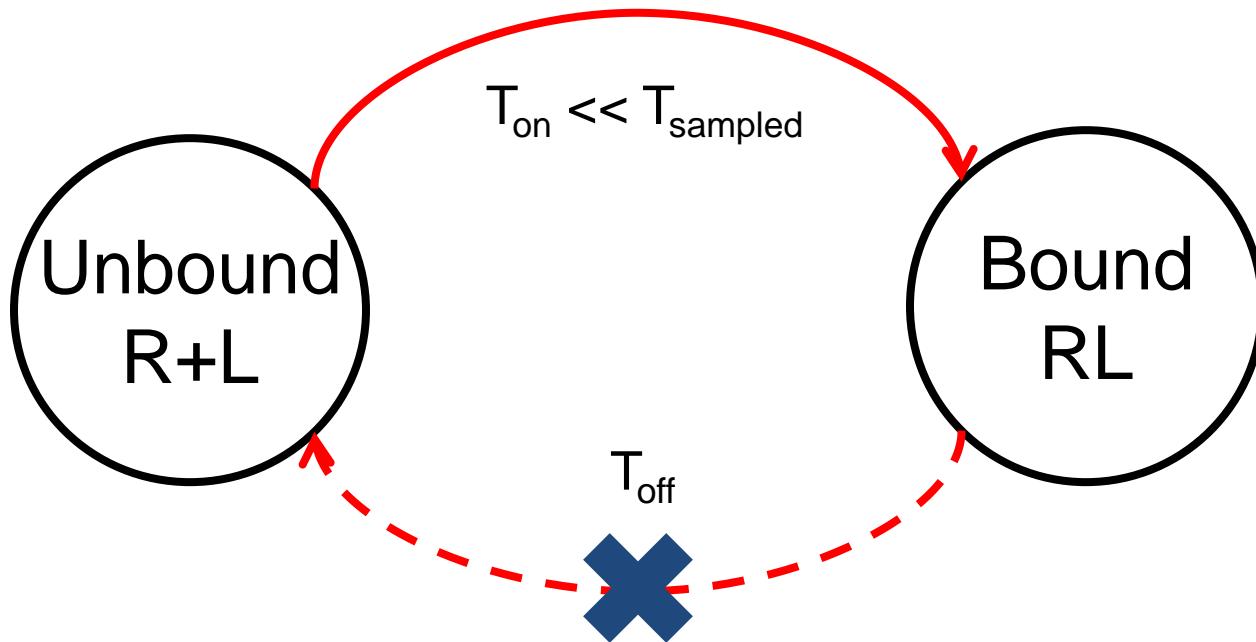
- Video: *unbiased binding simulation of pYEEI to Lck SH2 domain.*
- *Protein surface and ligand's atoms color-coded by local flexibility, from blue to red (lowest to highest RMSF).*
- *As the peptide binds into its native conformation ($\text{RMSD} < 2 \text{ \AA}$ from crystal structure), it loses flexibility; the same happens to the protein.*
- **By construction, we are not at equilibrium.**

<http://goo.gl/ZL142>

also https://youtu.be/Xhof-pf_Eo8

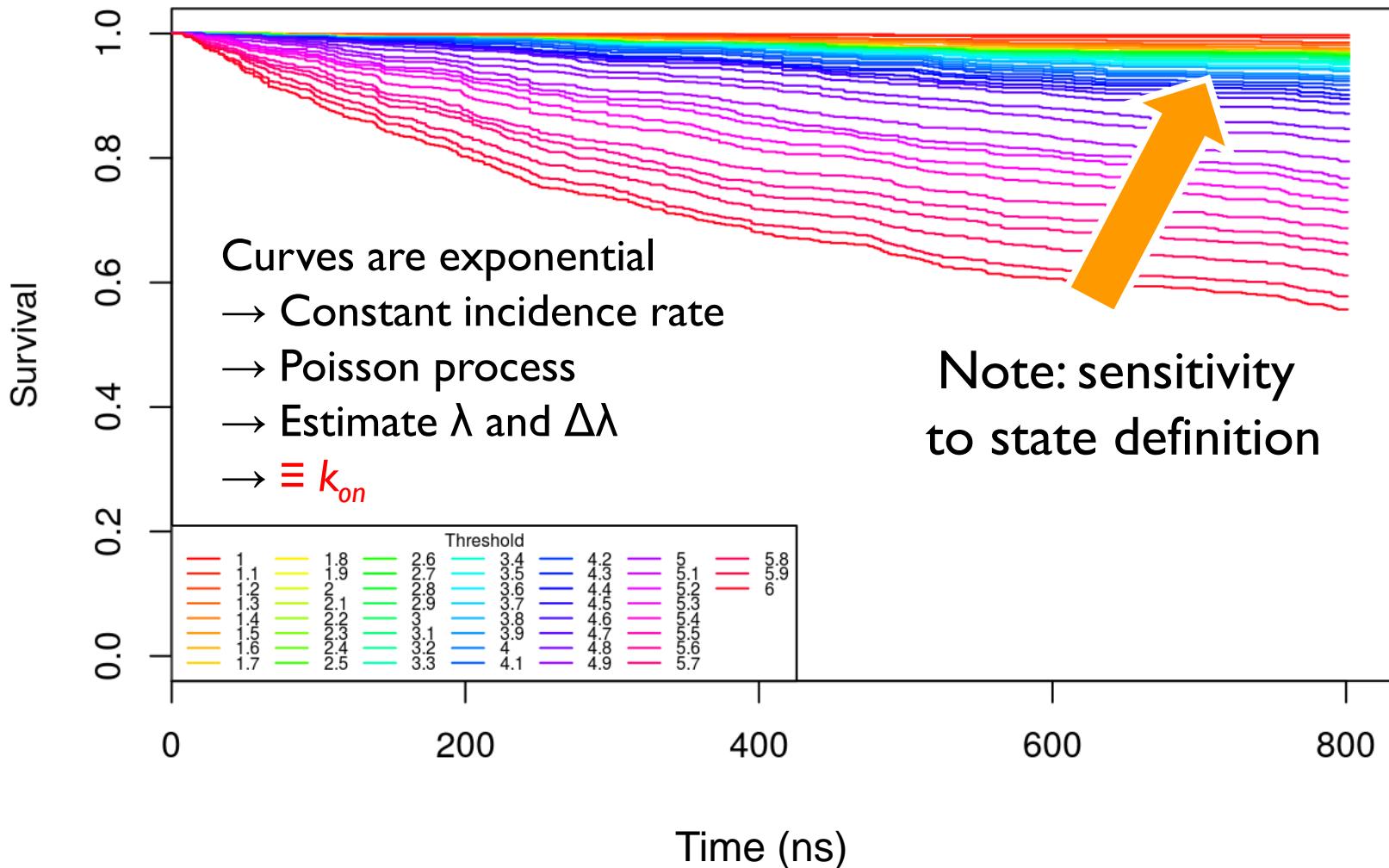
Binding as an *irreversible* two-state model

On the sampled timescales...



$T_{event} \ll T_{on} \ll T_{sampled} \ll T_{off}$

Non-bound-fraction curves



Motivation for a more general solution

This 2-state method is not sufficient:

1. One needs to define the *bound state*
 - which is not available if we lack a crystallographic or NMR structure.
2. We need to observe a sizeable number of *full binding events*.
 - Association processes are slow.
They may be (far) outside our sampling power.
3. We get no information about *off rates*.
4. We get no information about *metastable states*.

Going multi-state: Discrete-time Markov Chains

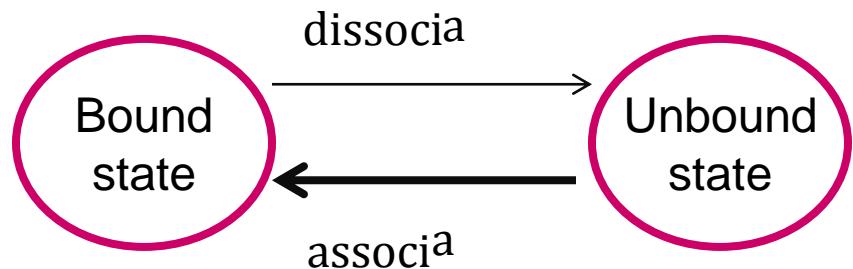
Why do we care?

Because equilibrium probabilities are the free energy of binding (ΔG).

ΔG is an important determinant of drug potency.

$$\begin{aligned} K_D &= p_U/p_B (\propto [P][L]/[PL]) \\ &= \exp(-\Delta G/RT) \end{aligned}$$

$$\Delta G = -RT \ln K_D$$



Kinetics

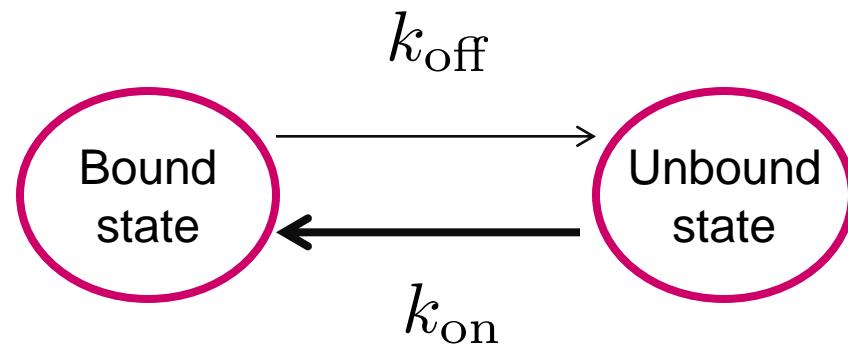
But the real breakthrough is that we can use the *dynamic* information, computing the *kinetics* as well. This is why we are doing *unbiased dynamics* in the first place!

Binding kinetics is important for rational drug design as well (e.g. related to the time spent in complex with receptors).

$$K_D = k_{\text{off}}/k_{\text{on}}$$

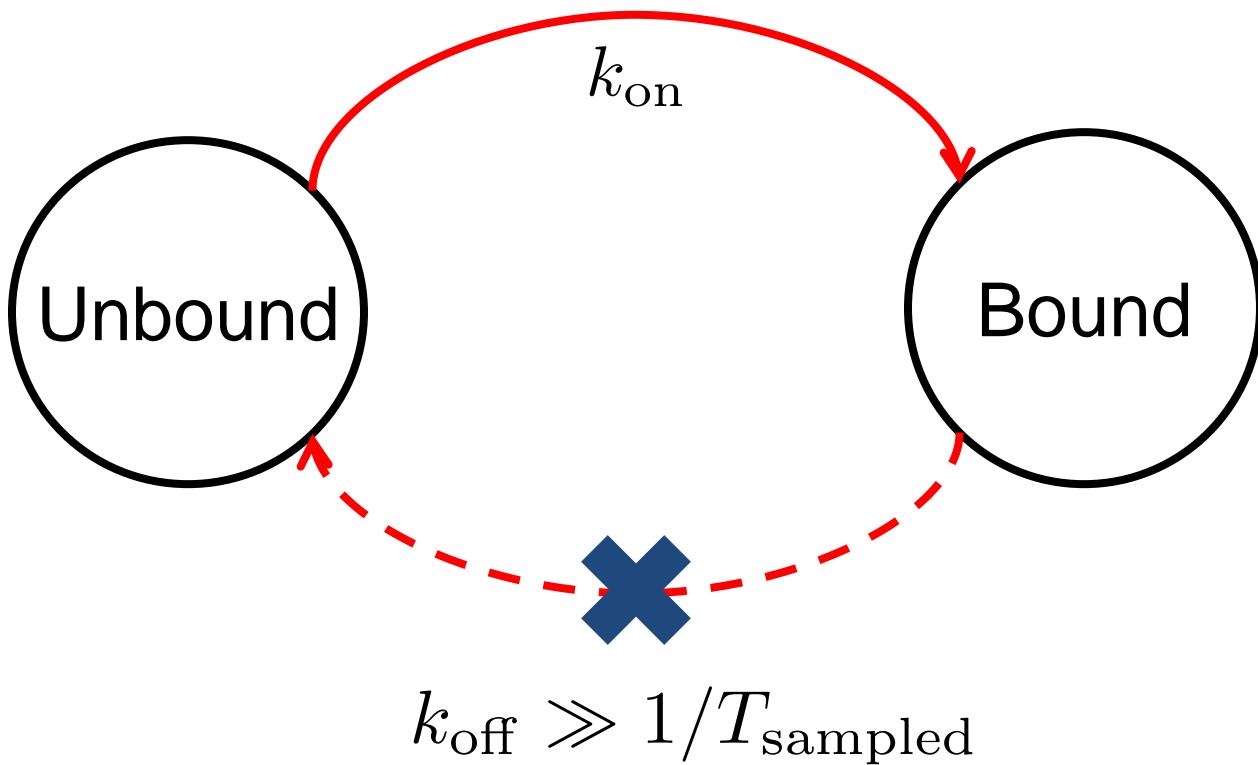
$$\propto \tau_{\text{on}}/\tau_{\text{off}}$$

$$\sim t_U^{(B)}/t_B^{(U)}$$



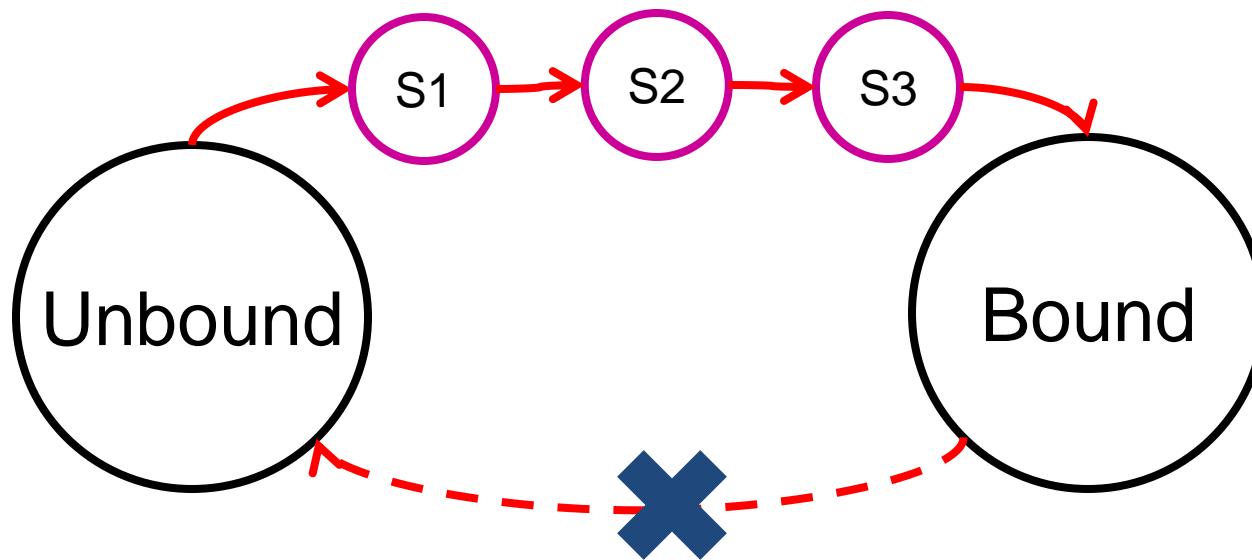
Let's move beyond this model...

On the sampled timescales...



Idea

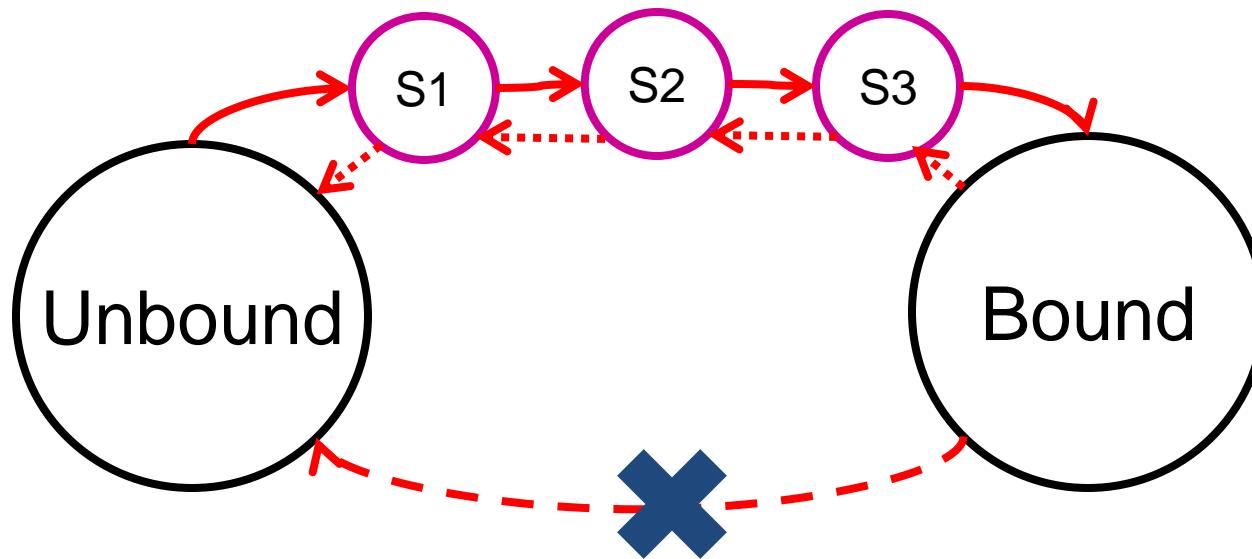
- We introduce several *intermediate steps*



k_{on} reconstructed combining rates of different intermediate steps

Idea

- We introduce several *intermediate steps*
- Compute MFPT



k_{on} reconstructed combining rates of different intermediate steps

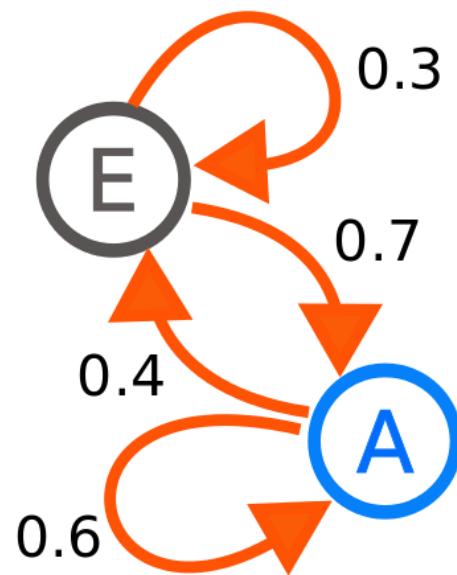
k_{off} may also be reconstructed, if states are fine-grained!

Going multi-state: Discrete-time Markov Chains



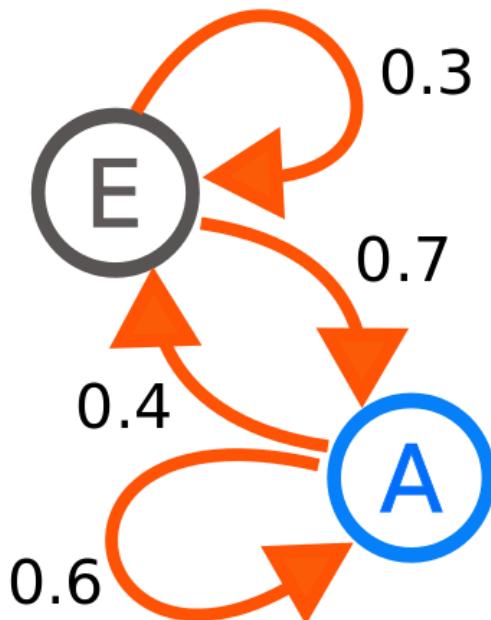
Discrete Time Markov Chains

Andrei Markov
1856-1922



- A **random** process.
- The system's state is a **discrete** variable.
- It undergoes transitions between states at uniformly-spaced (**discrete**) time points.
- Transition probabilities do not depend on the previous history of states (**memorylessness**).

Transition probability matrix

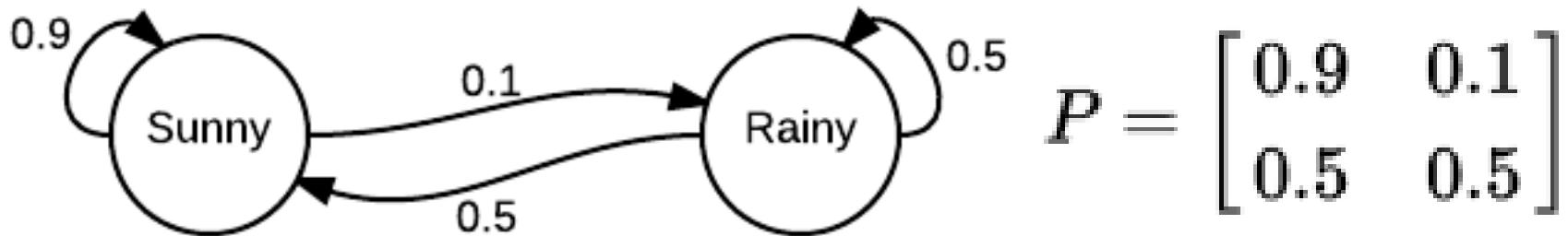


P_{ij} = Probability to change state
from i to j at each time point

$$P_{ij} = P(X_t = j | X_{t-1} = i)$$

		j	
		A	E
		0.6	0.4
i	A	0.6	0.4
	E	0.7	0.3

First example



Assume a deterministic initial condition:

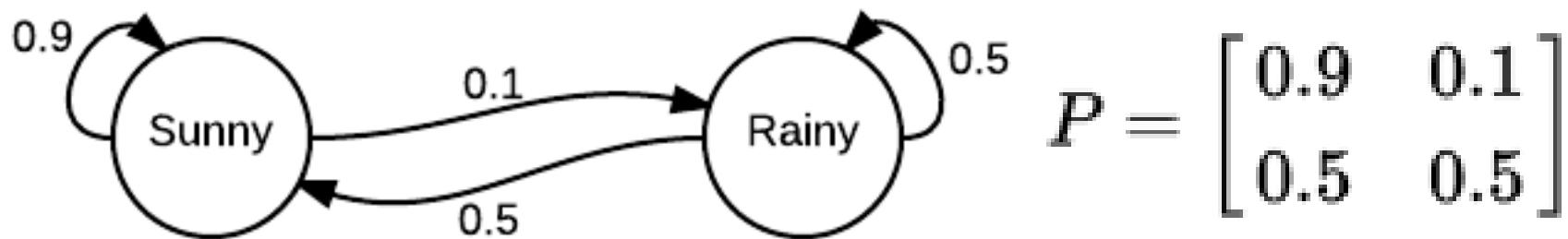
$X_0 = \text{Sunny}$ with certainty; i.e.,

$$\begin{aligned} p(\text{Sunny} | t=0) &= 1 \quad \text{and} \\ p(\text{Rainy} | t=0) &= 0 \quad \text{i.e.} \end{aligned}$$

$$s_0 = [1, 0]$$

...now, what is s_1 ?

First example

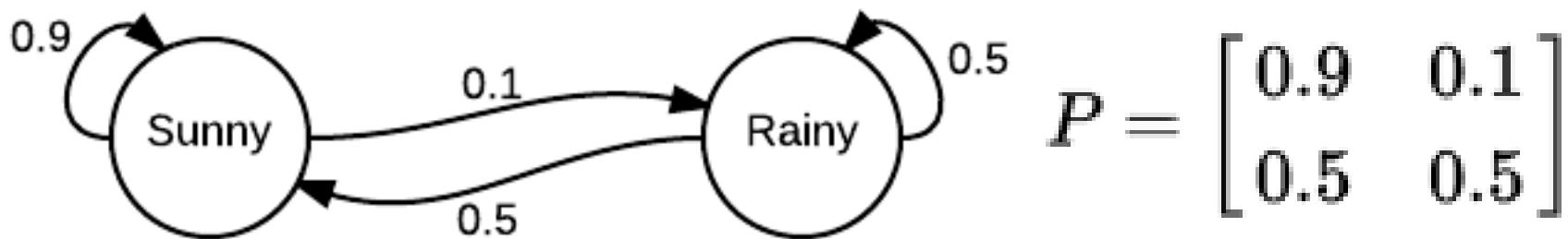


What is s_1 ?

$$s_1 \quad \left\{ \begin{array}{l} p(\text{Sunny} \mid t=1) = 0.9 \\ p(\text{Rain} \mid t=1) = 0.1 \end{array} \right.$$

...now, what is s_2 ?

First example



What is s_2 ?

$$s_2 \left\{ \begin{array}{l} p(\text{Sunny} | t=2) = 0.9 \quad p(\text{Sunny} | t=1) + 0.5 \quad p(\text{Rainy} | t=1) = 0.86 \\ p(\text{Rainy} | t=2) = 0.1 \quad p(\text{Sunny} | t=1) + 0.5 \quad p(\text{Rainy} | t=1) = 0.14 \end{array} \right.$$

In matrix form...

$$\mathbf{s}_2 = \mathbf{s}_1 P$$

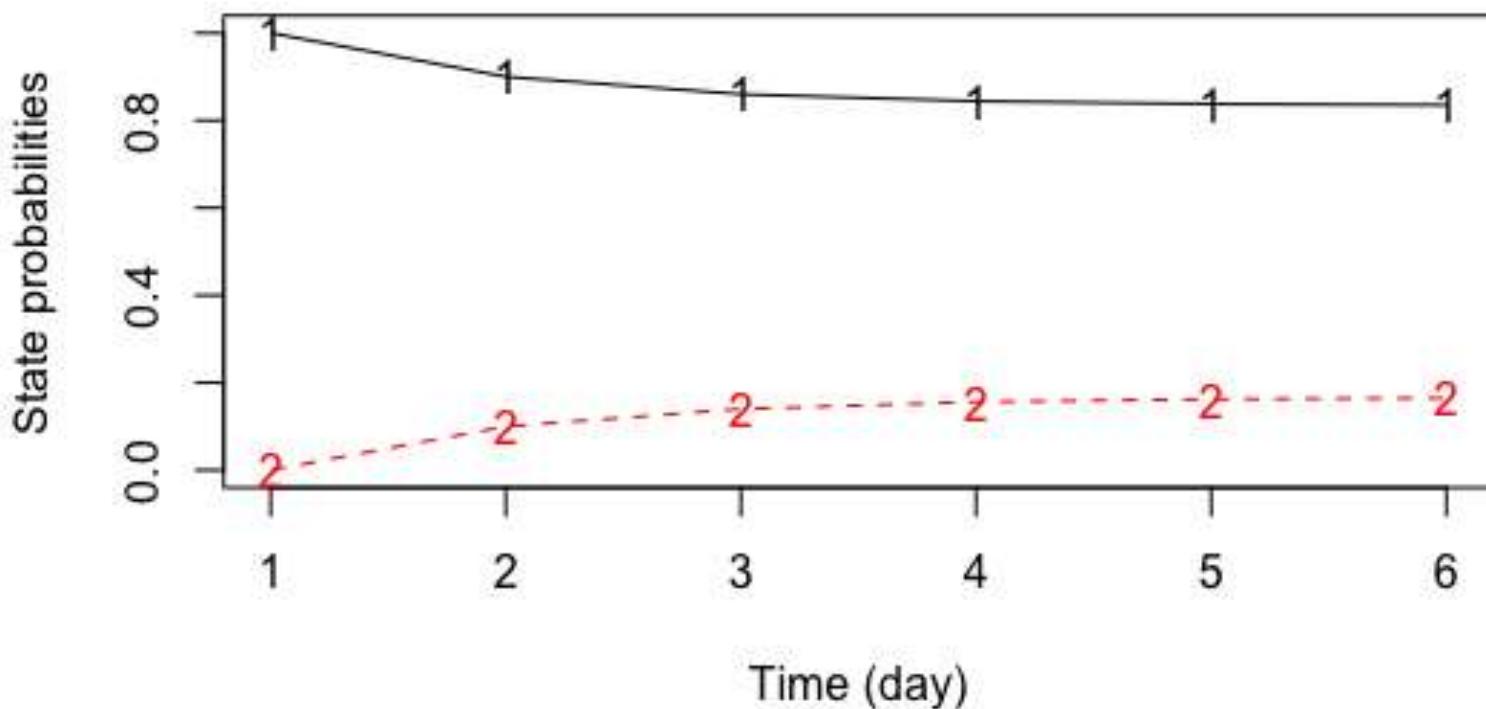
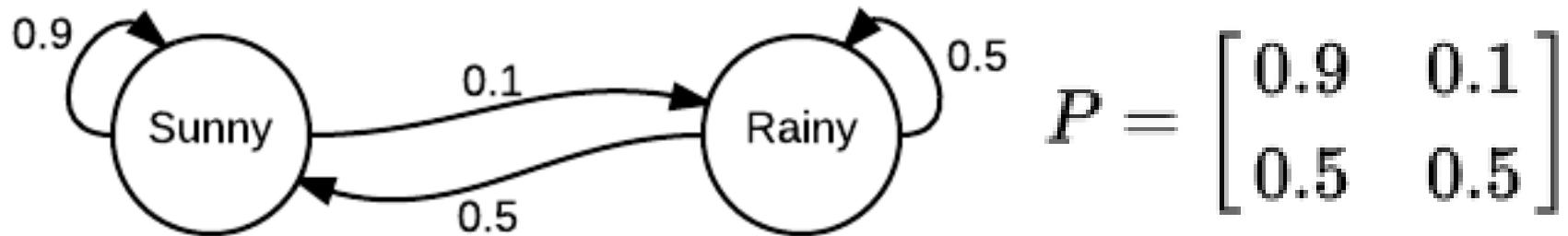
And in general...

$$\mathbf{s}_{t+1} = \mathbf{s}_t P$$

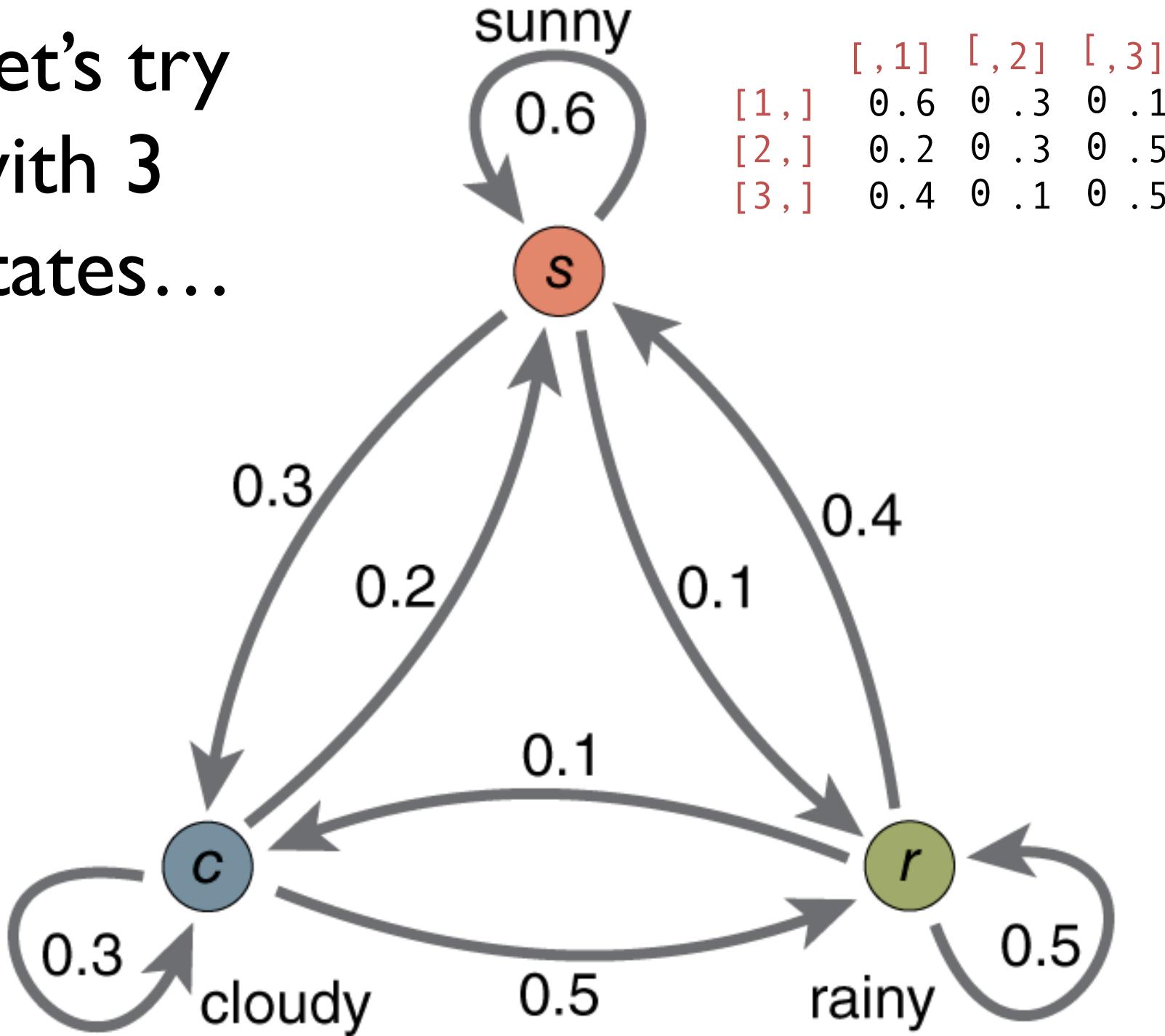
Meaning...

$$\mathbf{s}_t = \mathbf{s}_0 P^t$$

Let's do a numerical test...



Let's try
with 3
states...

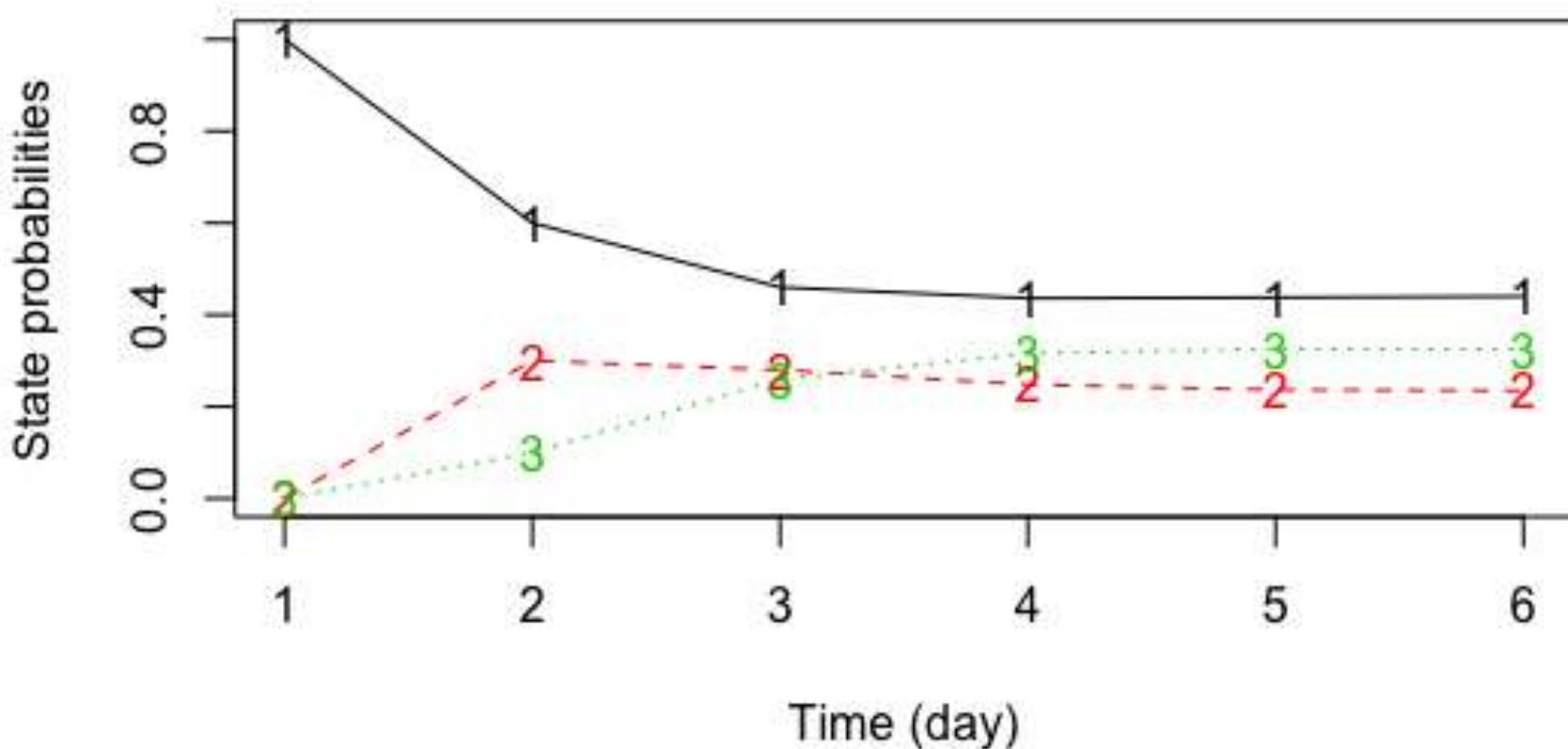


Starting from
Sunny...



	[, 1]	[, 2]	[, 3]
[1 ,]	0 . 6	0 . 3	0 . 1
[2 ,]	0 . 2	0 . 3	0 . 5
[3 ,]	0 . 4	0 . 1	0 . 5

Initial state: $s_0 = [1,0,0]$

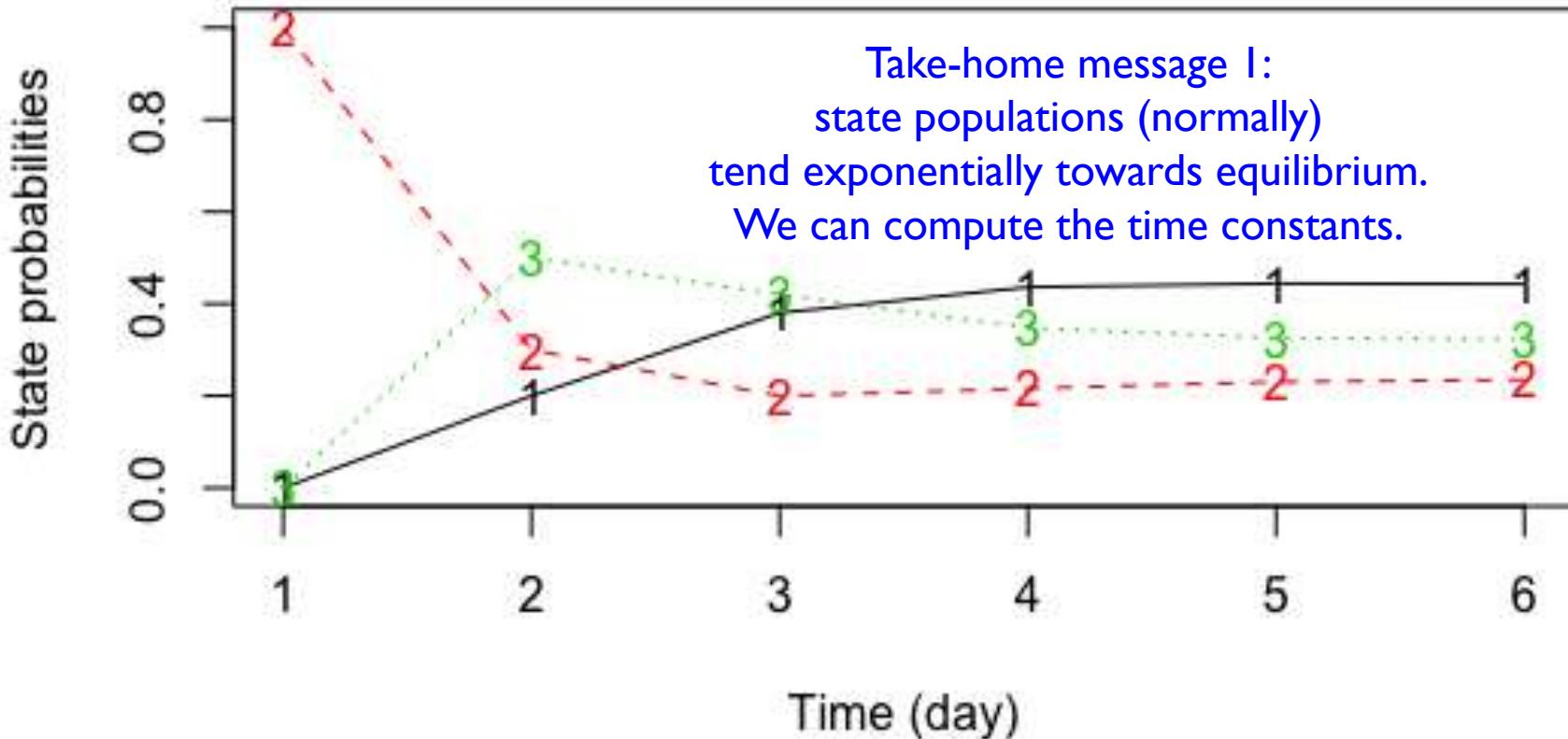


Starting from
Cloudy...



	[, 1]	[, 2]	[, 3]
[1 ,]	0 . 6	0 . 3	0 . 1
[2 ,]	0 . 2	0 . 3	0 . 5
[3 ,]	0 . 4	0 . 1	0 . 5

Initial state: $s_0 = [0,1,0]$



A couple of definitions

- State j is **accessible** from i ($i \rightarrow j$) if the system has a probability $\neq 0$ of going from one to the other
- j **communicates** with i iff $i \rightarrow j$ and $j \rightarrow i$
- A set of states C is a **communicating class** if every pair of states in C communicates with each other.
- A Markov chain is said to be **irreducible** if its state space is a single communicating class.
- An **absorbing state**, once entered, cannot be left.
- (... and many more)

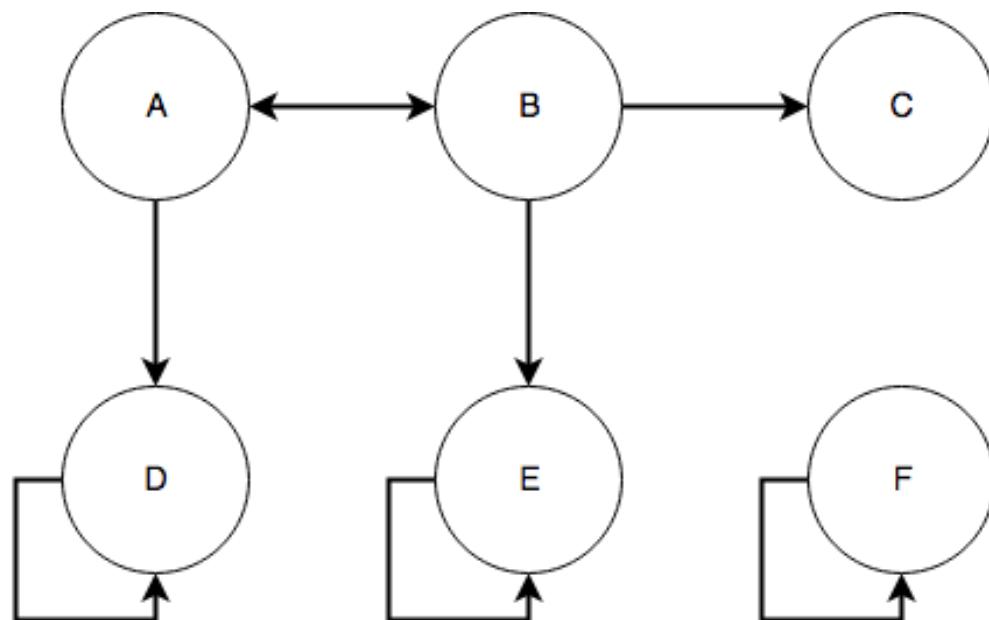


Homogeneity

- If the transition probabilities do not change with time, we have an *homogeneous* Markov chain
- Example of non-homogeneous MC: weather transition probabilities Markovian, but dependent on the season.

Absorbing states

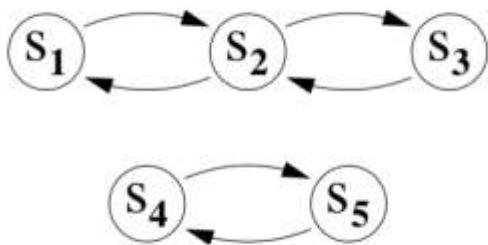
- A state is called **absorbing** if it is impossible to leave it.
- In other words, $p_{jj}=1$ (implying $p_{ij}=0$ for $i \neq j$)



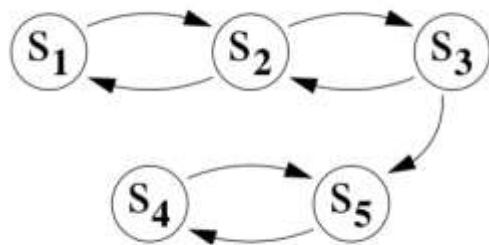
A, B are recurrent
C, D, E, F are absorbing

Irreducibility

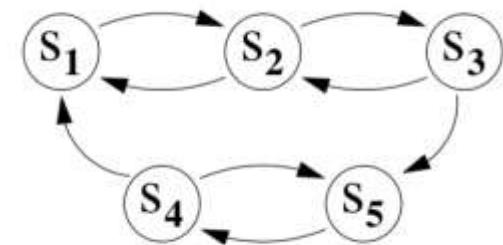
- A Markov chain is said to be **irreducible** if it is possible to get from any state to any state



Not irreducible



Not irreducible



Irreducible

Important quantities we can compute

- Stationary distribution (\rightarrow eq. probabilities)
- Relaxation times
- Mean-first passage times (\rightarrow kinetic rates)
- And others we won't discuss:
 - Commitor probabilities
 - Fluxes
 - ...

**Learning transition matrices
from trajectories
(of discrete states)**

Markov models

Total sampled time (e.g. 100 μs)

Define the *discrete* state of the system in
discrete time (e.g. via reaction coordinates)



Lag time τ (here: 4 time units)

Compute the transition probability matrix
sliding a window of lag time τ

	To		
	0	0	0
From	0	0	0
	0	0	0



	To		
	0	0	1
From	0	0	0
	0	0	0

Markov models

Total sampled time (e.g. 100 μs)



Define the *discrete* state of the system in
discrete time (e.g. via reaction coordinates)



Lag time τ (here: 4 time units)



Compute the transition probability matrix
sliding a window of lag time τ

		To	
		0	1
From	0	0	1
	1	0	0
	2	0	0



		To	
		0	1
From	0	0	1
	1	1	0
	2	0	0

Markov models

Total sampled time (e.g. 100 μs)



Define the *discrete* state of the system in
discrete time (e.g. via reaction coordinates)



..
Lag time τ (here: 4 time units)



Compute the transition probability matrix
sliding a window of lag time τ

		To	
		0	1
From	0	0	1
	1	1	0
	2	0	0



		To	
		0	1
From	0	0	1
	1	1	0
	2	1	0

Markov models

Total sampled time (e.g. 100 μs)



Define the *discrete* state of the system in
discrete time (e.g. via reaction coordinates)



Lag time τ (here: 4 time units)

Compute the transition probability matrix
sliding a window of lag time τ

To

	Red	Green	Blue
Red	3	0	2
Green	1	3	0
Blue	1	2	1

From



Repeat until the end of the trajectory.

Note the Markovian assumption:
transition probabilities **do not** depend
on history (neither, for us, on time).

(Note how “Early” and “late” events
are squashed in the same matrix.)

Important: you can accumulate counts
from **different** trajectories in the **same** matrix!

It's all valid, independent sampling
of the same underlying system.

Example: multiple weather stations in the same city.
They may be active at different times!

Transition counts

		To		
		3	0	2
From	1	3	0	0
	2	1	2	1

Transition probabilities

		To		Σ_j	
		3/5	0	2/5	1
From	1	$\frac{1}{4}$	$\frac{3}{4}$	0	1
	2	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{1}{4}$	1

Normalize
by rows

$$P_{ij}$$

Probability vector

1	0	0
---	---	---

$$\times \begin{pmatrix} & & \\ & & \\ & & \\ \text{3/5} & 0 & 2/5 \\ \frac{1}{4} & \frac{3}{4} & 0 \\ \frac{1}{4} & 2/4 & \frac{1}{4} \end{pmatrix}$$

$$s_i$$

$$P_{ij}$$

Evolved (after τ) state

3/5	0	2/5
-----	---	-----

$$s'_j$$

Probability vector

1	0	0
---	---	---

s_i

$$\times \begin{pmatrix} 3/5 & 0 & 2/5 \\ 1/4 & 3/4 & 0 \\ 1/4 & 2/4 & 1/4 \end{pmatrix} =$$

P_{ij}

Evolved state

3/5	0	2/5
-----	---	-----

s_j'

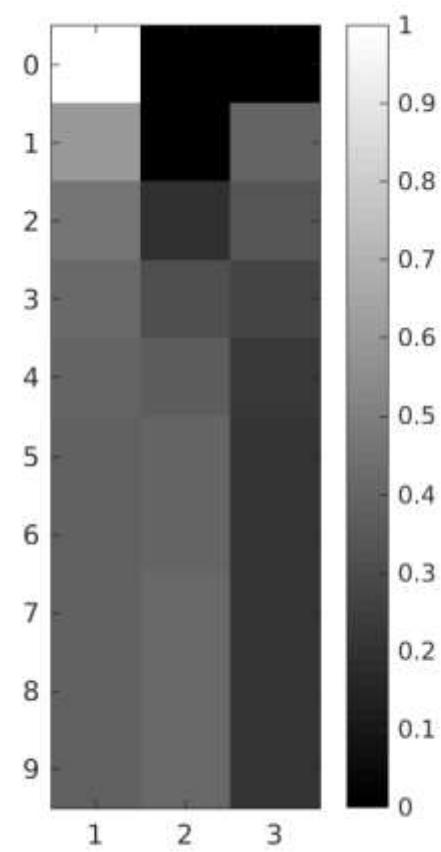
$$s' = sP$$

$$s'' = (sP)P = sP^2$$

$$s^{(n)} = sP^n$$

1	0	0
.60	0	.40
.46	.20	.34
.41	.32	.27
.39	.37	.23
.39	.40	.22
...

time



Probability vector

1	0	0
---	---	---

x

3/5	0	2/5
¼	¾	0
¼	2/4	¼

Evolved state (after tau)

3/5	0	2/5
-----	---	-----

s_i

P_{ij}

s'_j

$$s^\infty P = s^\infty$$

$$s' = sP$$

$$s'' = (sP)P = sP^2$$

$$s^{(n)} = sP^n$$

...

$$s^\infty = ?$$

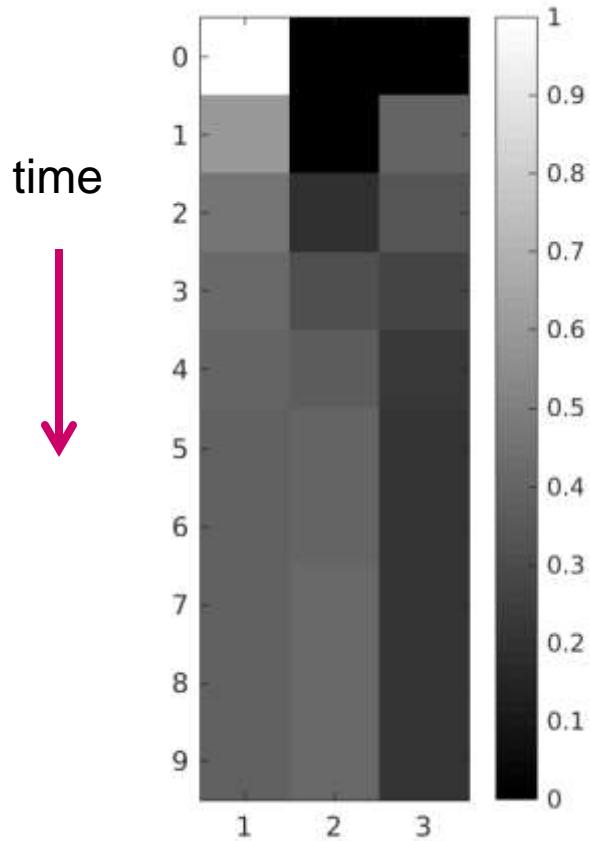
Left eigenvector of P (eigenvalue 1)

Is the stationary state
= equilibrium probabilities
= the free energy surface

[a,b]=eig(P')
a(:,3)/sum(a(:,3))
= 0.385 0.410 0.205 = [5/13 16/39 8/39]

Relaxation towards equilibrium

Equilibrium is reached within typical **relaxation times***.



Do an eigen-decomposition...

$$\mathbf{x}VP = \Lambda \mathbf{x}V$$

... with $\Lambda = \text{diag}(\mu_i)$

$$\mathbf{y} \equiv \mathbf{x}V$$

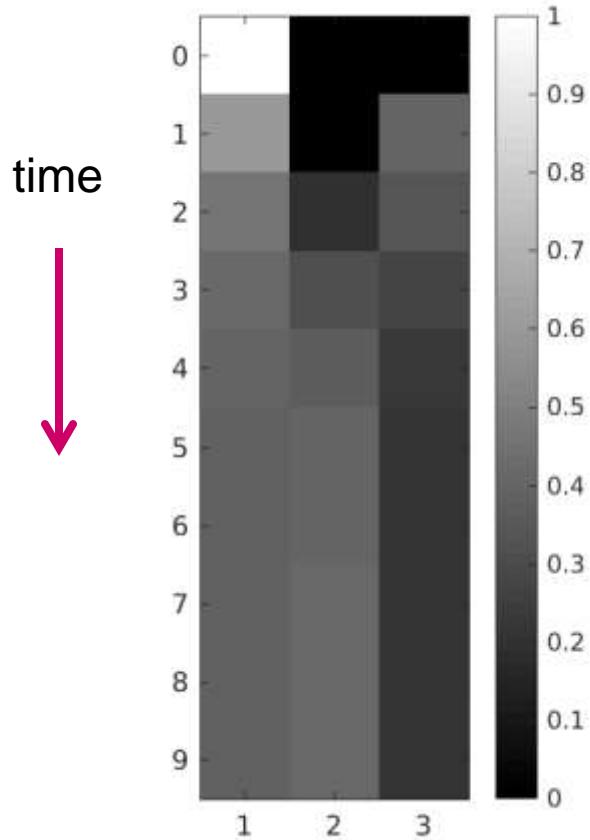
$$\mathbf{y}(t+1) = \Lambda \mathbf{y}(t)$$

$$y_i(t) = y_i(0)\mu_i^t = y_i(0)e^{t \log \mu_i}$$

*

Relaxation towards equilibrium

Equilibrium is reached within typical **relaxation times**.



They are computed from the eigenvalues $< 1 \dots$

$$\tau_k = -\frac{\tau}{\ln \mu_k(\tau)}$$

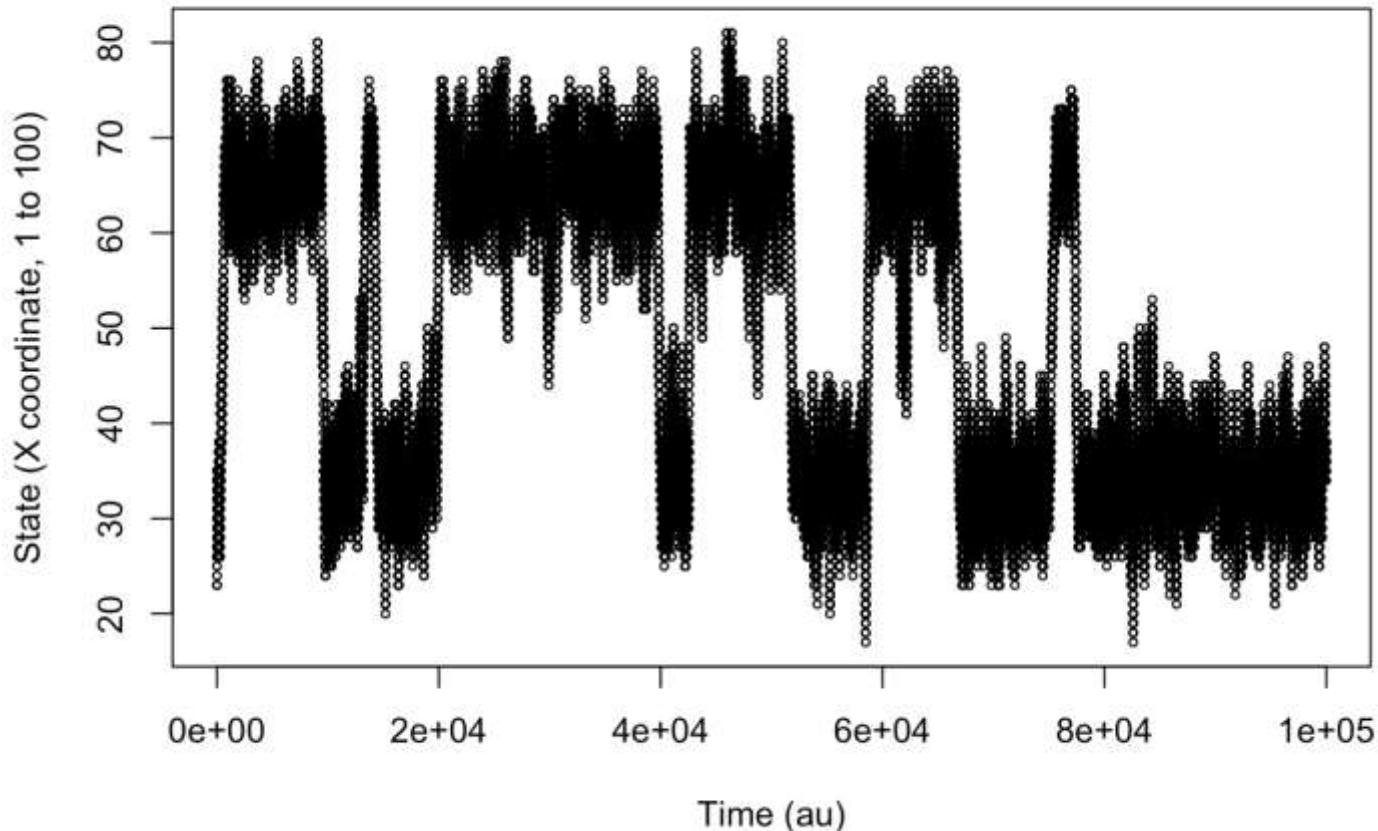
...and called *implied timescales* (they depend on τ)

Markov modeling a 1D trajectory

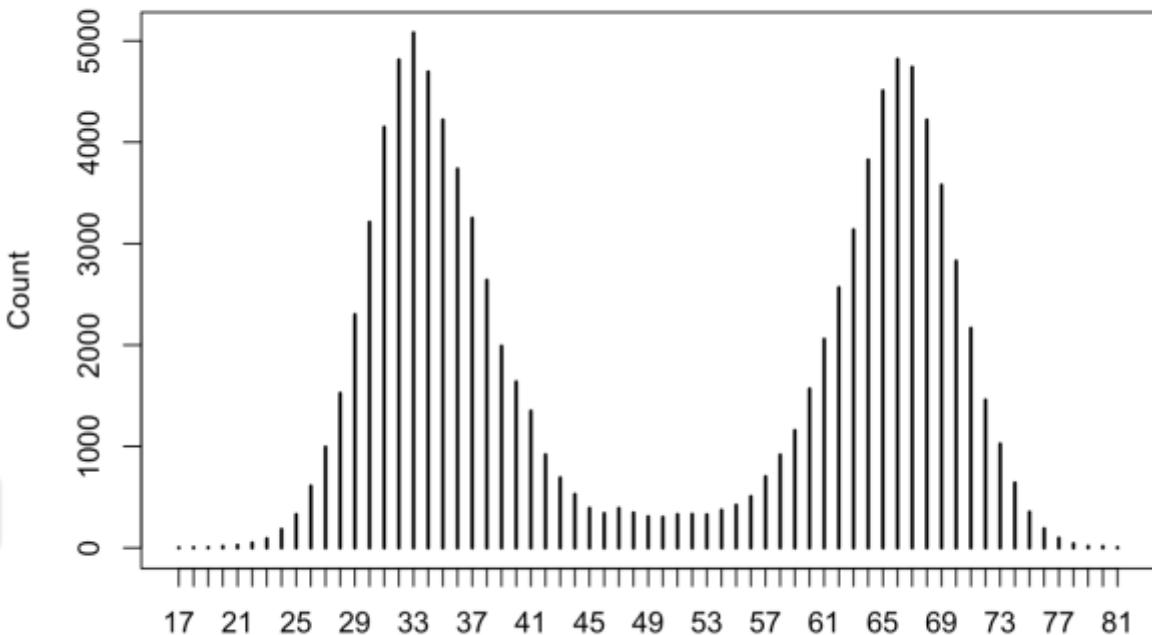
(Please find the extended version online,
“Markov state models of a 1D trajectory”,
with R code)

Start with a 1-D trajectory

- Already discretized in 100 bins

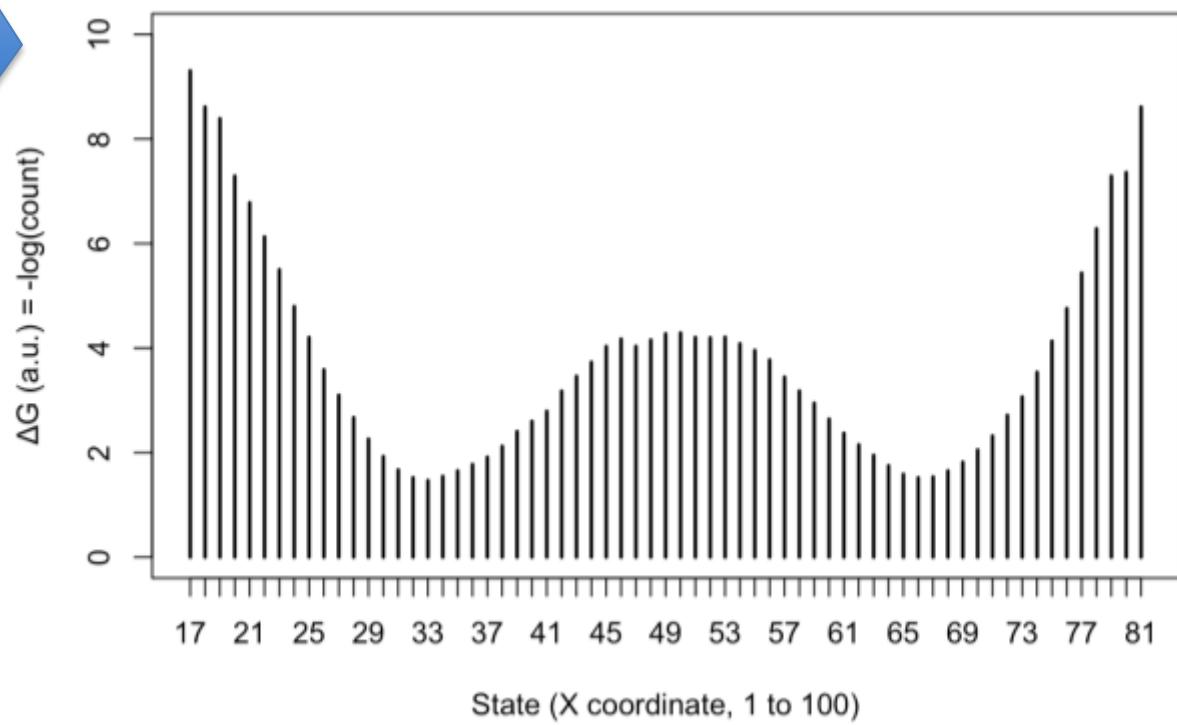


Histogram



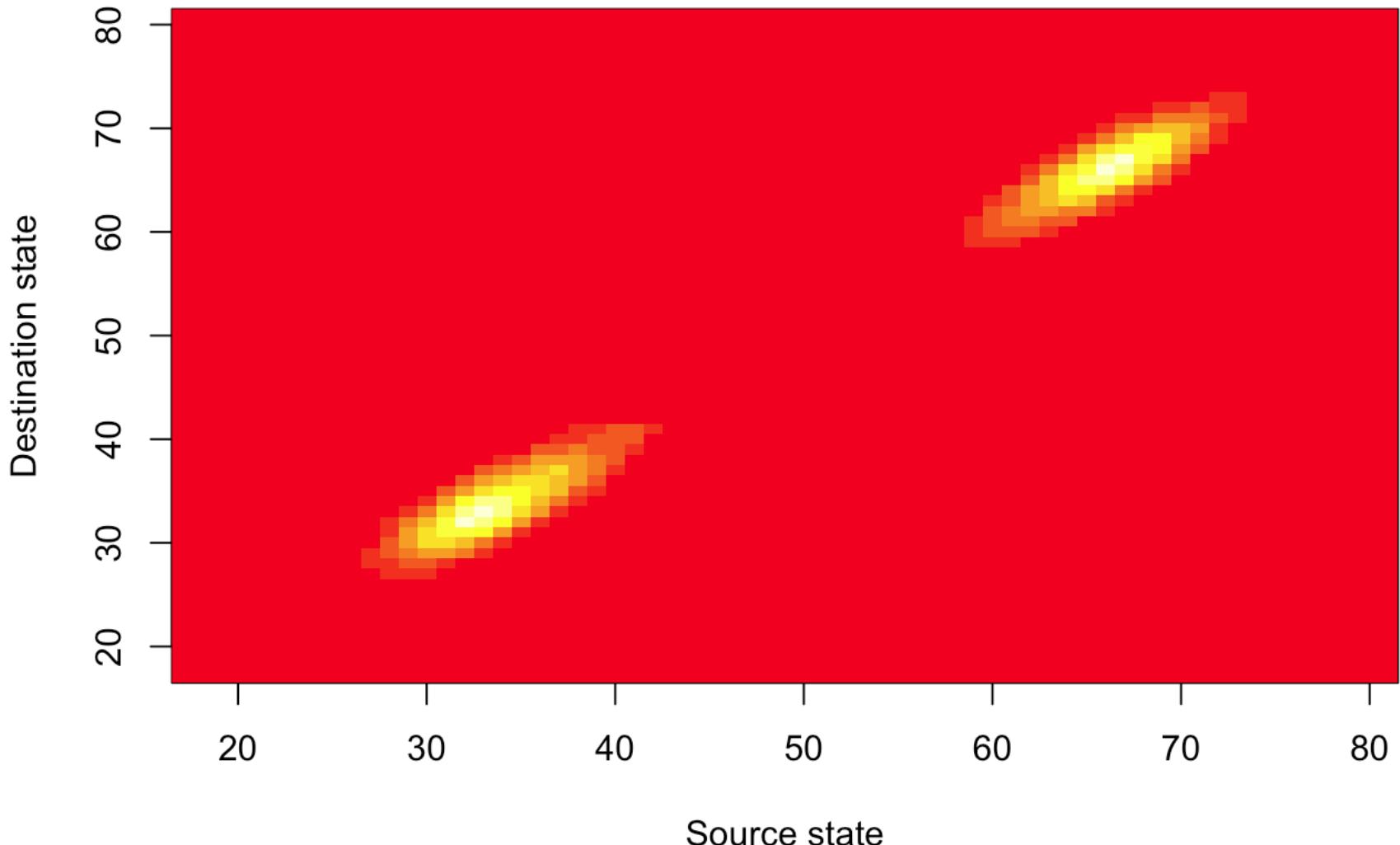
Boltzmann inversion

-log(count)



The transition count matrix*

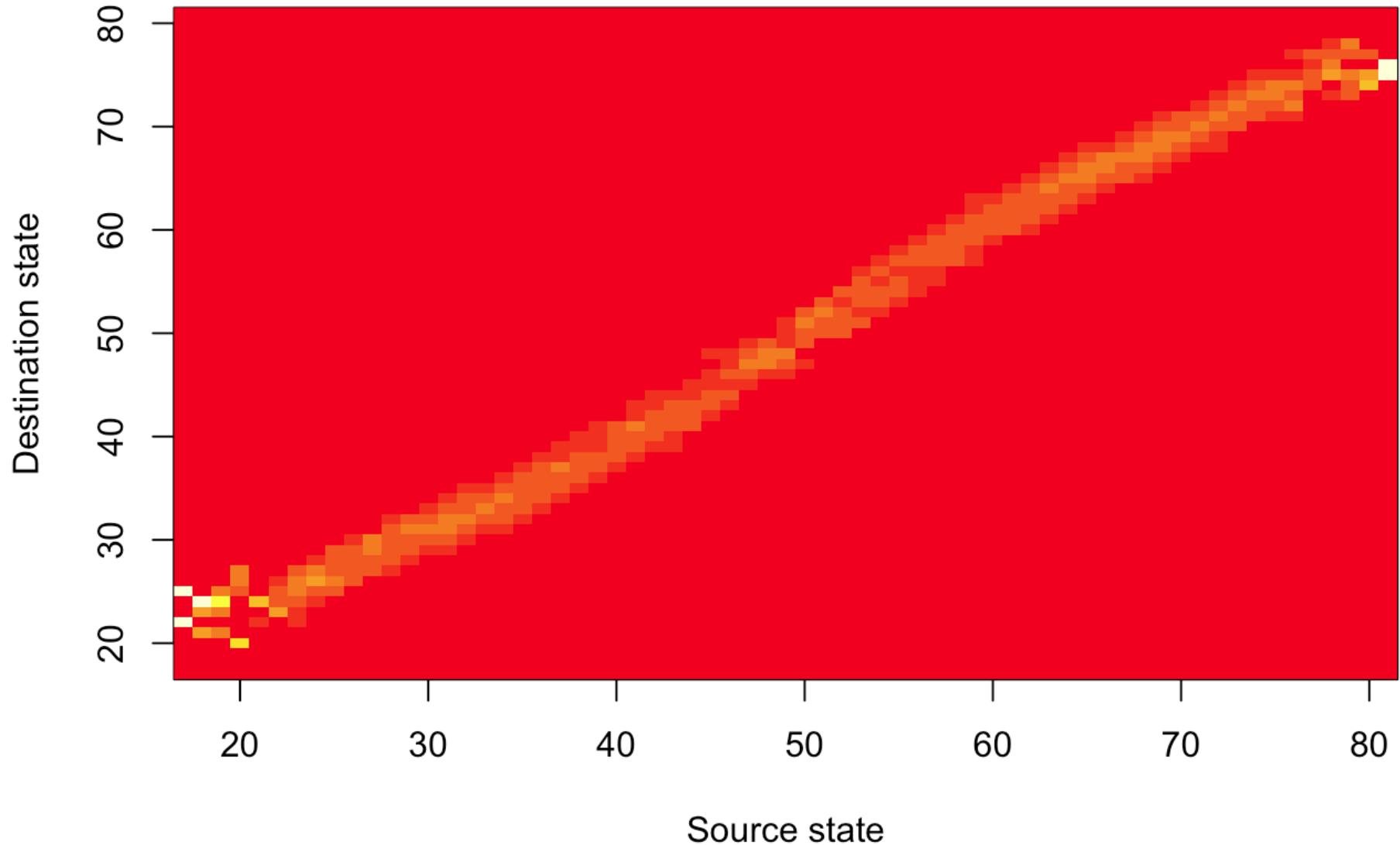
How many times we have seen state i going to j after $\tau=10$ time units



* shown as an image for compactness

The transition *probability* matrix

Rows normalized to sum to 1



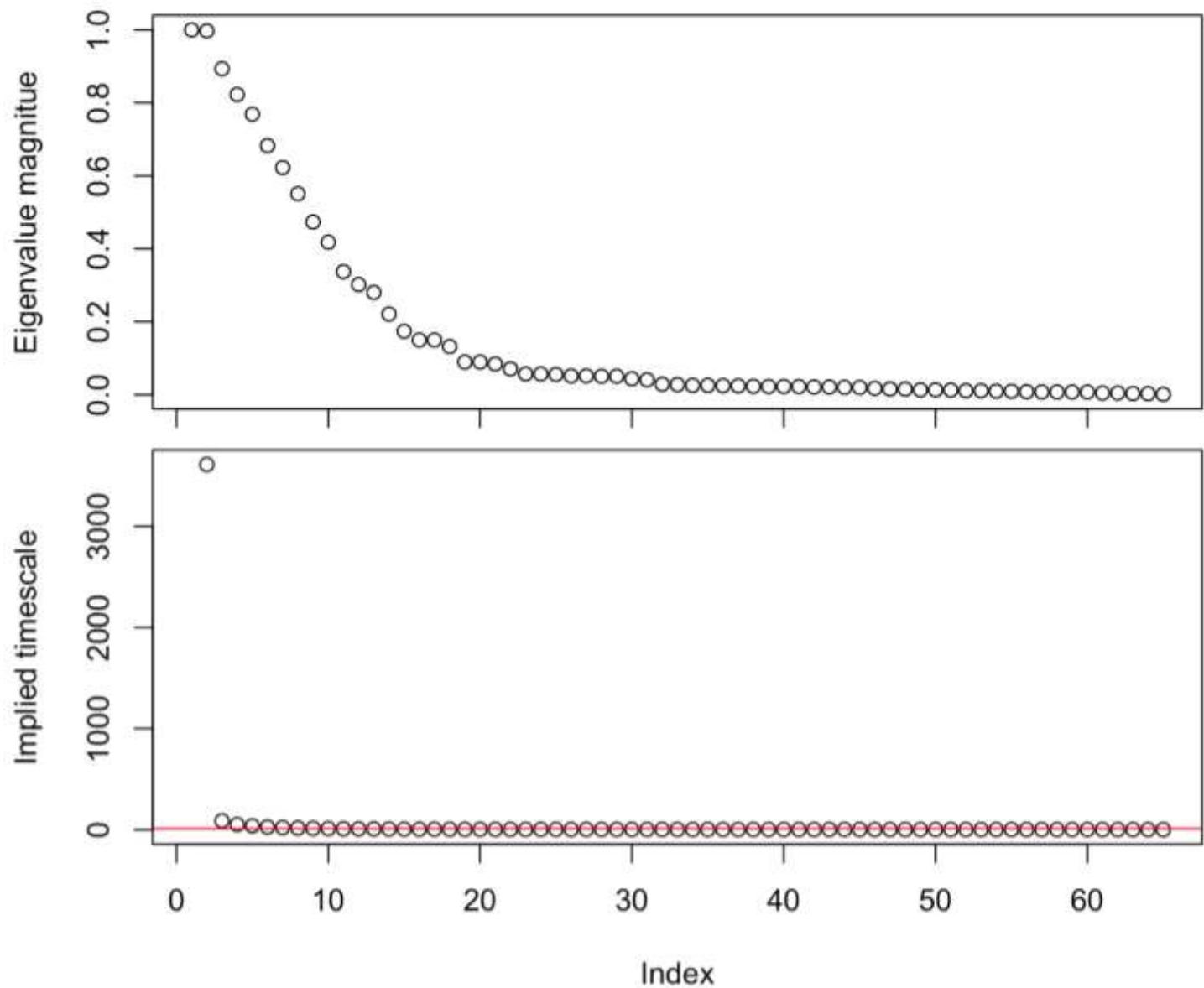
Eigenvalues (at $T=10$)

Eigenvalues μ_i



Implied
timescales

$-T/\log(\mu_i)$



Eigenvalues (at $T=10$)

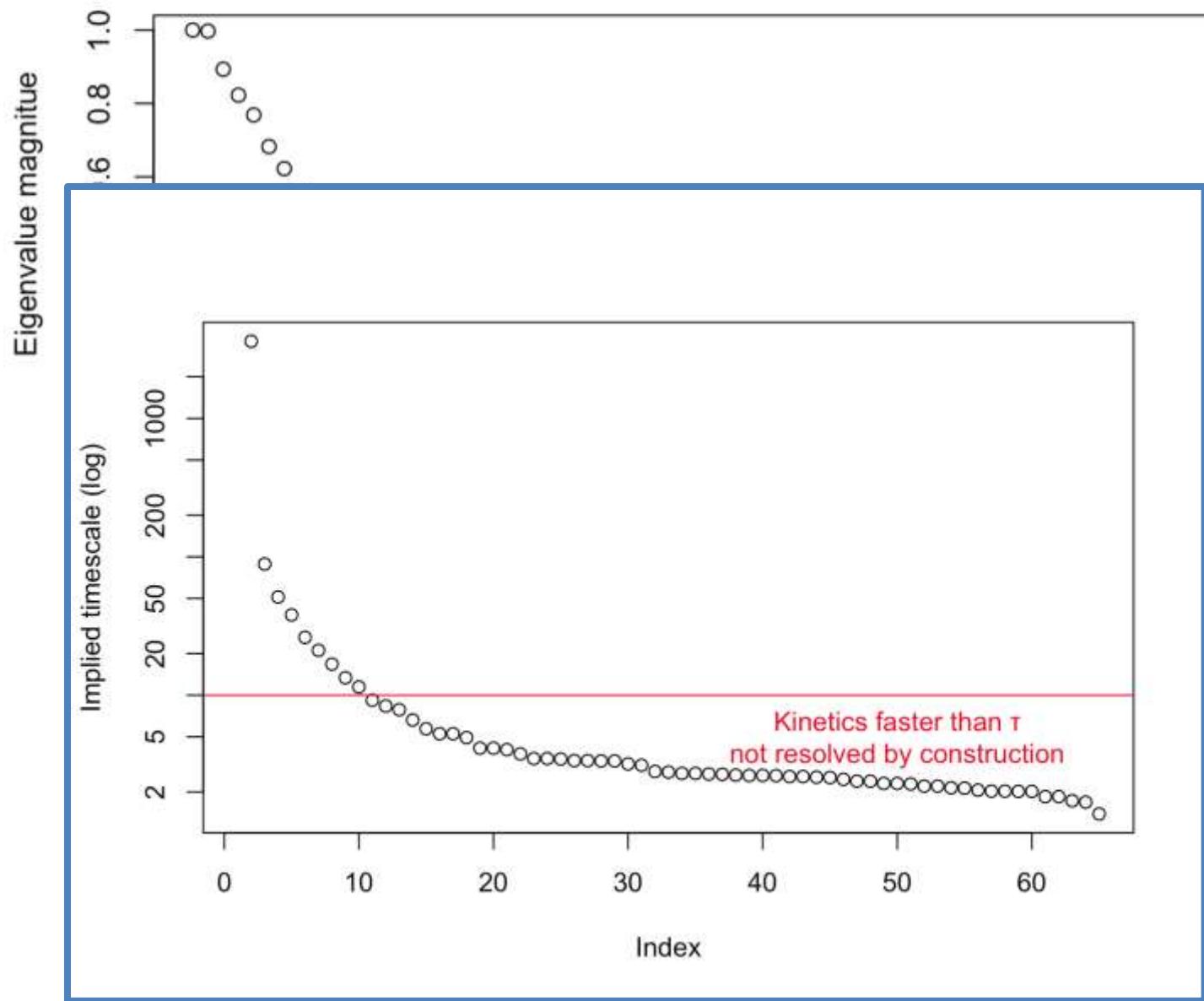
Eigenvalues μ_i



Implied timescales

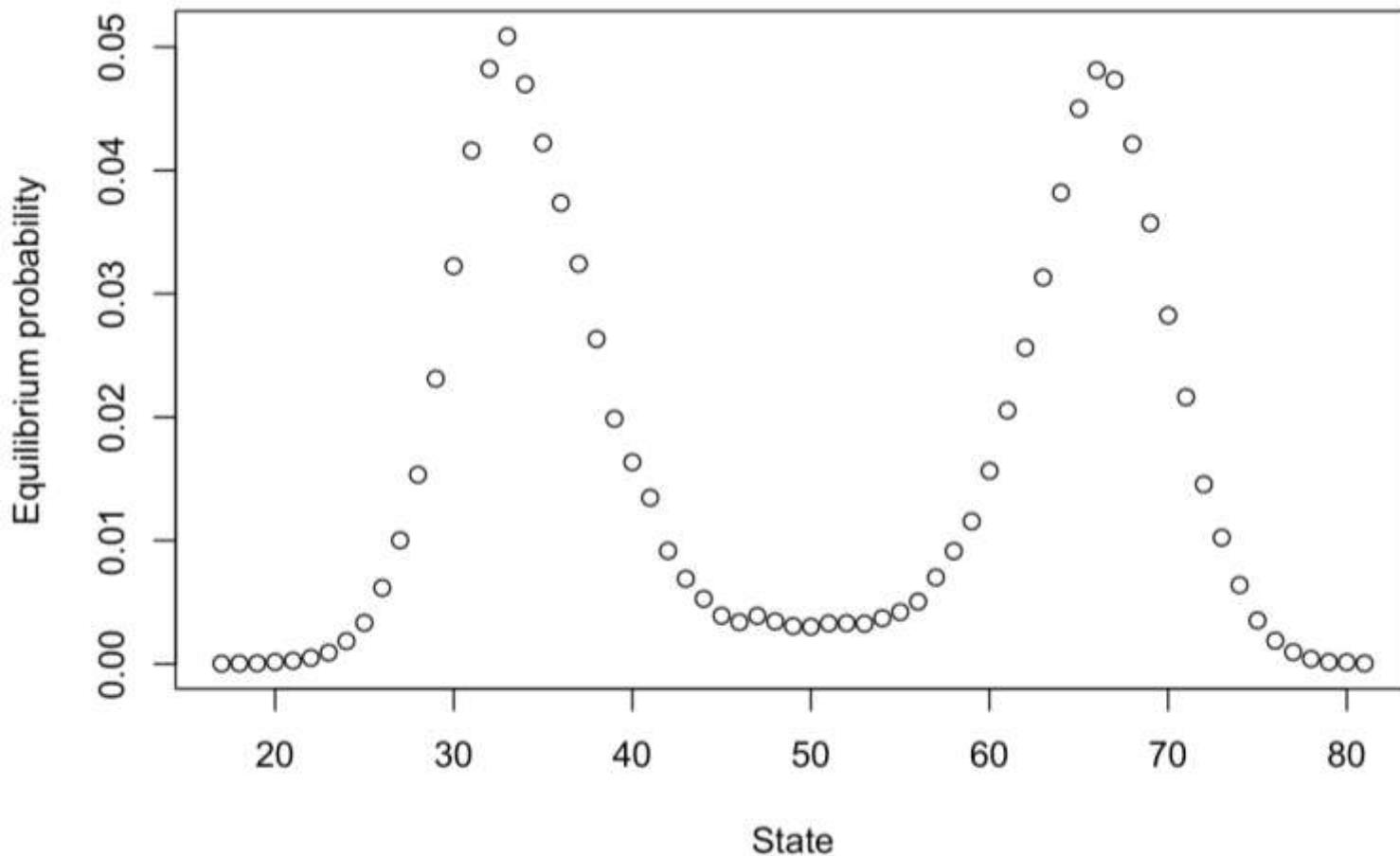
$-T/\log(\mu_i)$

Now in log scale!



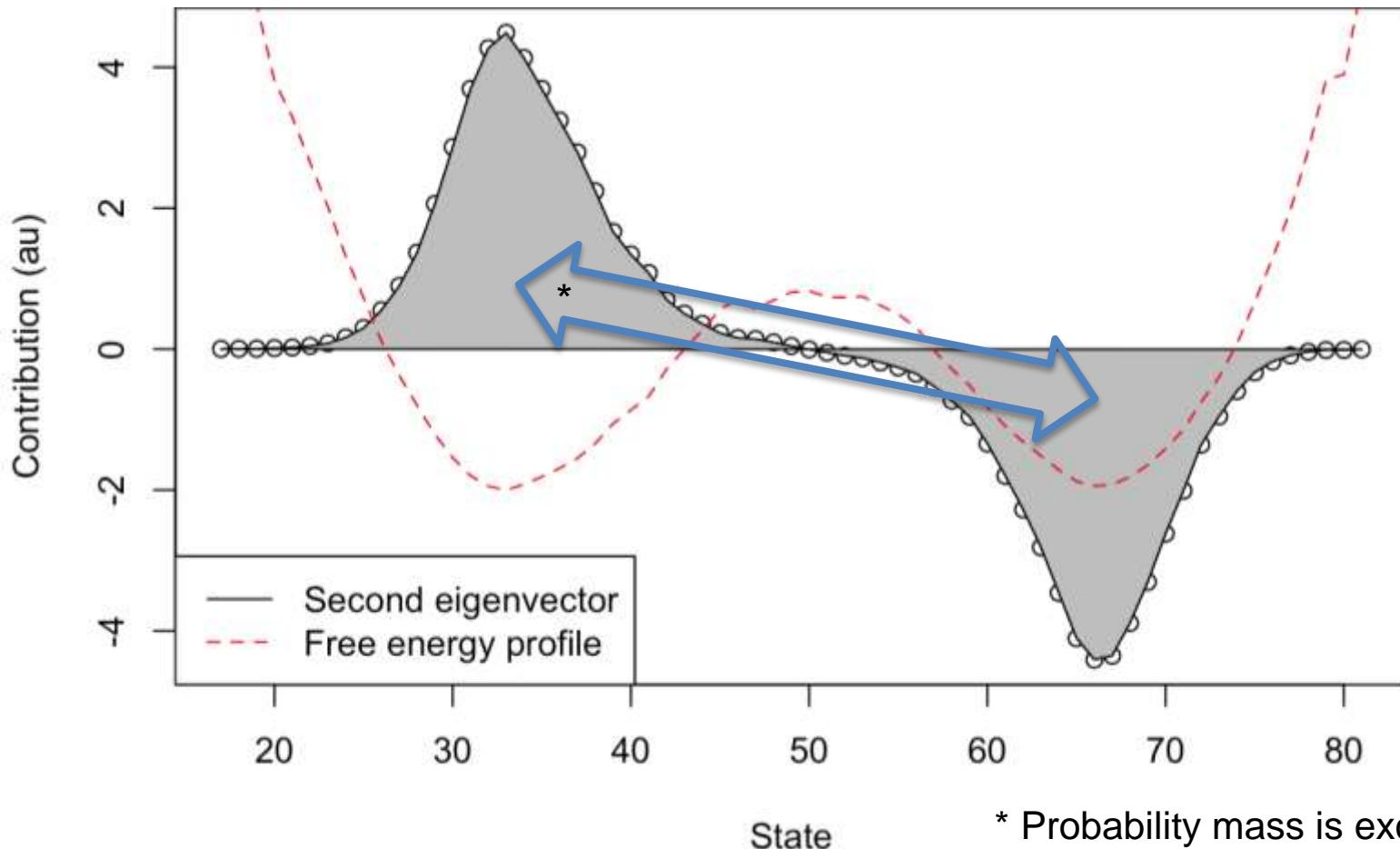
First eigenvector ($\mu_1=1$)

This is the stationary state (normalize so it sums to 1)



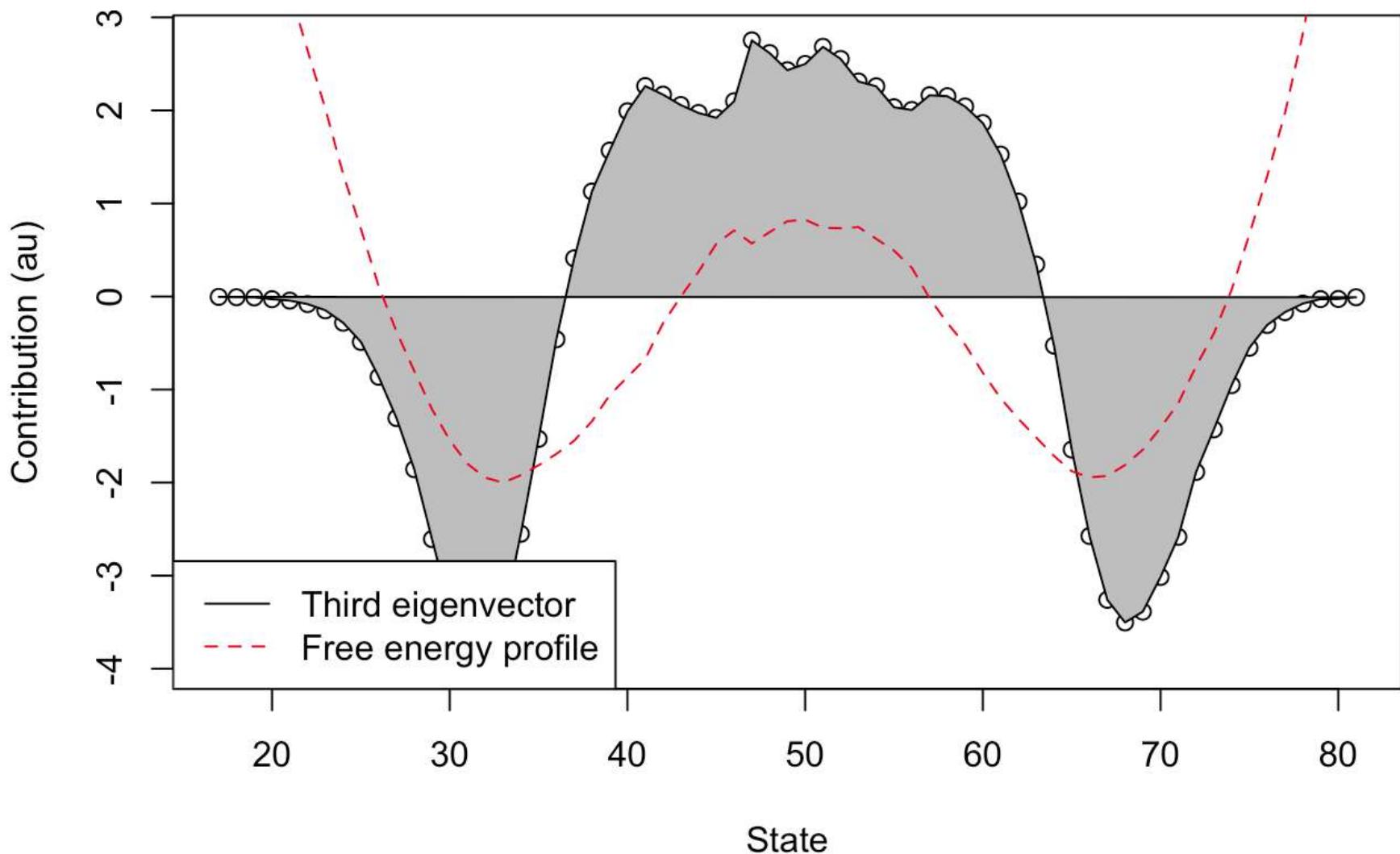
Second eigenvector ($\mu_2=0.997$)

This is the slowest relaxation mode: ITS $\tau_2 = 3610$ time units



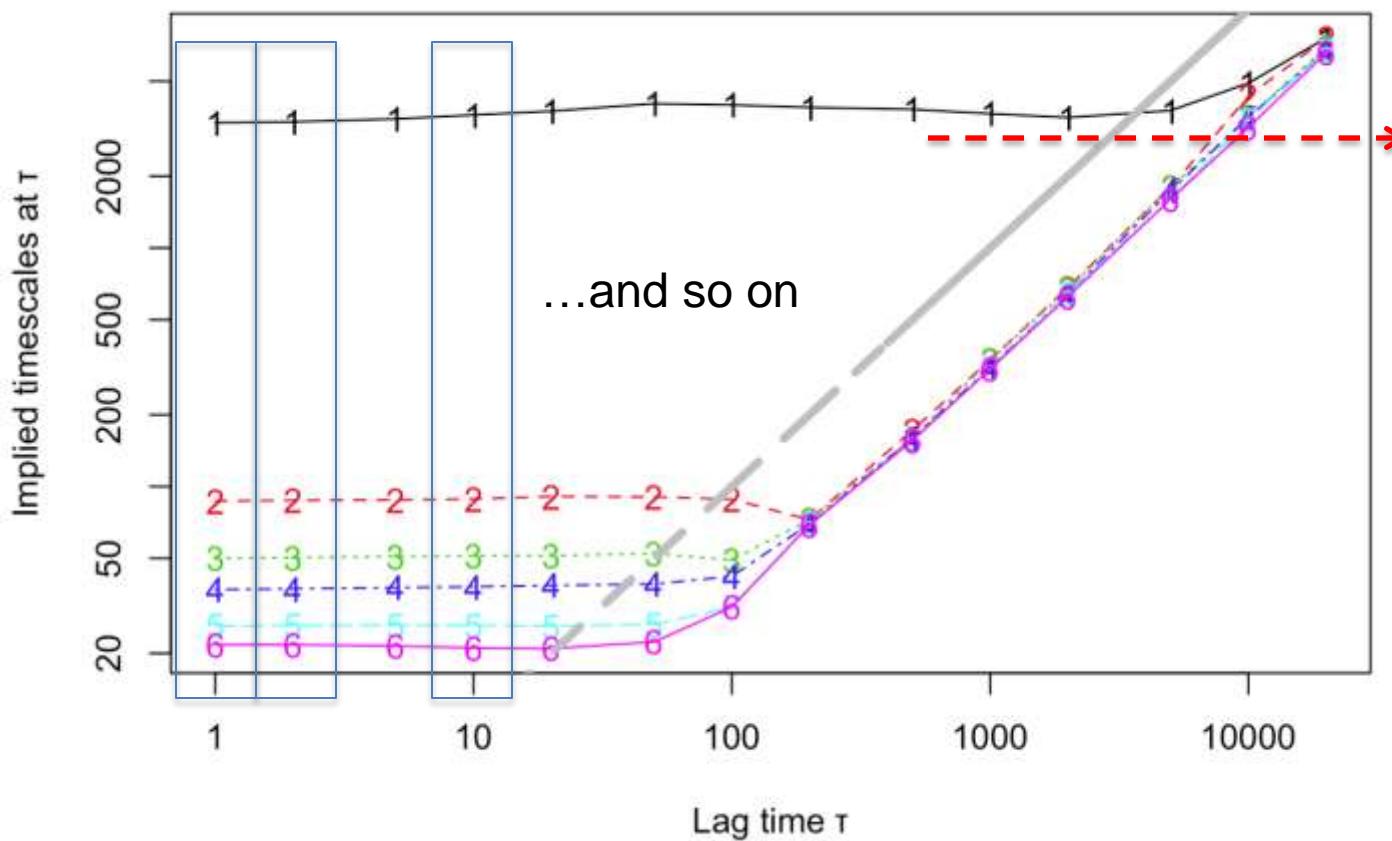
3rd eigenvector ($\mu_3=0.893$)

This is the slowest relaxation mode: ITS $\tau_3 = 88.7$ time units



Implied timescales plot

Repeat the eigenvalues determination for several lags. Check convergence.



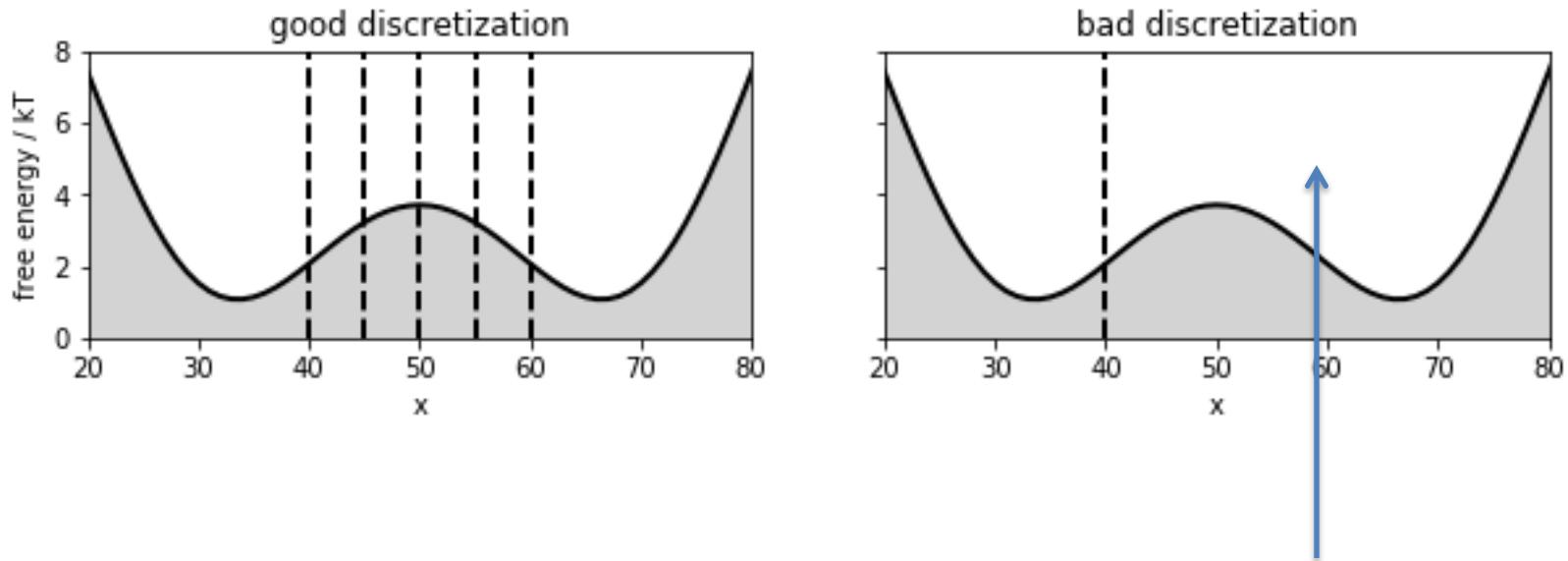
Here,
convergence is
achieved very
early.

Reasons:

- (a) true two-state dynamics;
- (b) absence of orthogonal degrees of freedom;
- (c) fine space discretization

Macrostates

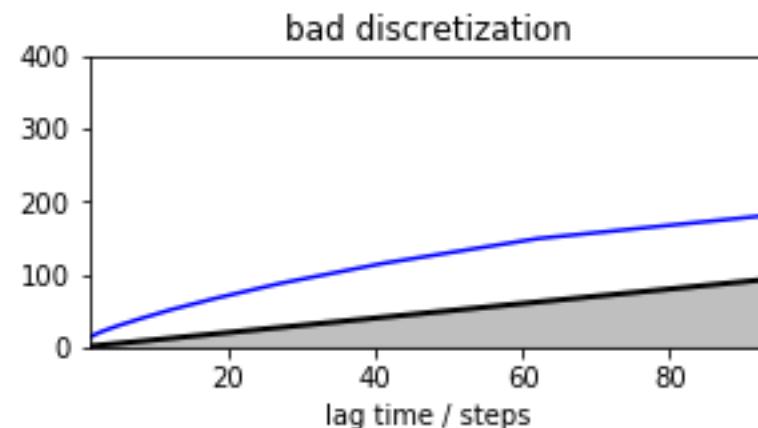
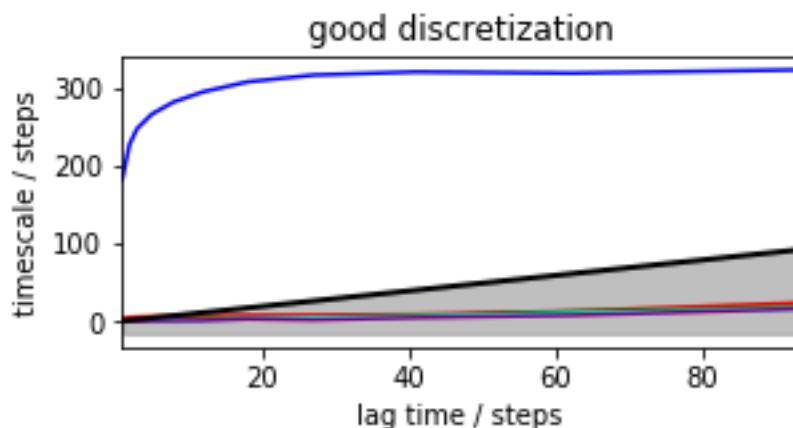
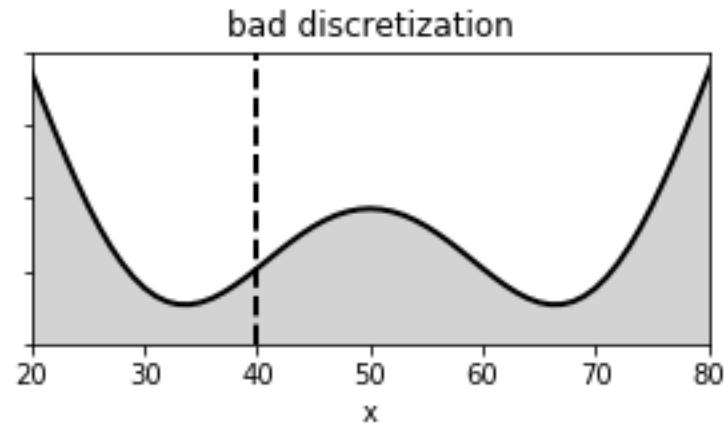
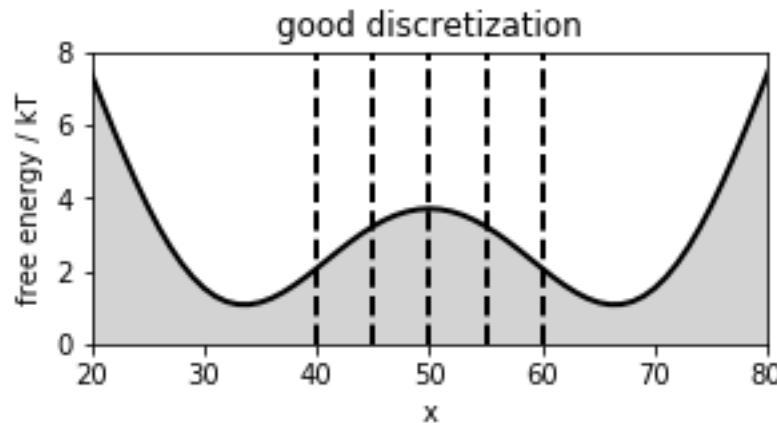
A bad choice of the discretization breaks the Markovianity assumption



In the “bad discretization” case, the barrier is embedded in one of the states. This generates a “long term memory” effect: the rightmost state could actually be short-lived (if we are on the left of the barrier) or long-lived (if we are on its right). These two cases are convoluted into the same, so that the present state information itself is not sufficient to predict the “future” of the system any more.

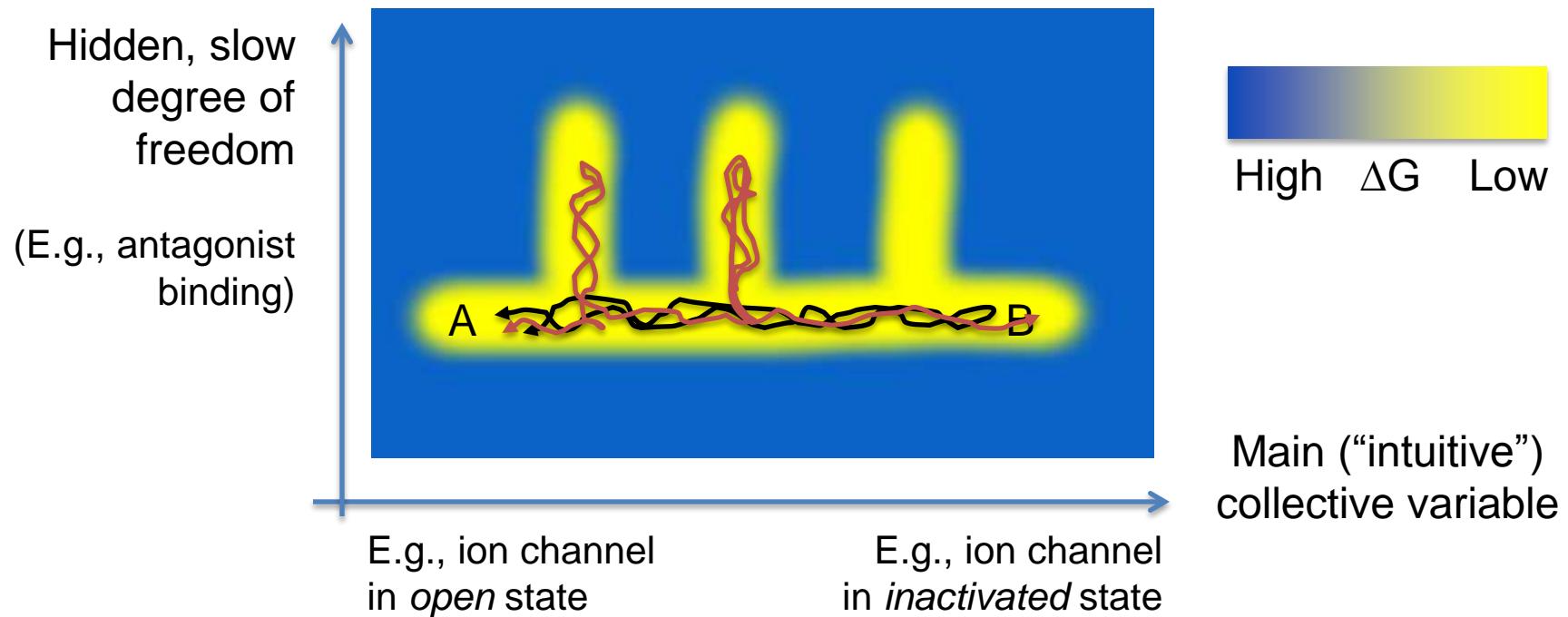
Macrostates

A bad choice of the discretization breaks the
Markovianity assumption



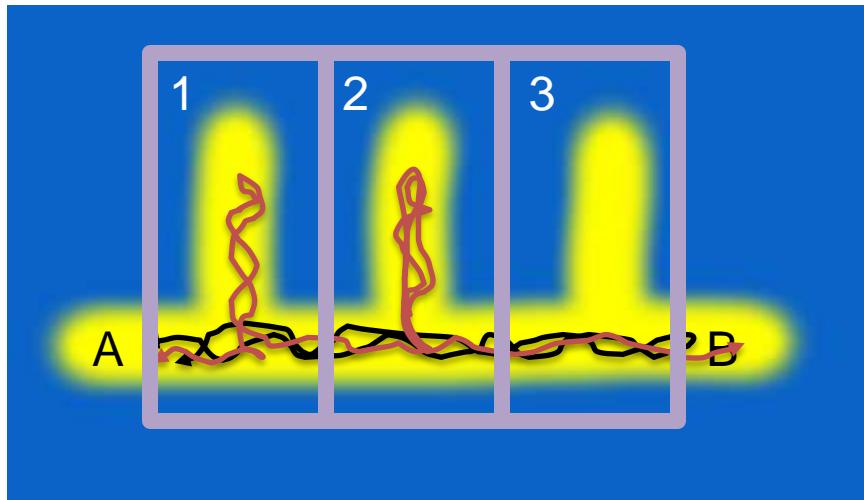
Orthogonal degrees of freedom

- Ideally, an intuitive collective degree of freedom would be a good reaction coordinate
- This is unfortunately rarely the case. E.g...

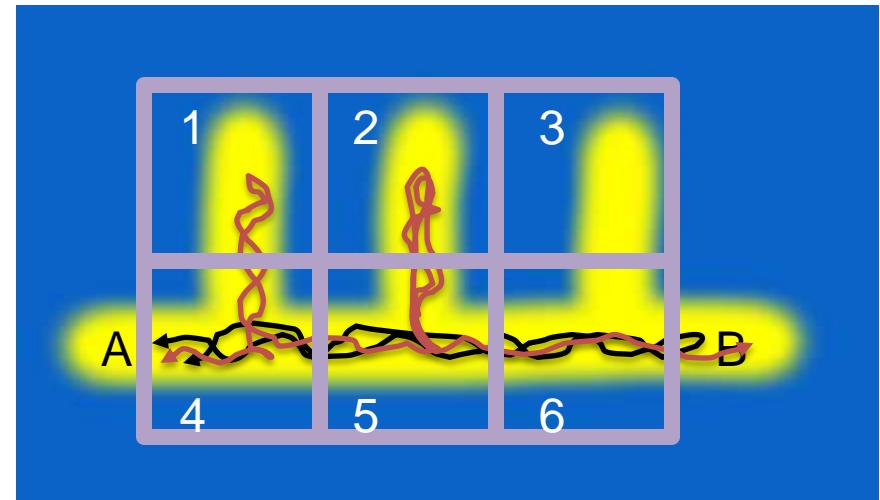


Orthogonal degrees of freedom

- Biased sampling: include bias along the “vertical” axis (i.e. make it part of the reaction coordinate)
- MSM: make distinct states
 - Otherwise, implied timescales plot will show non-convergence*



Not ok (violates Markovianity)

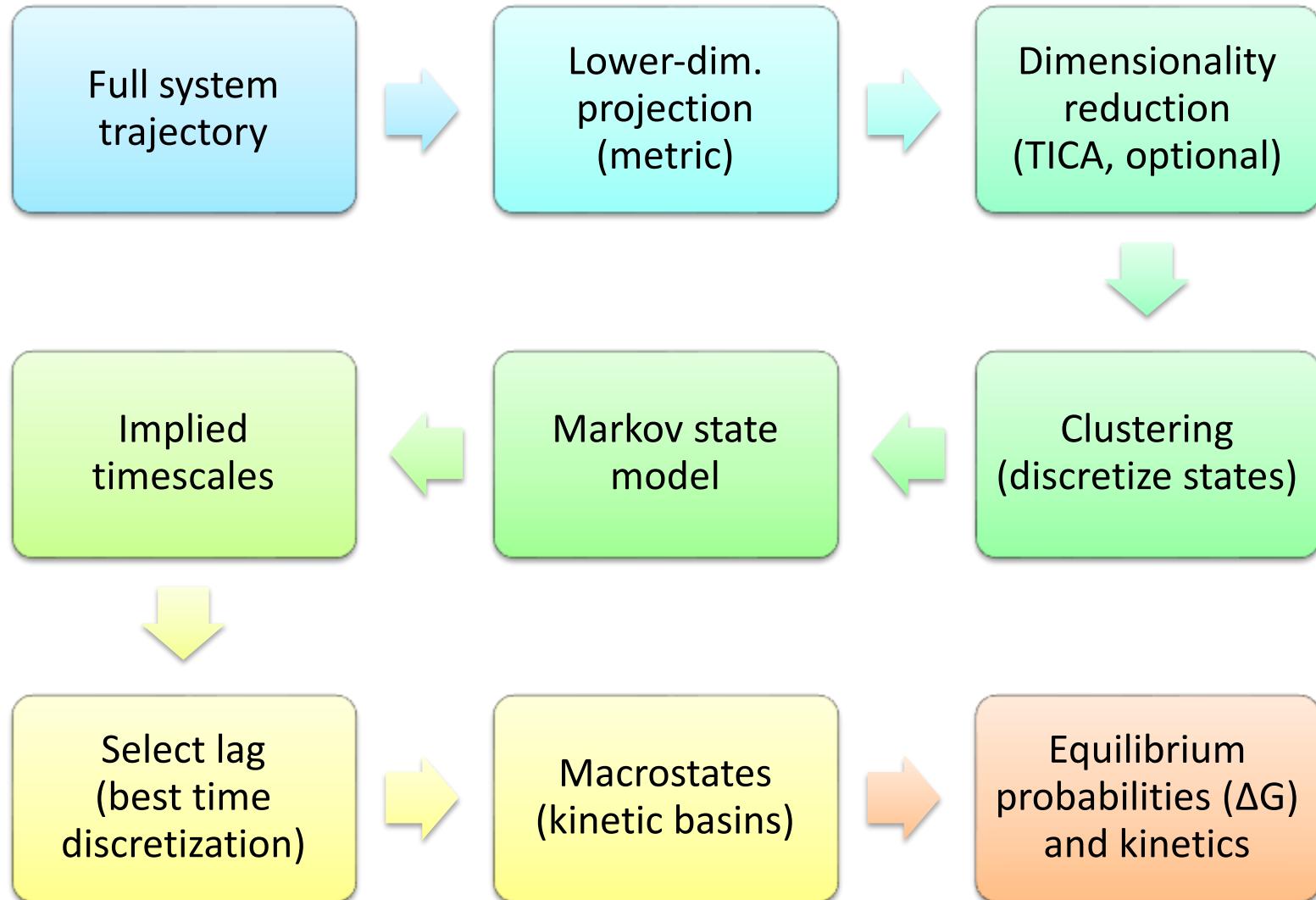


Hopefully better

* Better tests are available, e.g. Chapman-Kolmogorov

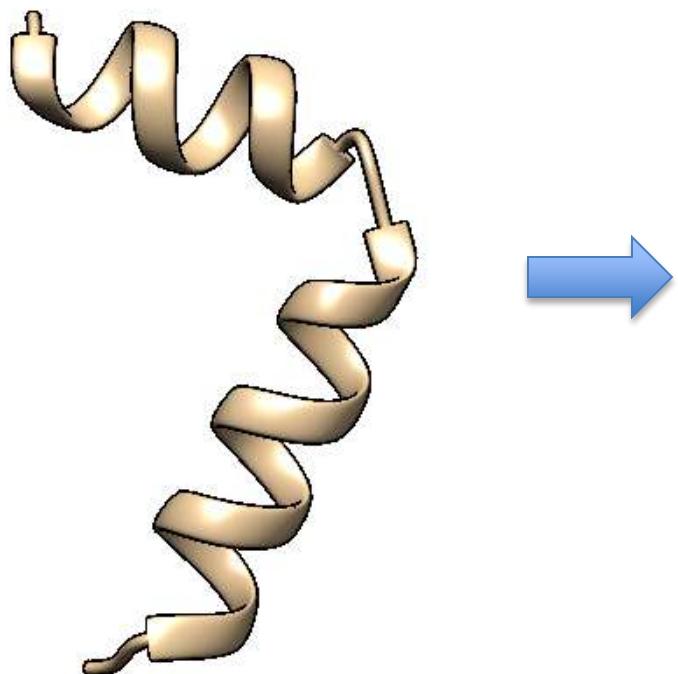
**Back to molecules:
>> 1 dimension**

MSM-based analysis overview



Metric projection

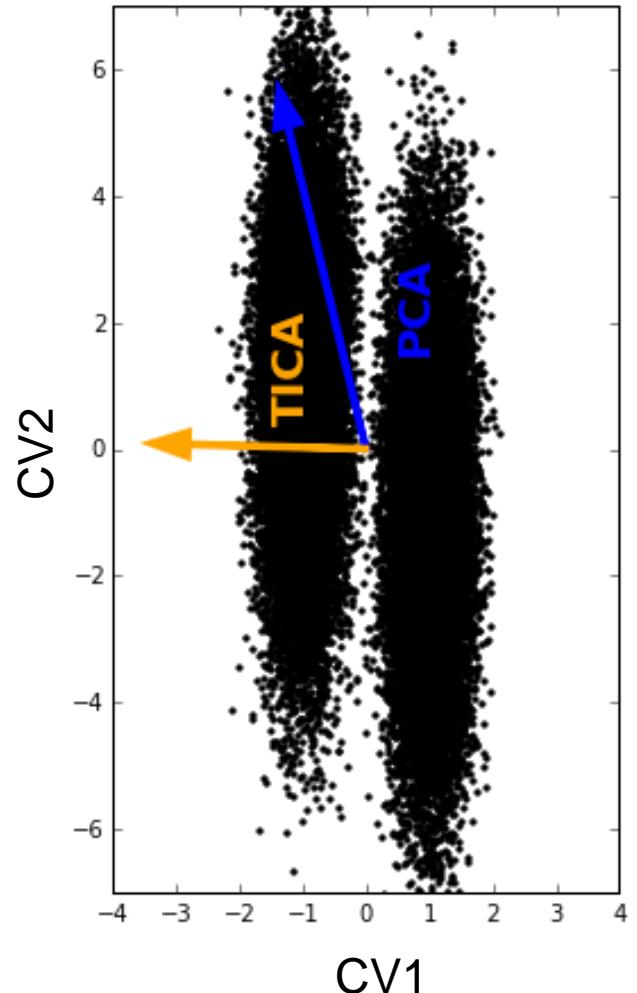
The first step is to project the system state in a lower-dimensional space (“metric”). Many choices are available, e.g.



- Manually chosen distances
- Atom coordinates
- N phi/psi Ramachandran angles
- Distance matrix
- Contact matrix
- ...

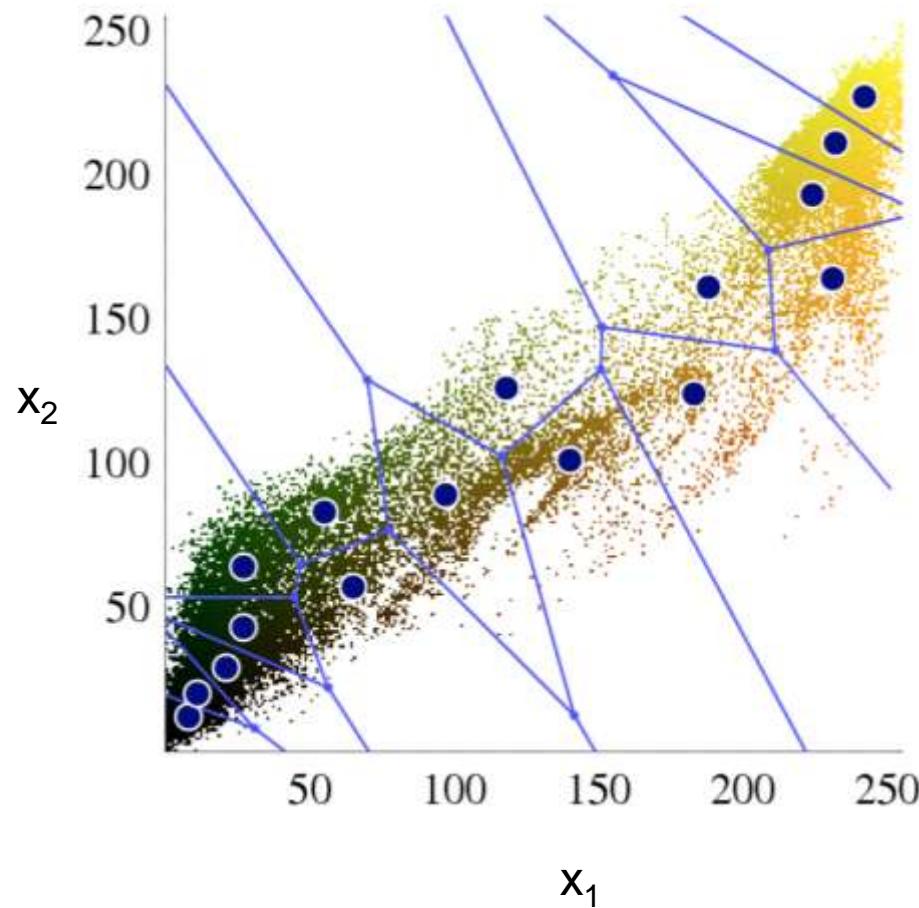
Time-lagged independent component analysis

- The feature space dimensionality may still be too high
- Do a further low-dimensional projection on the “slow” degrees of freedom:TICA
- It is based on lagged autocorrelation
- Contrast with PCA, which fits the “most elongated” ellipsoid, ignoring time

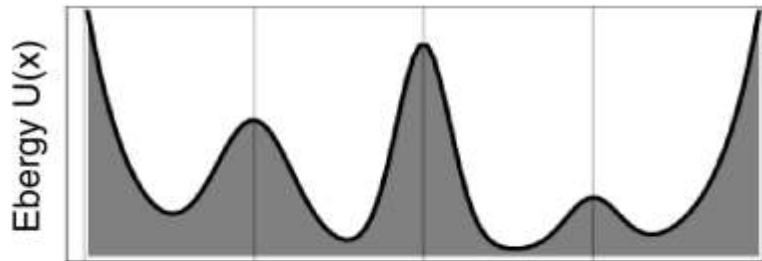


Clustering (1st level)

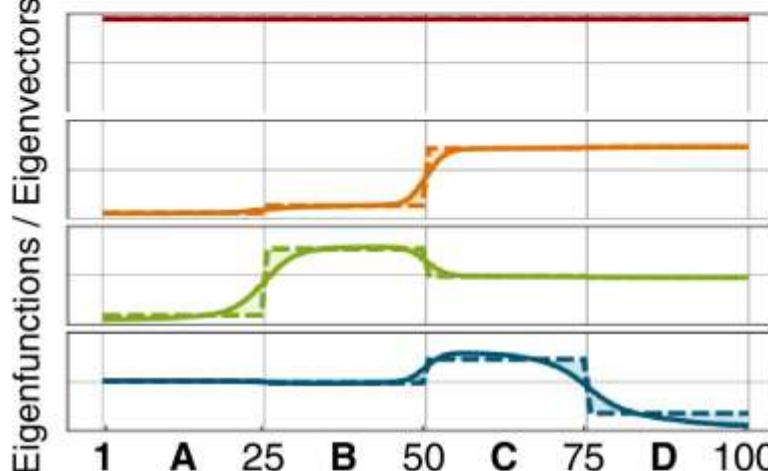
- Now move to a discrete state space
- Reduction from the low-dimensional space is done by *clustering*
 - Usually called “microstates”
- Several algorithms are implemented in MSM packages (e.g. grid; k-means; etc. We won’t discuss them)



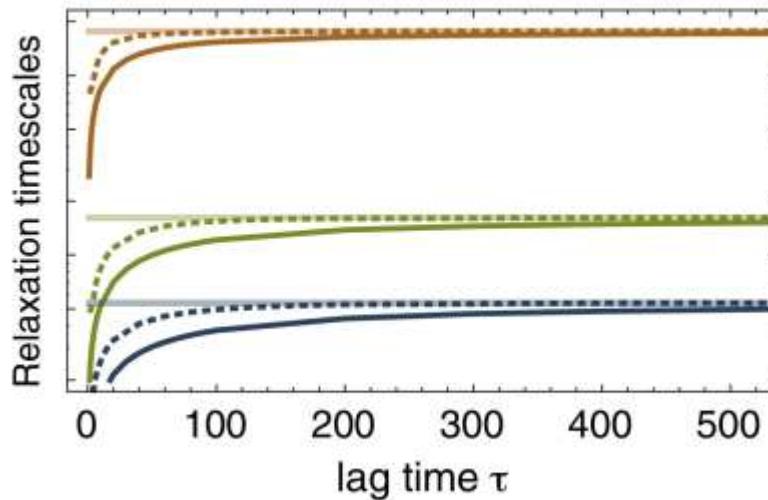
(a)



(b)

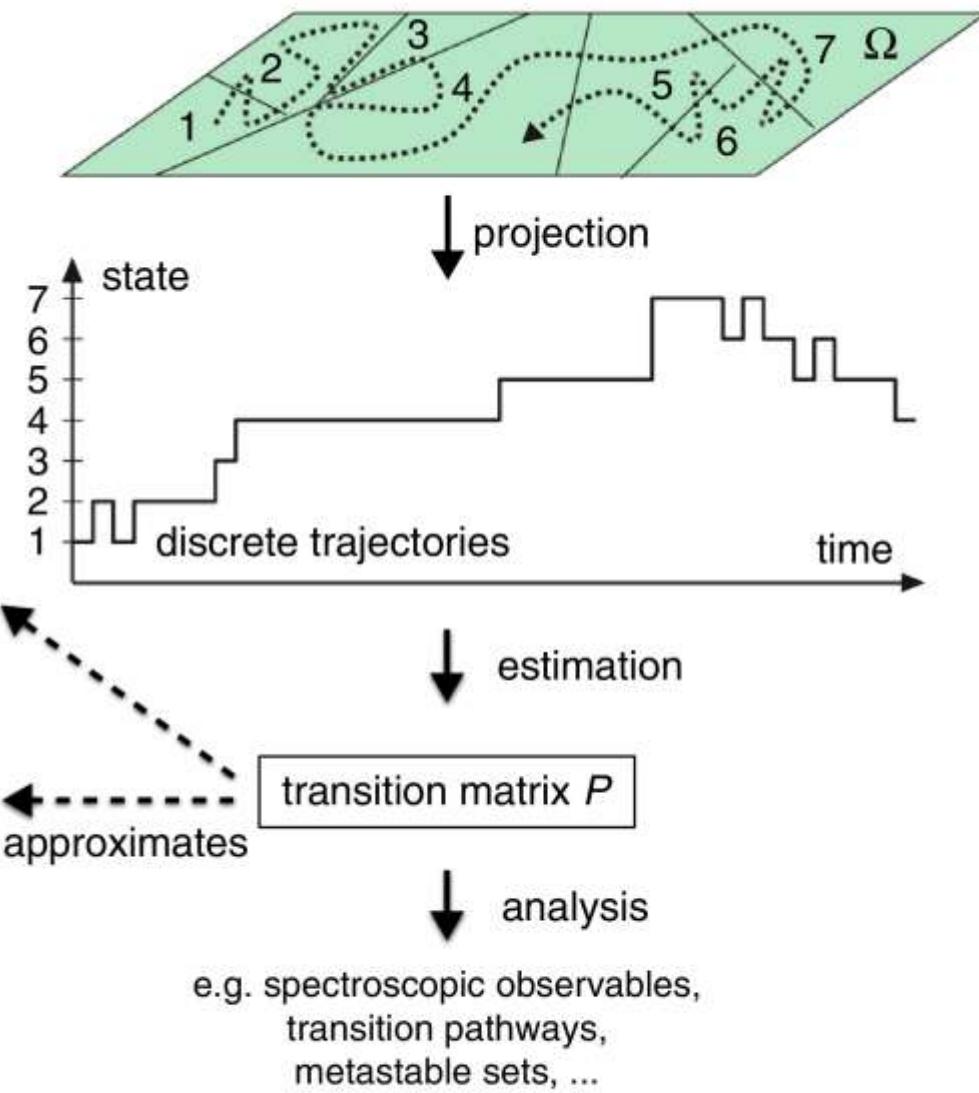


(c)



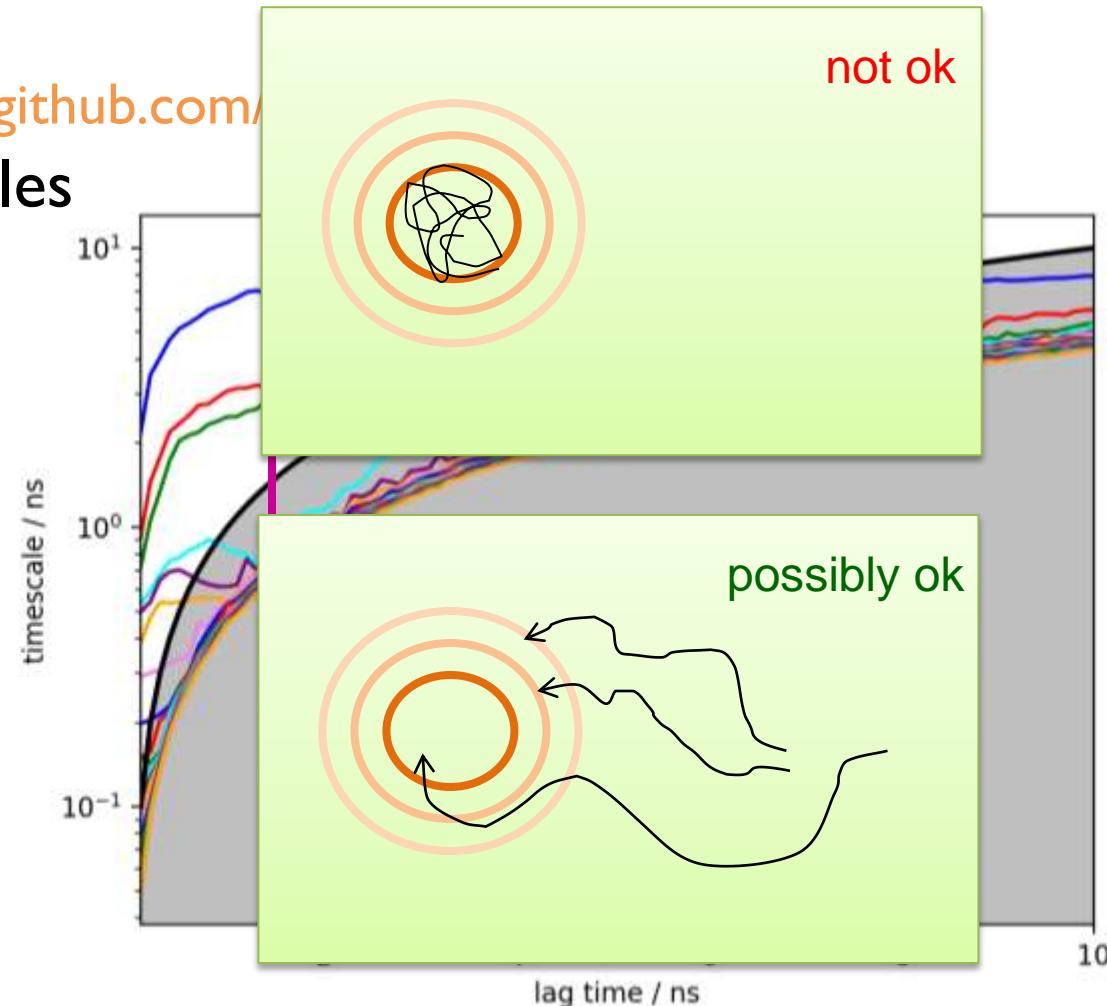
(d)

simulation trajectories



Even when merging multiple trajectories, making efficient use of sampling does not dispense us from sampling the phase space

- Ace-Ala3-Nme
- Part of examples at: github.com/omnisimulations/ace-alanine
- 6 Ramachandran angles
- The system is trapped
- How to «shoot» trajectories:
 - «bathtub» not ok
 - «shower» maybe
 - adaptive spawning
 - or your favourite string-like method



Examples from the literature

Complete reconstruction of an enzyme-inhibitor binding process by molecular dynamics simulations

Ignasi Buch, Toni Giorgino, and Gianni De Fabritiis¹

Computational Biochemistry and Biophysics Laboratory, Universitat Pompeu Fabra, Barcelona Biomedical Research Park, C/Doctor Aiguader 88, 08003 Barcelona, Spain

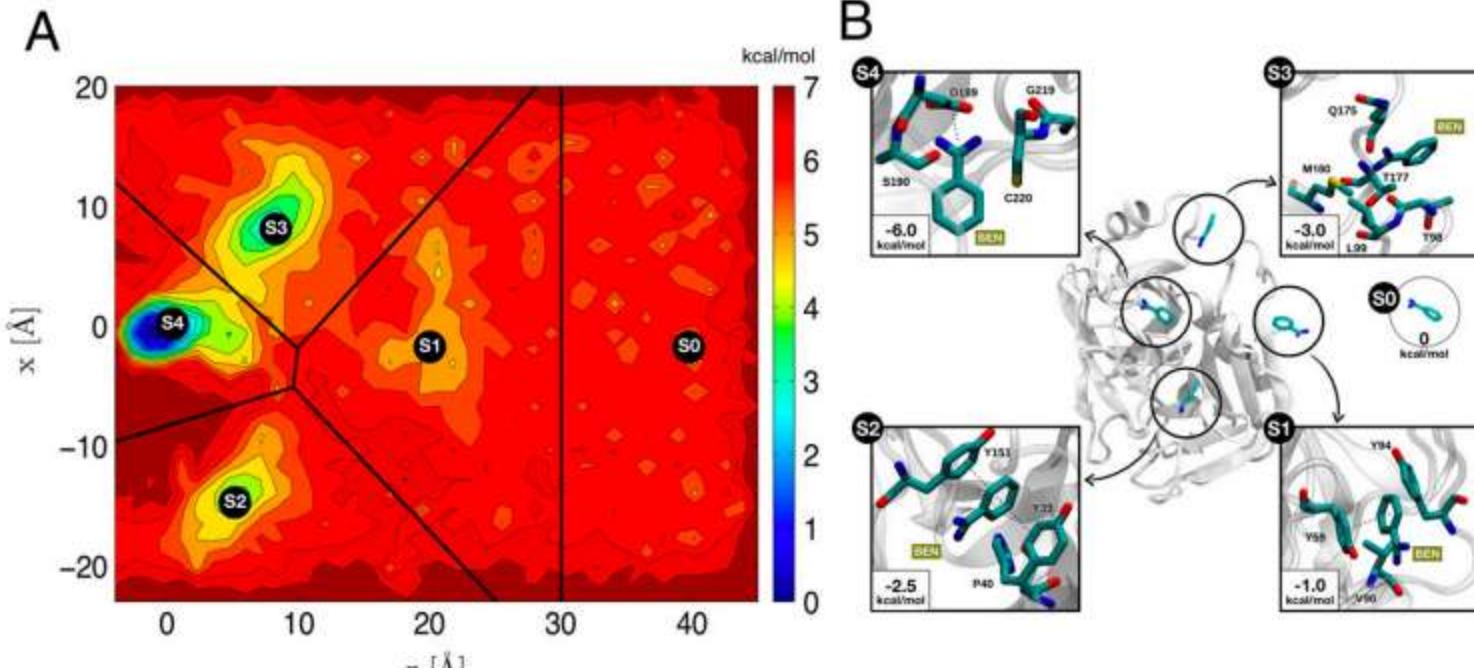
Edited by Arieh Warshel, University of Southern California, Los Angeles, CA, and approved May 11, 2011 (received for review March 4, 2011)

The understanding of protein–ligand binding is of critical importance for biomedical research, yet the process itself has been very difficult to study because of its intrinsically dynamic character. Here, we have been able to quantitatively reconstruct the complete binding process of the enzyme-inhibitor complex trypsin-benzamidine by performing 495 molecular dynamics simulations of the

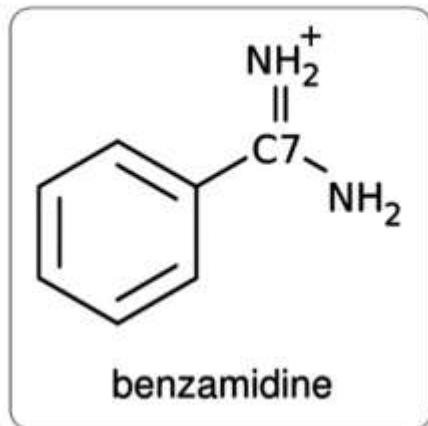
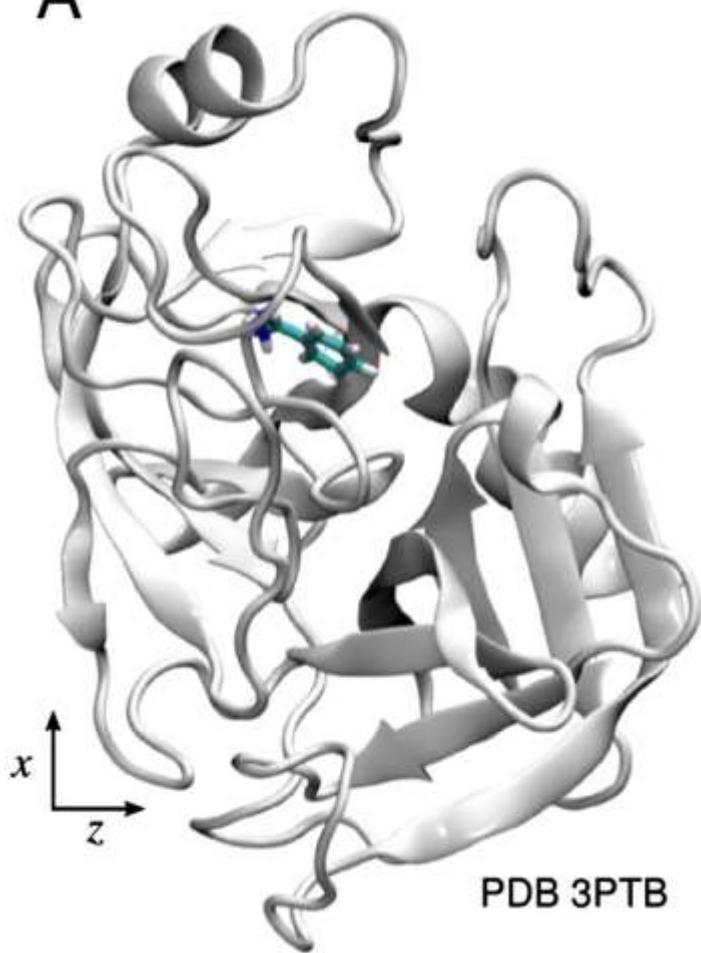
reproduce with atomic resolution the crystallographic mode of binding, but we also provide the kinetically and energetically meaningful transition states of the process.

Free ligand binding has been used in the past to describe computational experiments in which, typically, a ligand is placed at a certain distance from the target protein and first by diffusion and

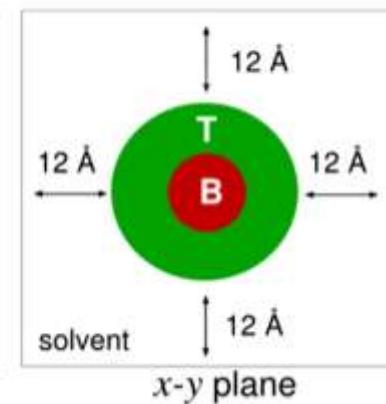
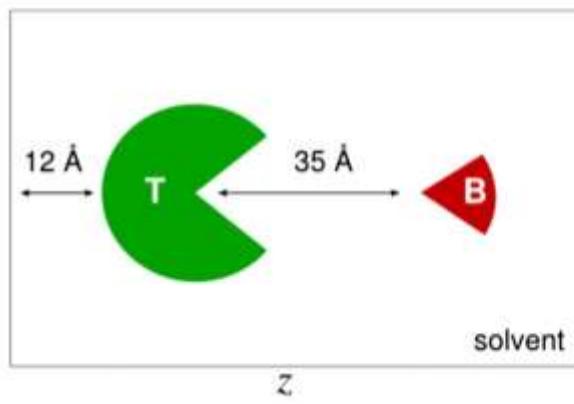
in the pro-

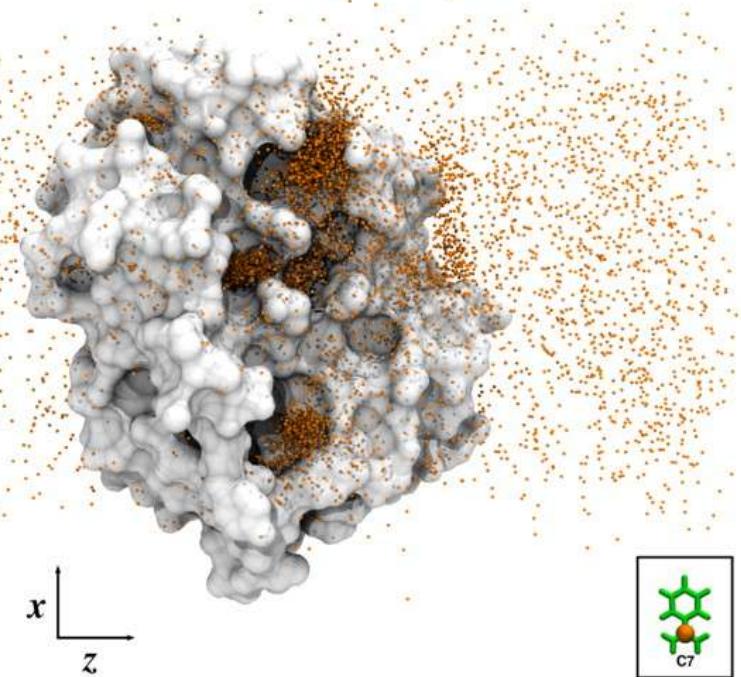
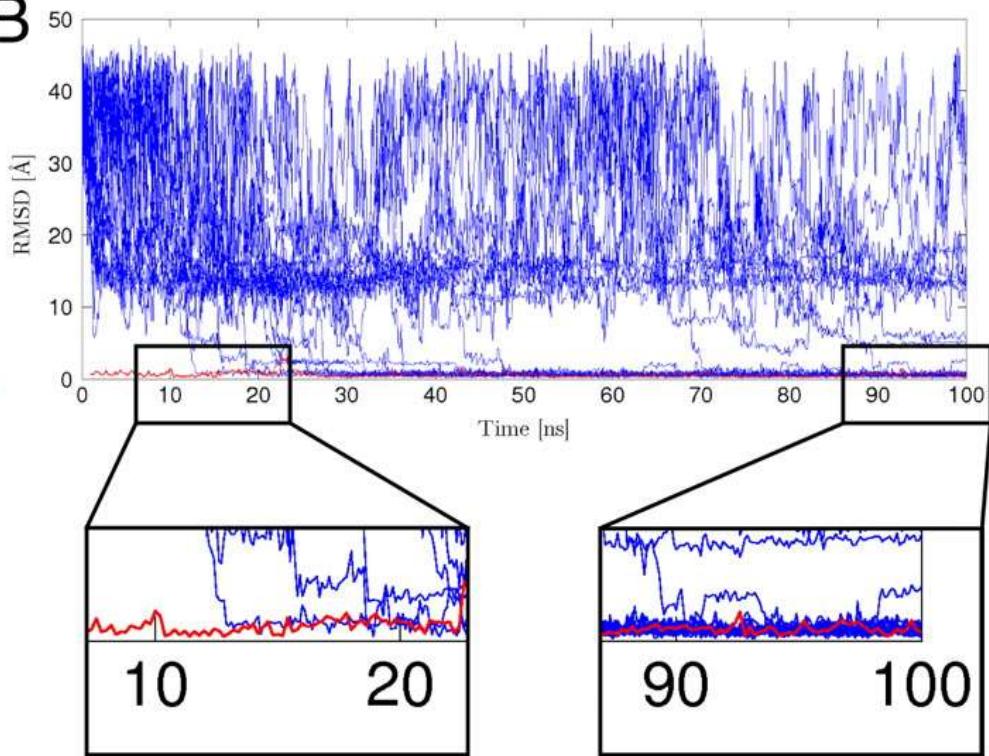


A

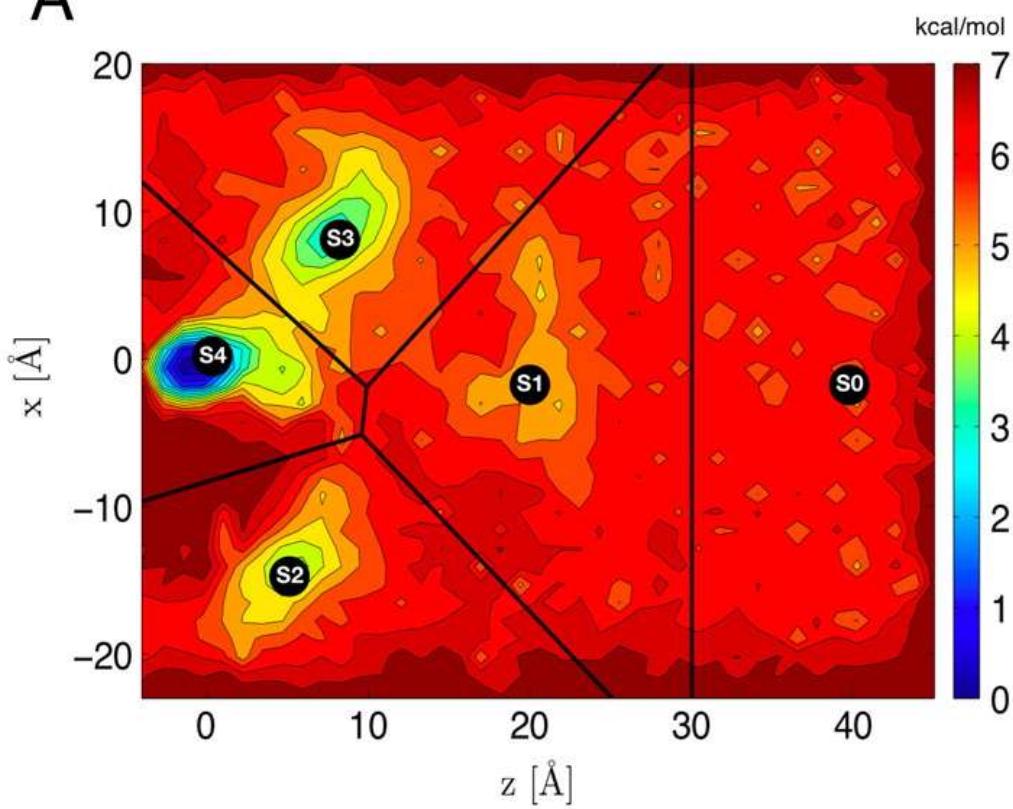


B

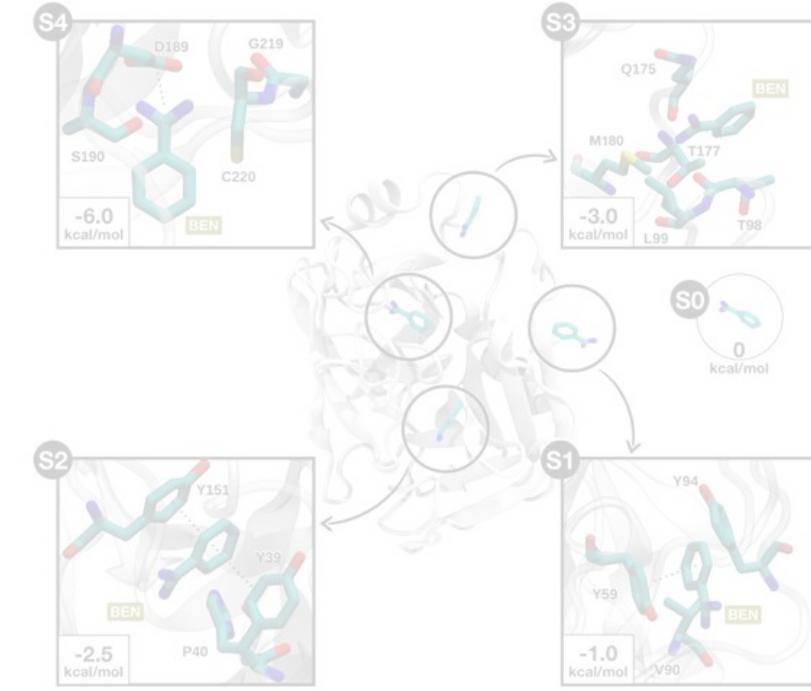


A**B**

A

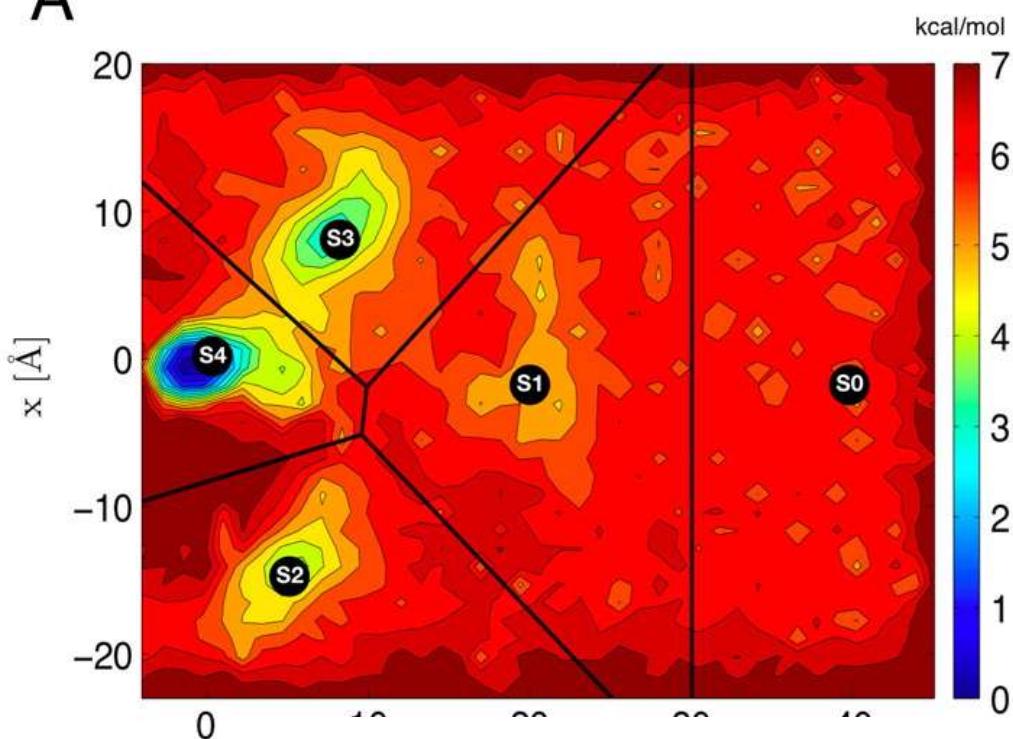


B



Identification of metastable states. (A) PMF in the xz plane. Five different metastable states can be identified from the different free-energy minima (S0 to S4). The relative free energy between the unbound state S0 and the bound state S4 is -6 kcal/mol. The most probable transition to the bound state S4 may be from S3 from the fact that the barrier between the two states is just 1.5 kcal/mol. (B) Structural characterization of metastable states. In states S1 and S2, benzamidine is stabilized by π - π stacking interactions with Y151 and Y39 side chains. In S3, a hydrogen bond may be formed between NH₂ groups of benzamidine (only heavy atoms shown for clarity) and Q175 side chain, or by a cation- π interaction between the Q175 side chain again, and benzamidine's benzene ring.

A



B

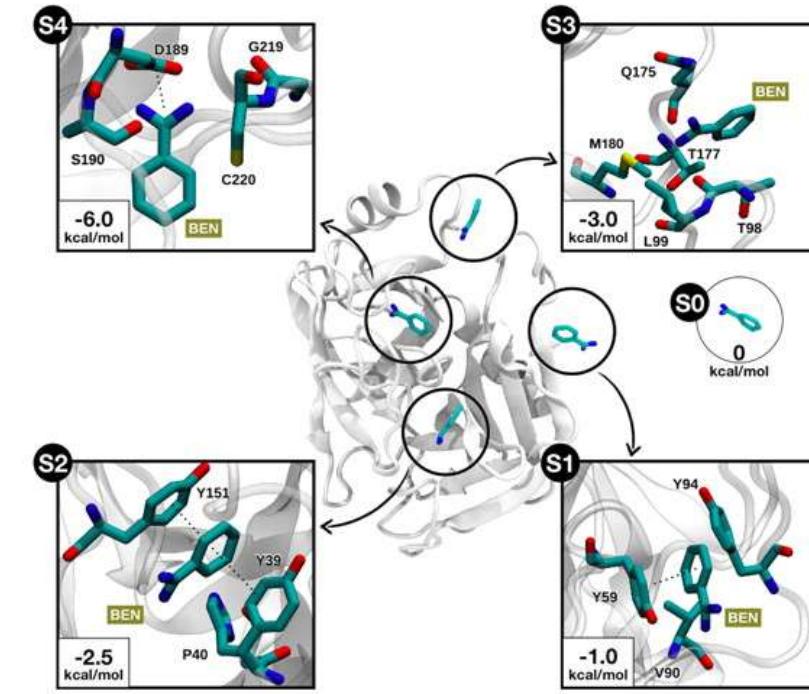


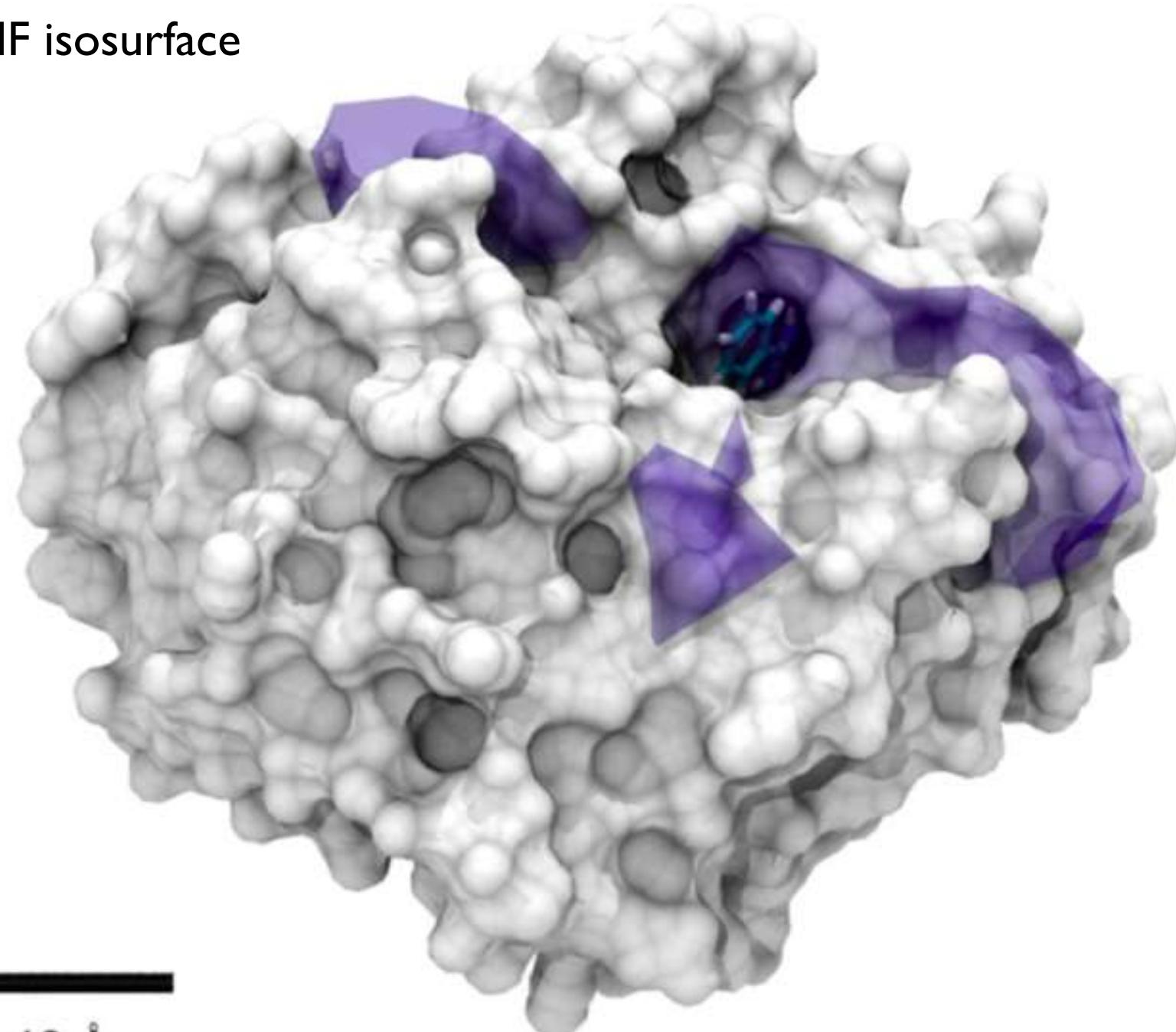
Table S1. Transition probabilities for the metastable states in the five-state coarse-grained model at a lag time of 50 ns

Identification of the identified from the state S0 and the from S3 from the characterization interactions with of benzamidine (between the Q17

	S0	S1	S2	S3	S4
S0	0.069	0.090	0.130	0.305	0.406
S1	0.066	0.094	0.124	0.300	0.416
S2	0.059	0.091	0.186	0.313	0.352
S3	0.073	0.096	0.103	0.361	0.366
S4	0.021	0.032	0.043	0.107	0.797

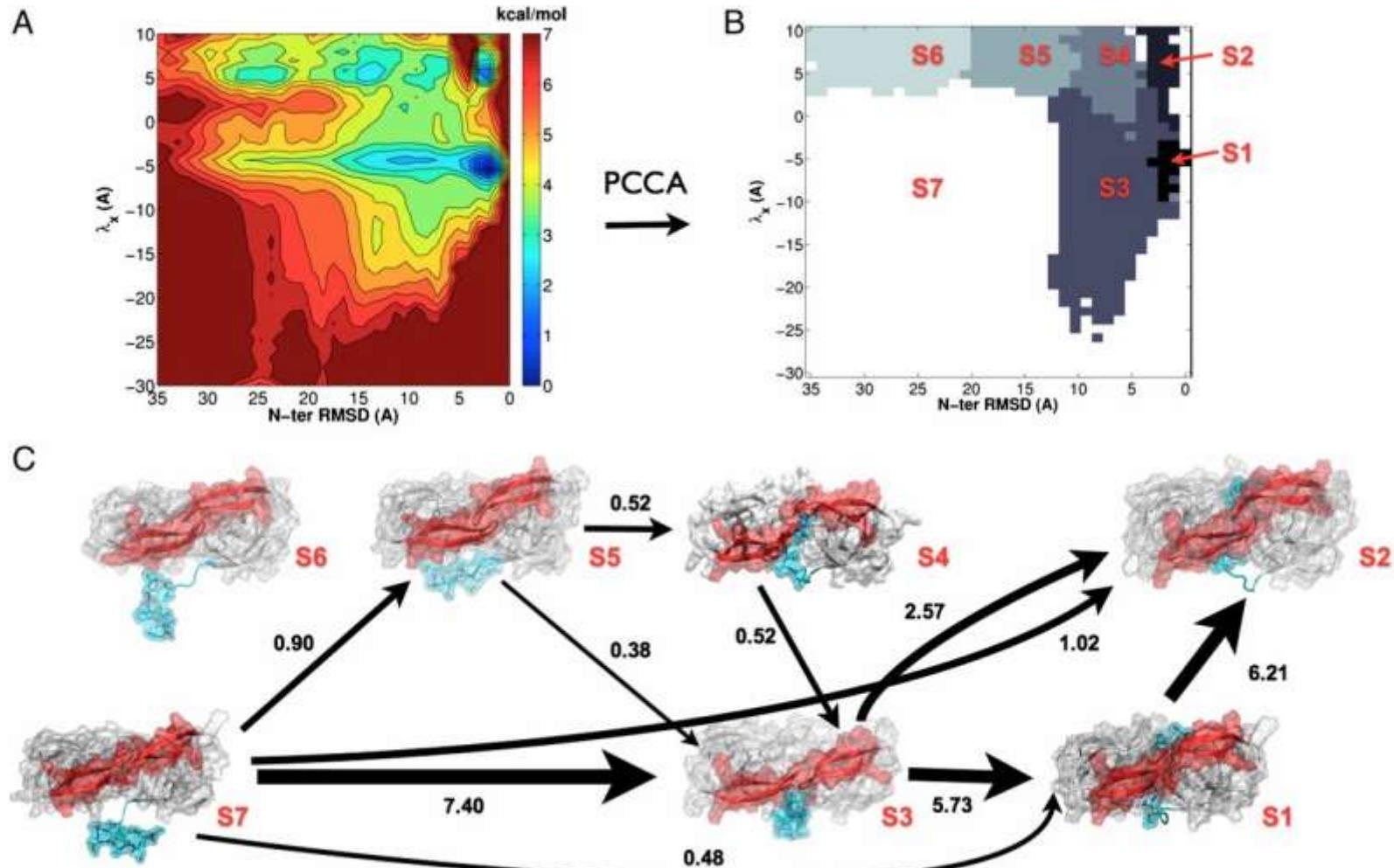
states can be seen in the unbound state S4 may be (B) Structural π-π stacking between NH₂ groups and π interaction

3D PMF isosurface



10 Å

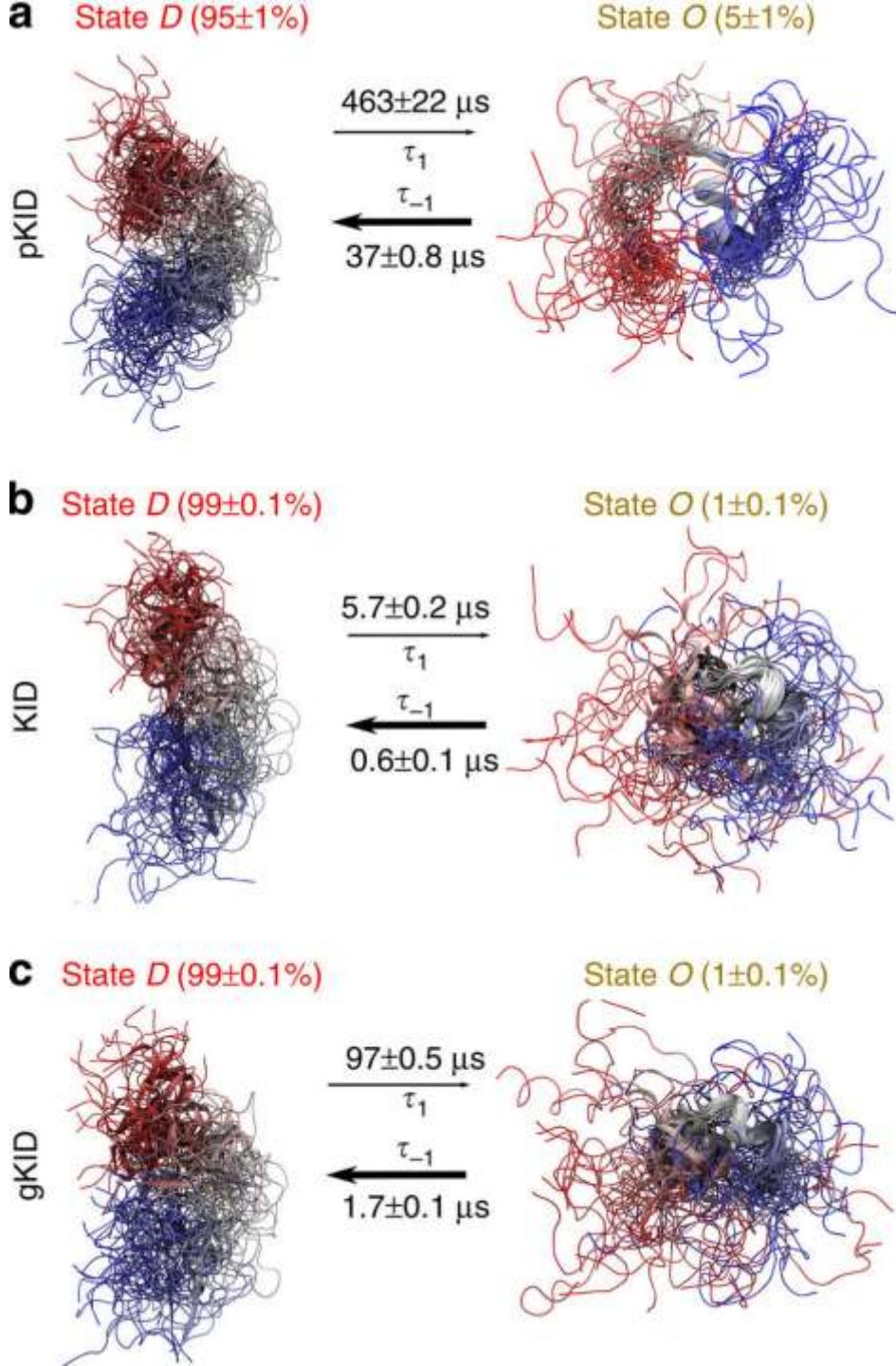
Kinetic characterization of the critical step in HIV-1 protease maturation



Kinetic modulation of a disordered protein domain by phosphorylation

N. Stanley, S. Esteban and G. De Fabritiis, Nat. Commun. 5, 5272 (2014)

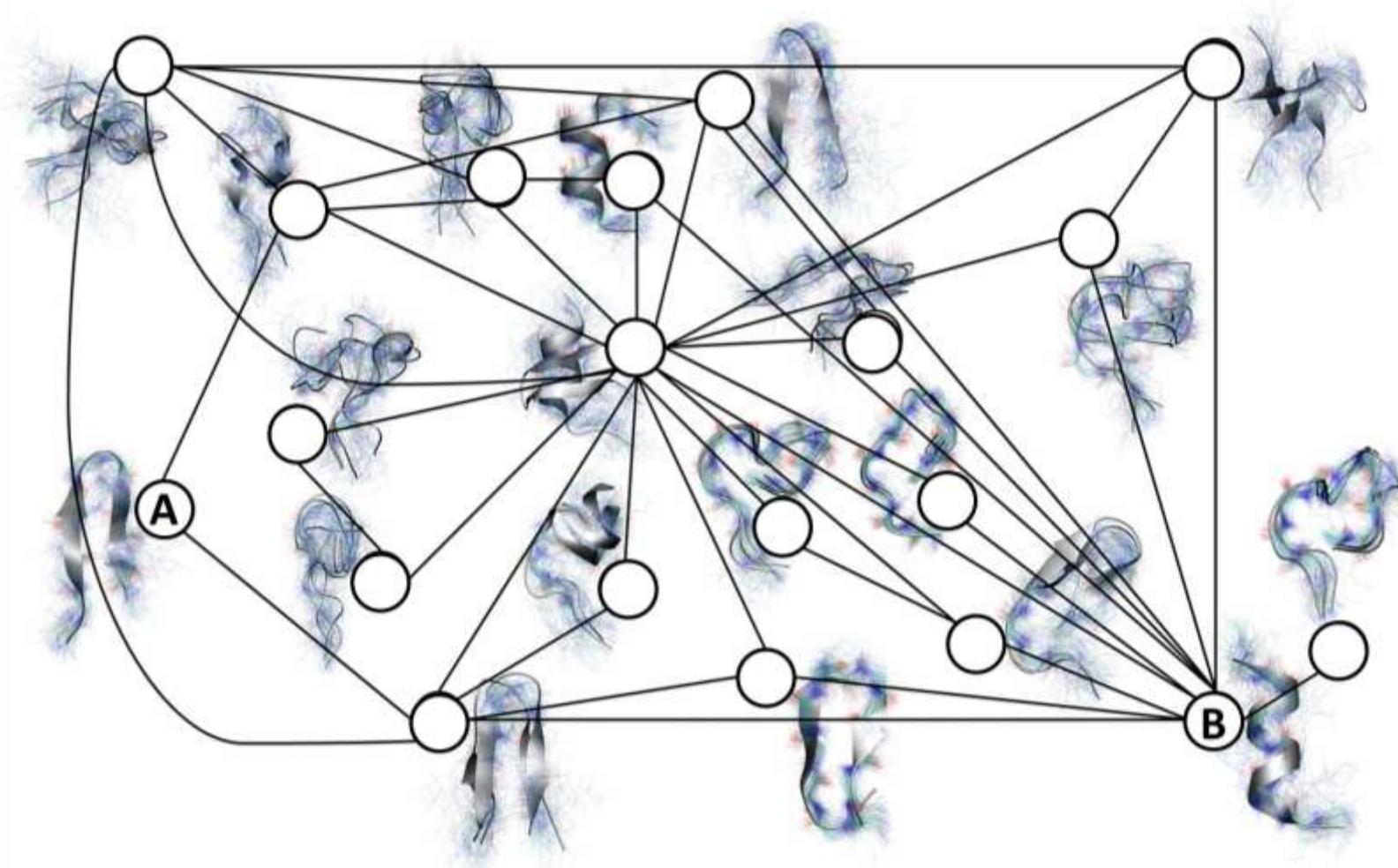
doi:10.1038/ncomms6272

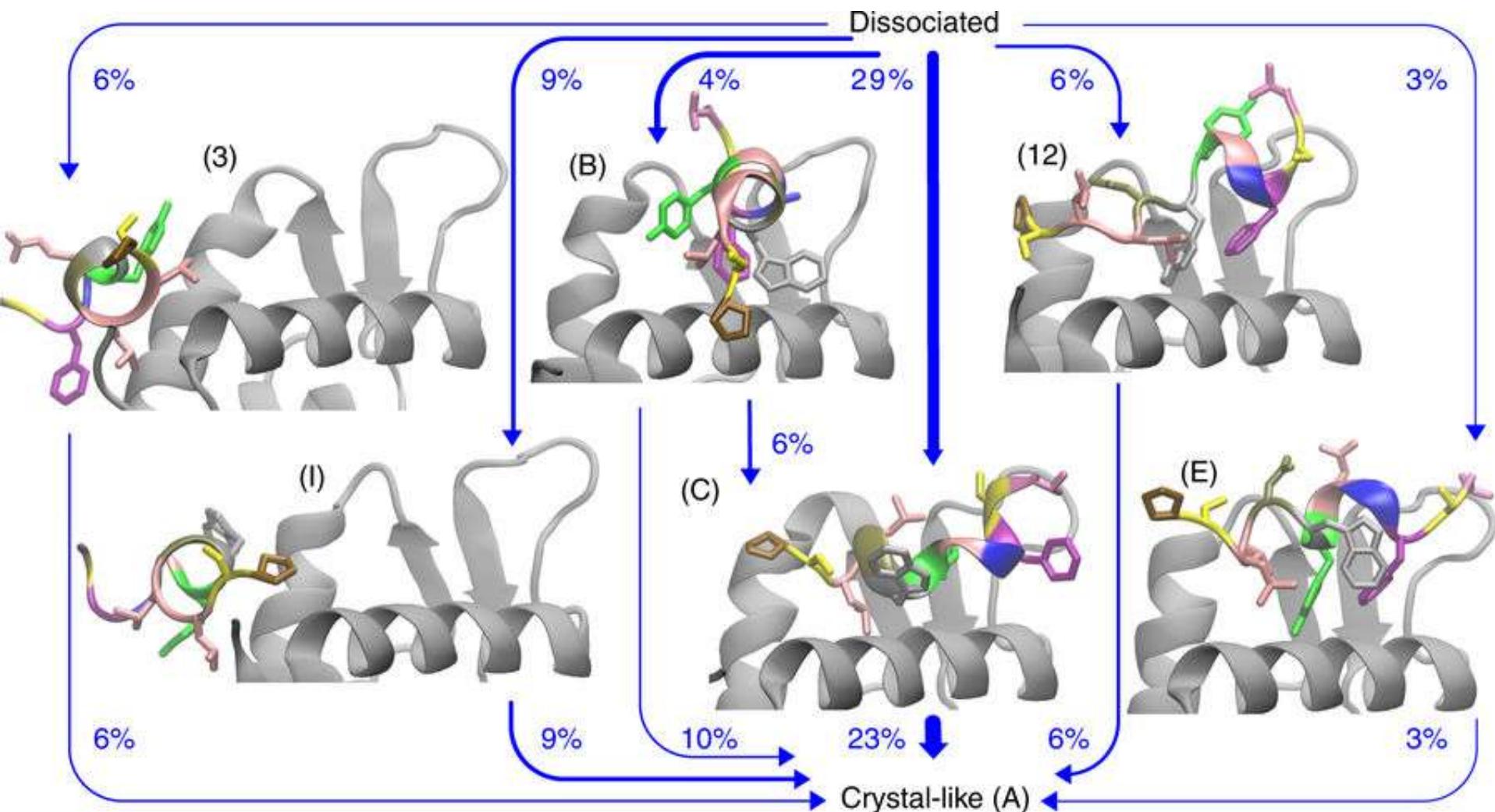


Folding kinetics

It is possible to model (small) protein folding processes at all-atom resolution.
However, the folding landscape is quite convoluted!

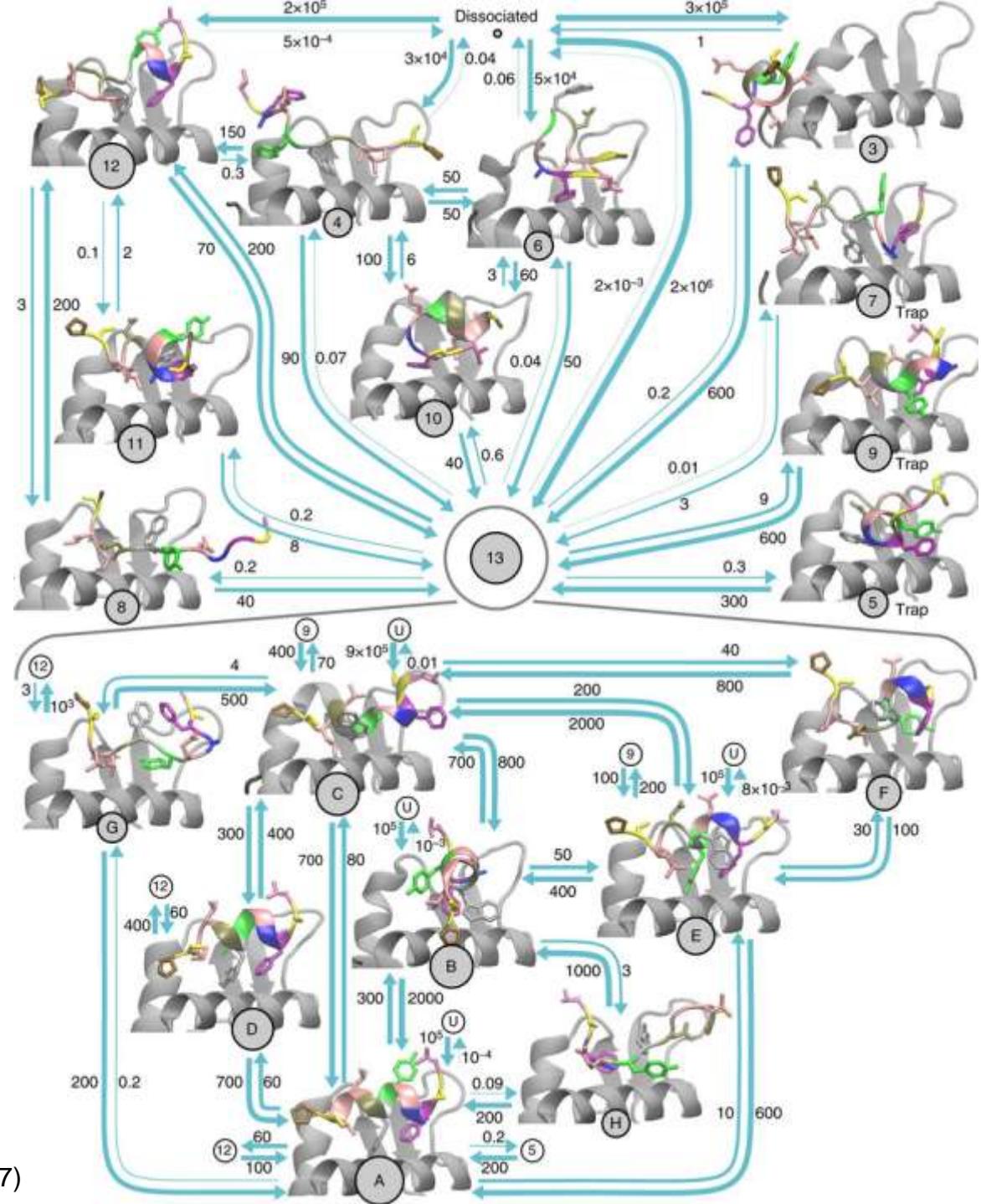
$$k_{on} = ?$$
$$k_{off} = ?$$





Binding mechanism comprised by the 60% most probable pathways. Structures of metastable (on-pathway) intermediates are shown, labels are as in Fig. 1. Arrows indicate the direction and relative magnitude of the reactive flux from the dissociated state to the crystal-like bound state. PMI residues that form PMI–Mdm2 contacts with at least a probability of 0.5 in a given macro-state are shown as sticks

Metastable states and transition rates for the binding of PMI to Mdm2. The PMI peptide is colored according to (). States are represented by discs with areas proportional to the natural logarithm of the equilibrium probability. Arrows indicate transitions with rate constants of at least 1 ms^{-1} in either direction. Numbers quantify transition rate constants in $\text{ms}^{-1} \text{ M}^{-1}$ for association events and in ms^{-1} for all other transitions. The definition of the states is hierarchical: between top-level states 0 and 13, transitions happen on timescales of $10 \mu\text{s}$ or slower. States in the lower part of the figure are sub-states of top-level state 13. There, PMI transitions between different states in the main binding pocket of Mdm2 on timescales of microseconds or slower (only states with large probabilities are shown)



Conclusions and Resources

Conclusions

- MSM methods may be an attractive formalism for medium-sized problems
- Make efficient use of *unbiased* sampling
- Still require *huge* (but not unreachable) amounts of sampling/simulation for biologically interesting systems
- Strong mathematical foundation, with good software available

Software + Tutorials

- All are Python-based
 - They include clear walkthroughs (highly recommended) with datasets
- PyEMMA – www.emma-project.org
- HTMD – www.htmd.org
 - Also analysis + system build + adaptive ...
 - Can aggregate large-scale datasets
- MSMBuilder - msmbuilder.org

Related topics

If you liked this, you may also like...

- Estimation of reversible transition matrices
- Hidden Markov Models doi:[10.1063/1.4828816](https://doi.org/10.1063/1.4828816)
- TICA
- Conformational fluxes
- Adaptive sampling
- Role of drug residence time in drug discovery

Literature

- Buch et al.,
10.1073/pnas.1103547108
 - Rigid ligand + protein association, simplest case
- Swinney, PMID:19152211
 - Importance of kinetics in drug design
- Chodera and Noe,
10.1016/j.sbi.2014.04.002
 - Excellent overview (1)
- Voelz et al., 10.1021/ja9090353
 - Full reconstruction of a millisecond folding
- Pande et al.,
10.1016/j.ymeth.2010.06.002
 - Excellent overview (2)
- Paul et al., 10.1038/s41467-017-0116-6
 - Peptide-peptide association
- Doerr et al.,
10.1021/ct400919u
 - Adaptive sampling
- Bryan et al., “The \$25 B eigenvector”
 - The original PageRank algorithm; e.g.
jdc.math.uwo.ca/M1600b/I/page_rank-1600.pdf

End