



THE UNIVERSITY
of EDINBURGH

Time Series Analysis with Markov State Models

Antonia Mey

School of Chemistry
University of Edinburgh
antonia.mey@ed.ac.uk



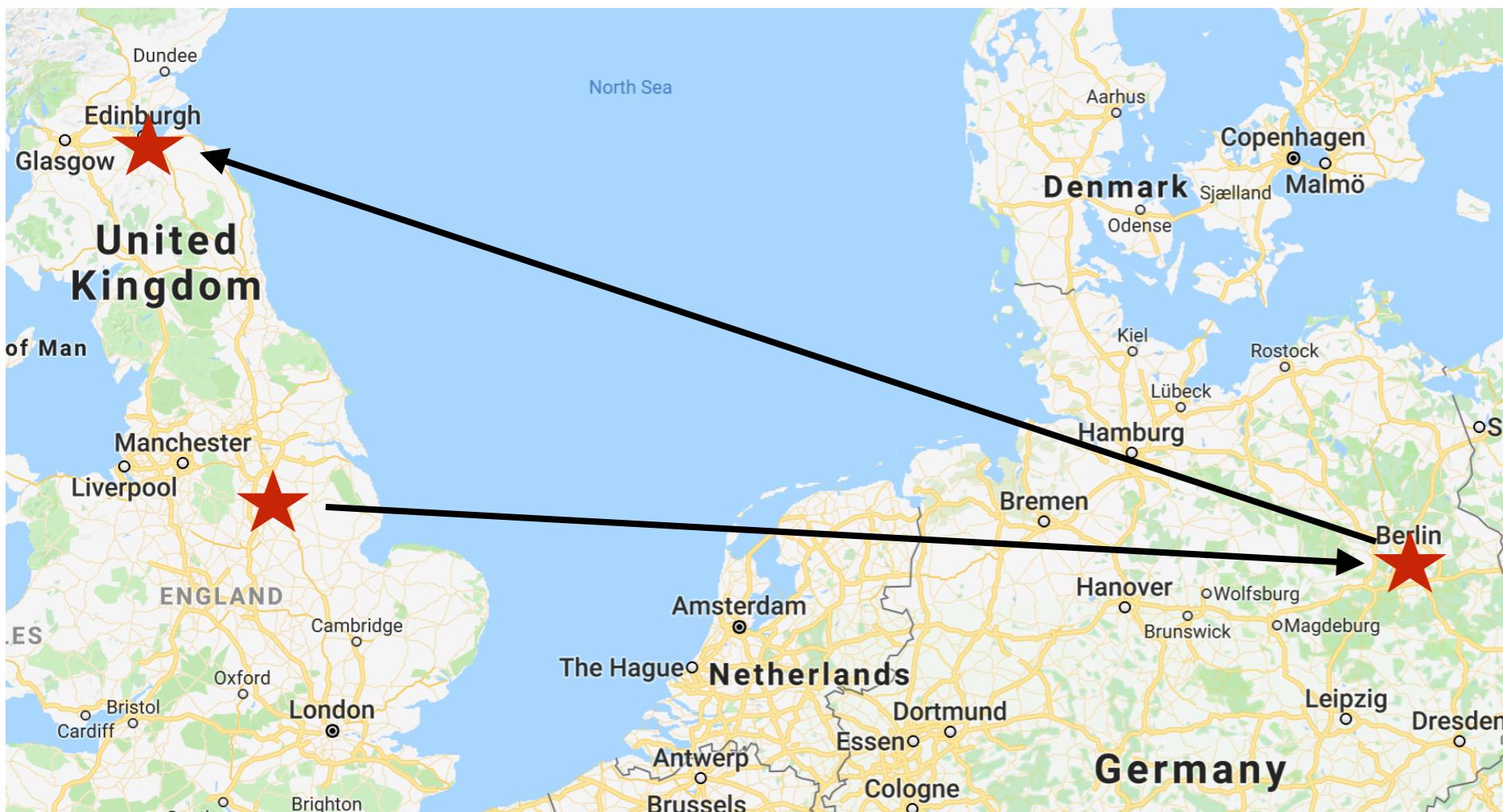
ppxasjsm
@ppxasjsm

PyData Meetup

06/09/2018

About me

- 2013 Ph.D. in Physics at the University of Nottingham
- 2013-2015 Researcher at Freie Universität Berlin (Department of Mathematics)
- 2015-now Researcher at University of Edinburgh (School of Chemistry)

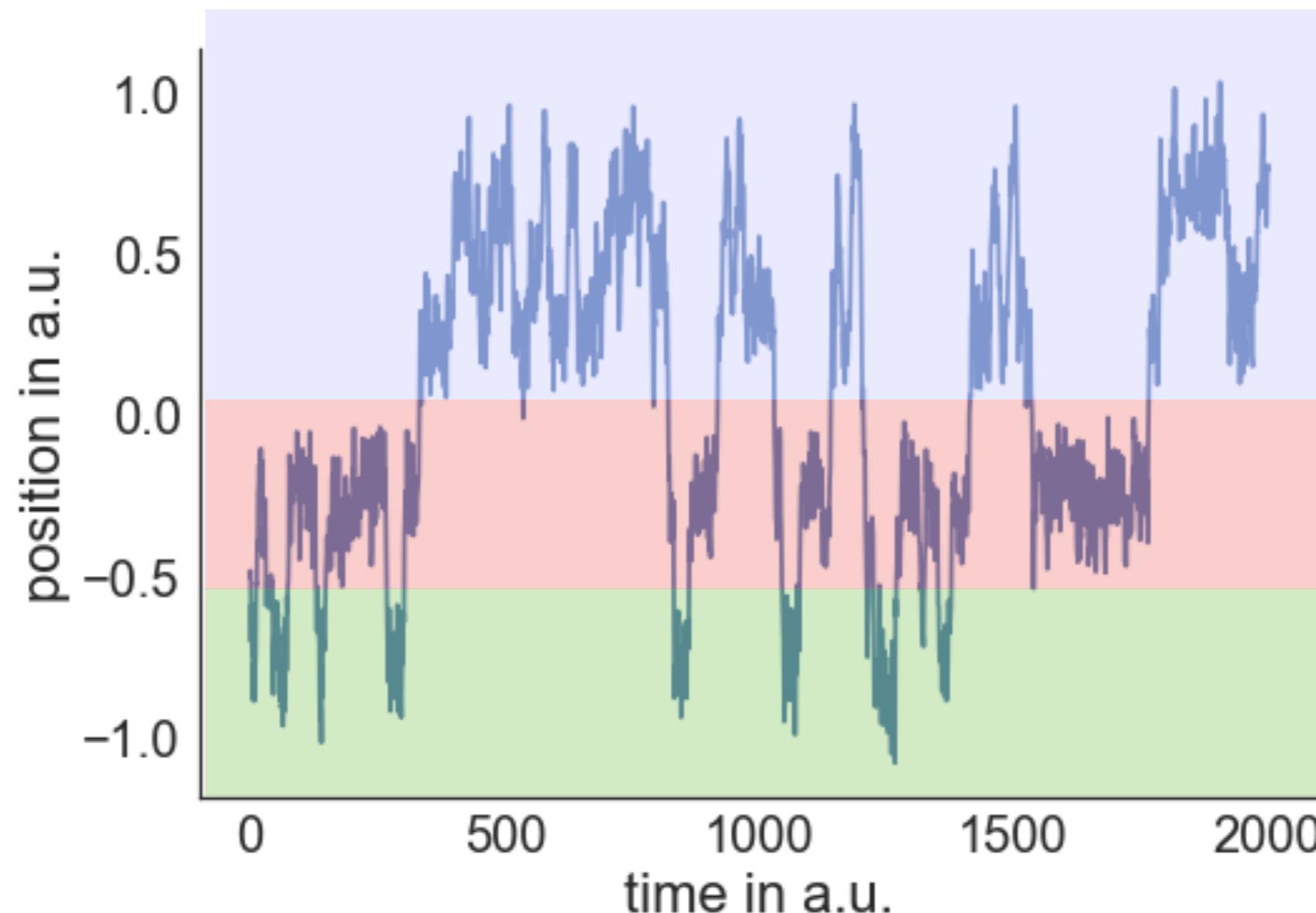


What am I?

A computational —
insert science
related word



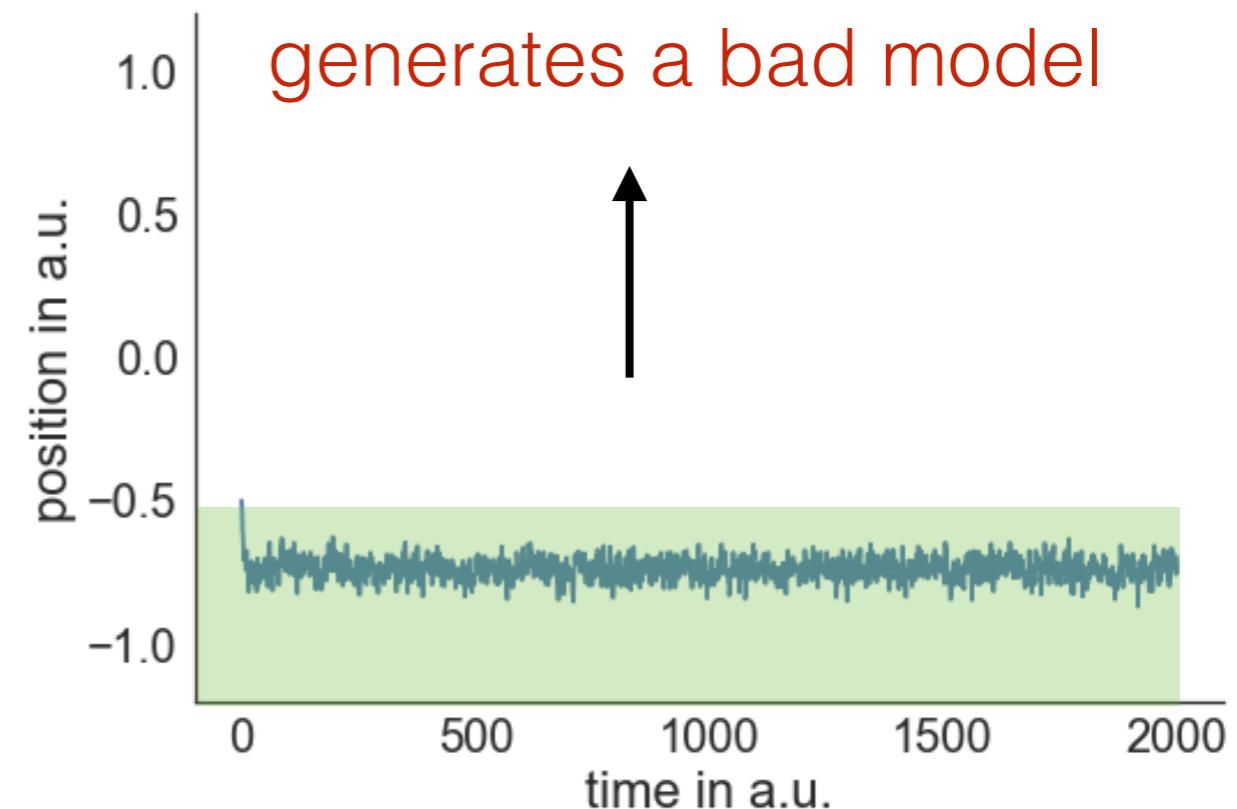
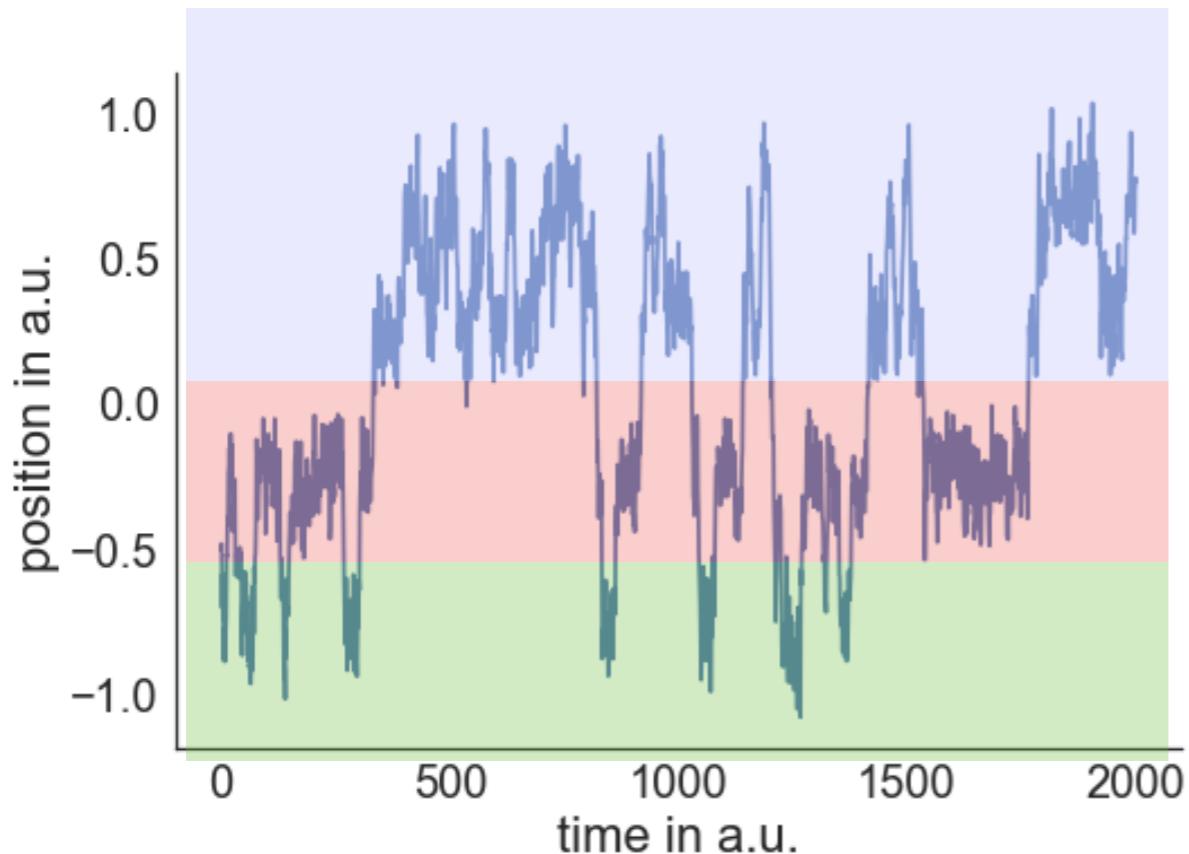
Timeseries, what now?



Suppose we have a timeseries and want to identify patterns and then use a model to make future predictions about this timeseries.

One option: model based analysis using a **Markov Model, or Hidden Markov Model**

Why a Markov Model approach?



- Estimate the residence time in metastable states
- Mean first passage time between state
- Transition path analysis allows to evaluate dominant paths in the system
- Many short timeseries can be used in an aggregate way avoiding bad models in systems with rare events

How to Markov Model

The basic assumption is that the timeseries can be represented by a Markov Jump process, meaning that the Markov property holds:

$$P(S_i \rightarrow S_j) = P[X_{t+1} \in S_i | X_t \in S_j]$$

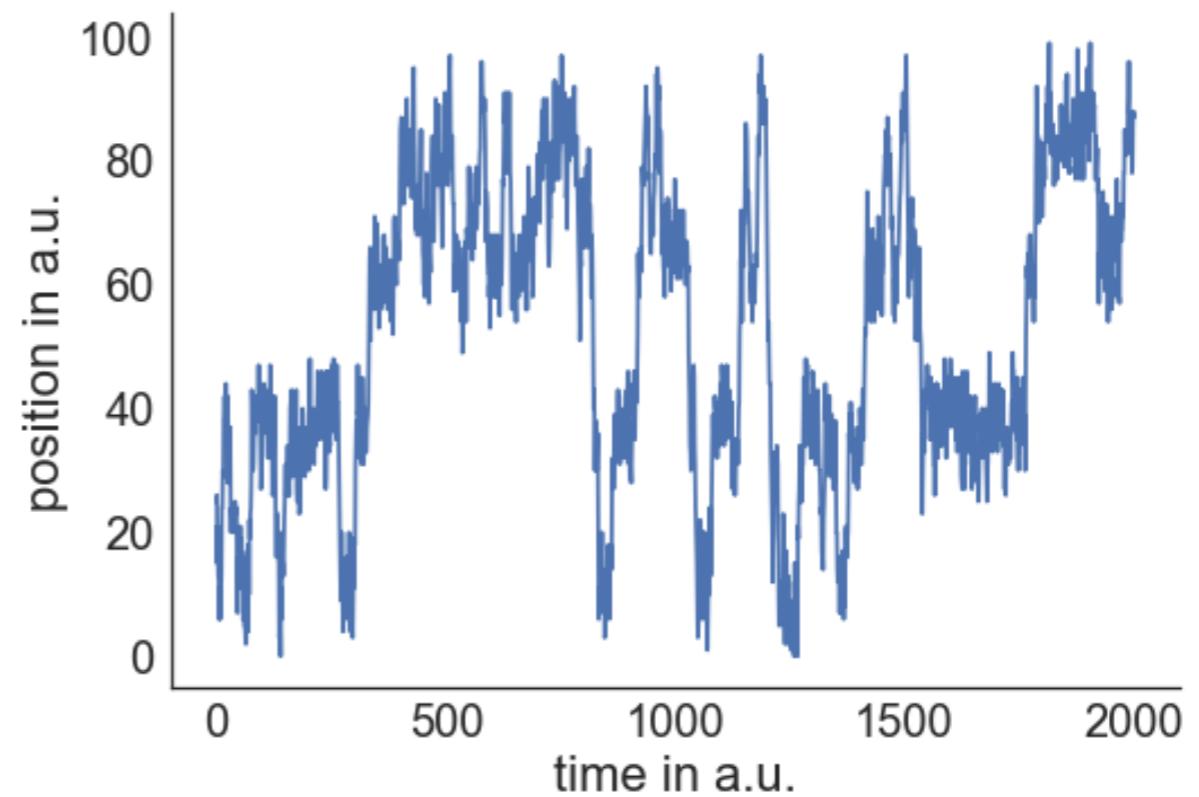
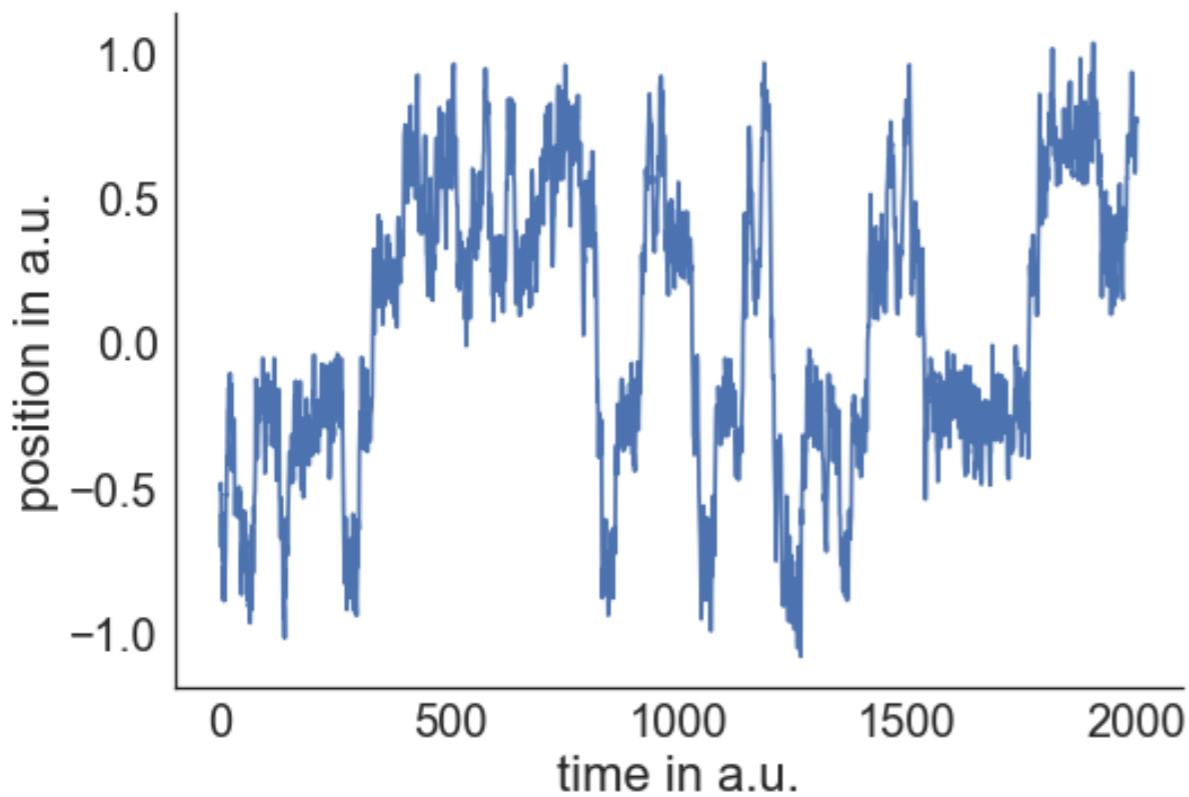
This leads to the following steps of building a Markov model from a timeseries.

1. Generate timesries/ acquire timeseries
2. Complexity reduction (optional if the system has many dimensions)
3. Discretisation
4. Transition matrix estimation
5. Transition matrix analysis

3. Discretisation

Discretisation

For illustrative purposes: Assume we discretise the position into 100 discrete states



4. Transition matrix estimation

Transition matrix estimation

$$\begin{matrix} \frac{c_{ii}}{\sum_i c_{ij}} & \frac{c_{ij}}{\sum_i c_{ij}} \\ \hline & = \mathbf{T} \end{matrix}$$

The transition matrix contains conditional probabilities, of going from state i to state j. Usually, a reversible estimation is used and detailed balance set as a constraint.

$$\pi_i T_{ij} = \pi_j T_{ji}$$

Reversible estimation

Objective: find the most likely **reversible transition** matrix, based on the **observed counts** using Bayes $\mathbb{P}(\mathbf{T}|\mathbf{C}) = \prod_{i,j} t_{ij}^{c_{ij}}$

We use log-likelihoods instead: $Q = \log \mathbb{P}(\mathbf{T}|\mathbf{C}) = \sum_{i,j} c_{ij} \log t_{ij}$

Maximise the log-likelihood, by taking its derivative and using the constraint, that detailed balance must hold, i.e. $\frac{\partial Q}{\partial x_{ij}} = 0$ and a variable transform $t_{ij} = \frac{x_{ij}}{\sum_k^n x_{ik}}$

$$\frac{\partial Q}{\partial x_{ji}} = \frac{c_{ij} + c_{ji}}{x_{ji}} - \frac{c_i}{x_i} - \frac{c_j}{x_j}$$

set to 0 and the entries can be iterated to convergence.

$$x_{ji} = \frac{c_{ij} + c_{ji}}{\frac{c_i}{x_i} + \frac{c_j}{x_j}} \longrightarrow \text{self consistent iterative update of } x_{ij}$$

$$x_i = \sum_k x_{ki} \quad \text{which is incidentally also the stationary probability of state i.}$$

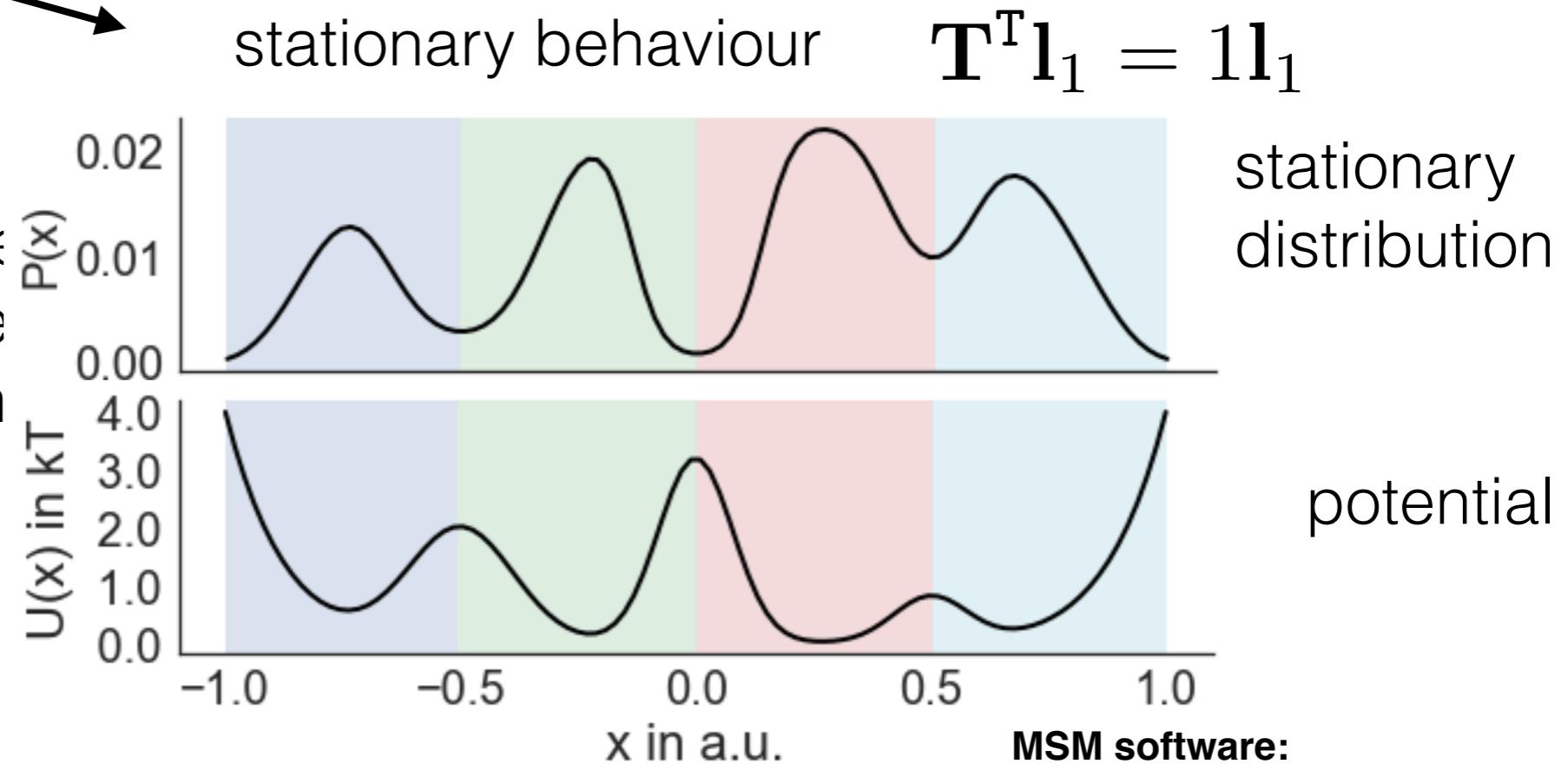
Markov Model Analysis I

Transition matrix analysis - stationary properties

For all transition matrices holds:

$$\lambda_1 = 1 > \lambda_2 > \lambda_3, \dots, > \lambda_n \longrightarrow \text{dynamics}$$

Also: *MFPT*, can be computed from the transition matrix and the stationary distribution.



MSM software:

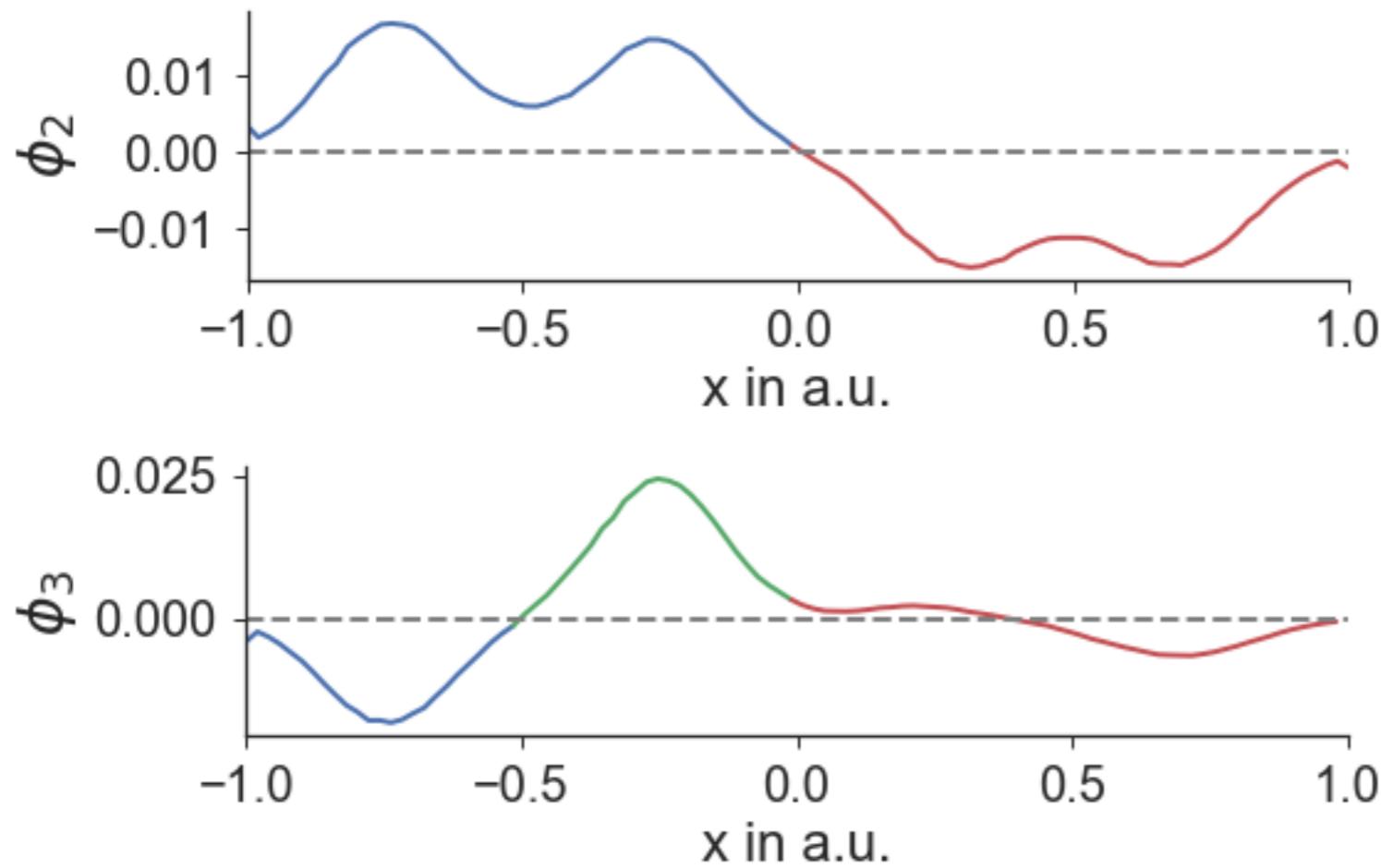
Senne, Trendelkamp-Schroer, Mey,
Schütte, Noé, JCTC, 8, 2223 (2012)
<http://simtk.org/home/emma>

Markov Model Analysis II

Transition matrix analysis - dynamic properties

For all transition matrices holds:

$$\lambda_1 = 1 > \boxed{\lambda_2 > \lambda_3, \dots, > \lambda_n} \longrightarrow \text{dynamics}$$

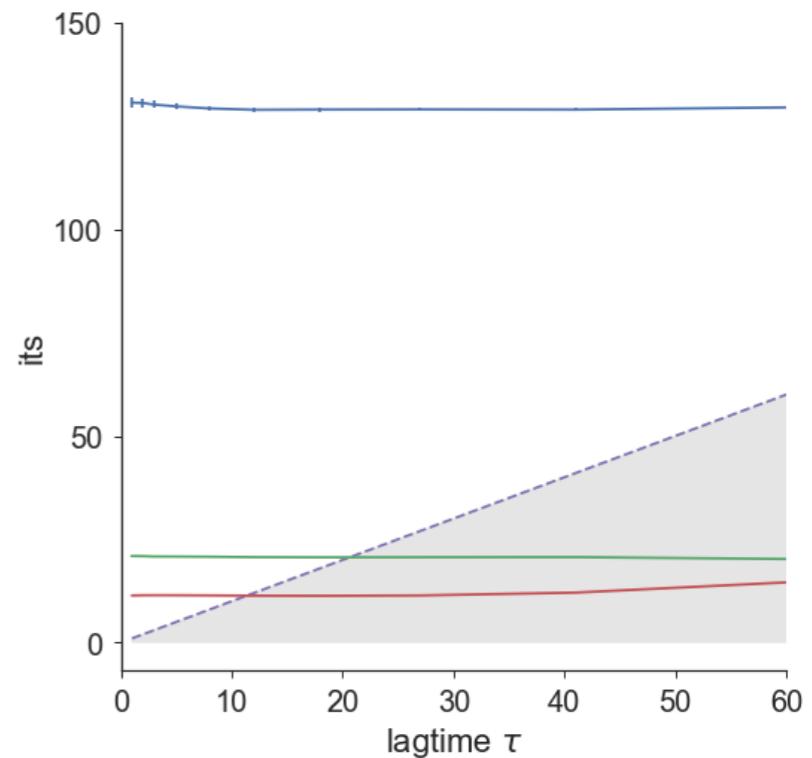


$$\mathbf{T}\mathbf{r}_i = \lambda_i \mathbf{r}_i$$

$$\mathbf{T}^\top \mathbf{l}_i = \lambda_i \mathbf{l}_i$$

$$t_i = -\frac{\tau}{\ln \lambda_i}$$

$$k = t_i^{-1}$$



MSM software:

Senne, Trendelkamp-Schroer, Mey,
Schütte, Noé, JCTC, 8, 2223 (2012)

<http://simtk.org/home/emma>

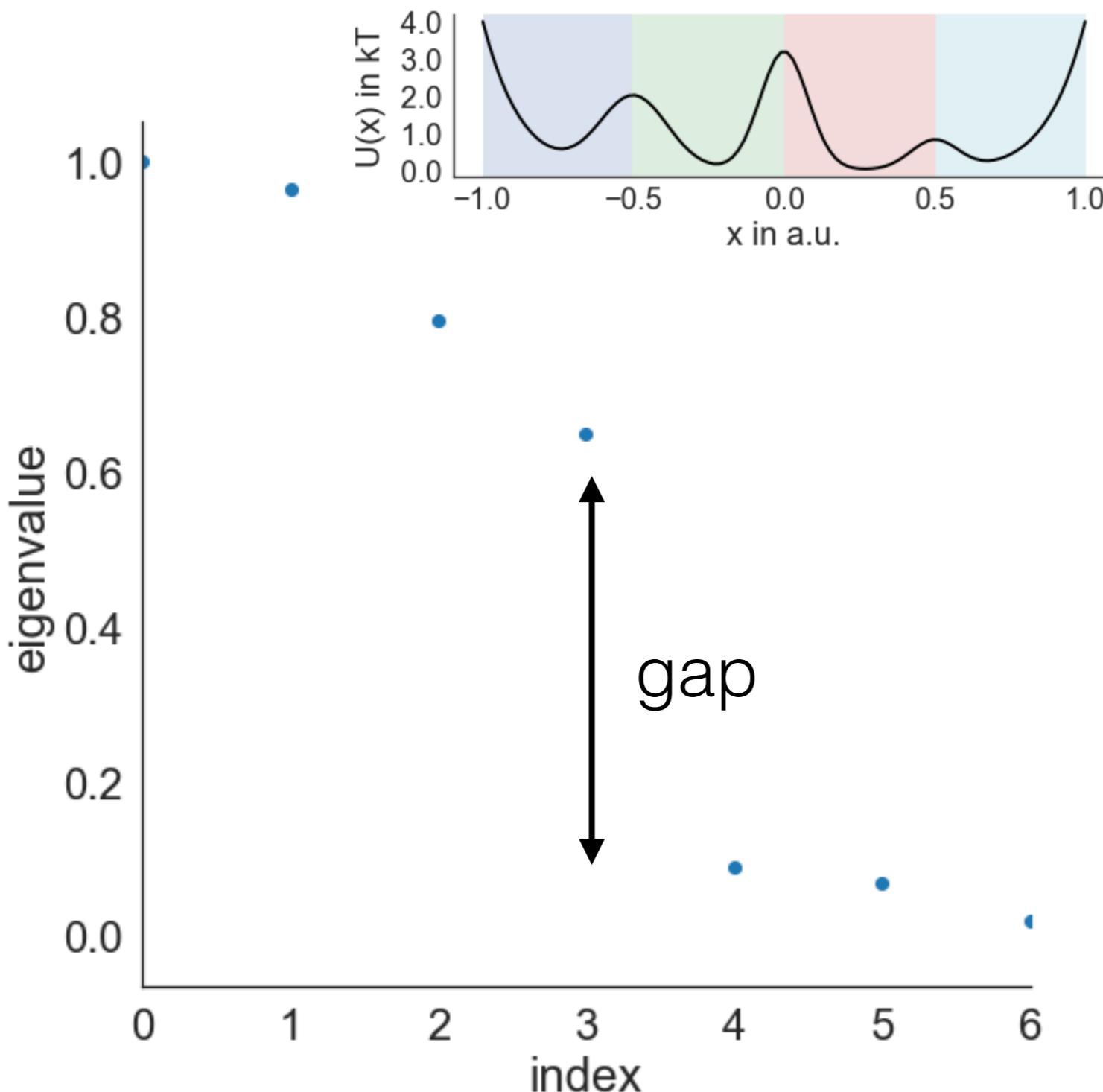
<https://github/markovmodel/pyemma>



PyEMMA demo

Go to Jupyter notebook

Coarse graining



System can be described in terms of slow dynamics between 4 wells.

HMM, PCCA or spectral clustering can be used to define coarse states.

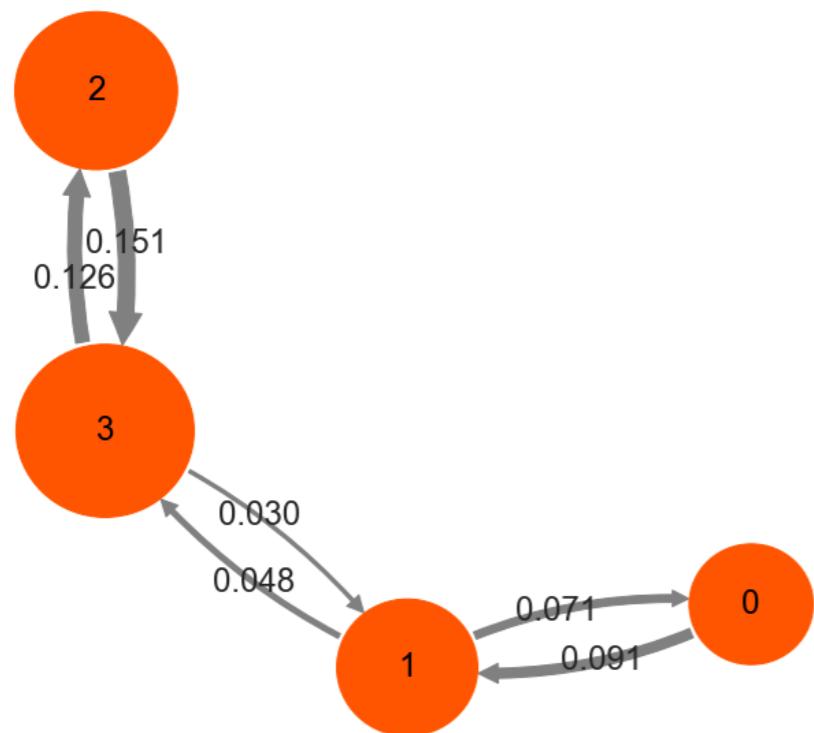
Mean first passage time

$$Z = (I - P + \Pi)^{-1}$$

$$m_{ij} = \frac{z_{jj} - z_{ij}}{\pi_j}$$

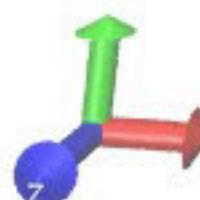
m_{ij} is the mean first passage time of state i to state j .

The inverse of the mean first passage time that is often experimentally measured in biological processes.

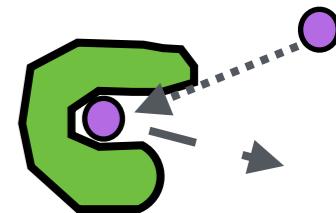


Dominant pathways between coarse grained metastable states can be found using transition path theory, allowing to compute net flux between states.

Molecular Dynamics



1. Timescale problem:



Molecular dynamics simulations are relying on integration time steps of 1-4 fs, for accuracy. $< 10^8$ integration time steps are needed to reach relevant timescales.

2. Complexity problem:

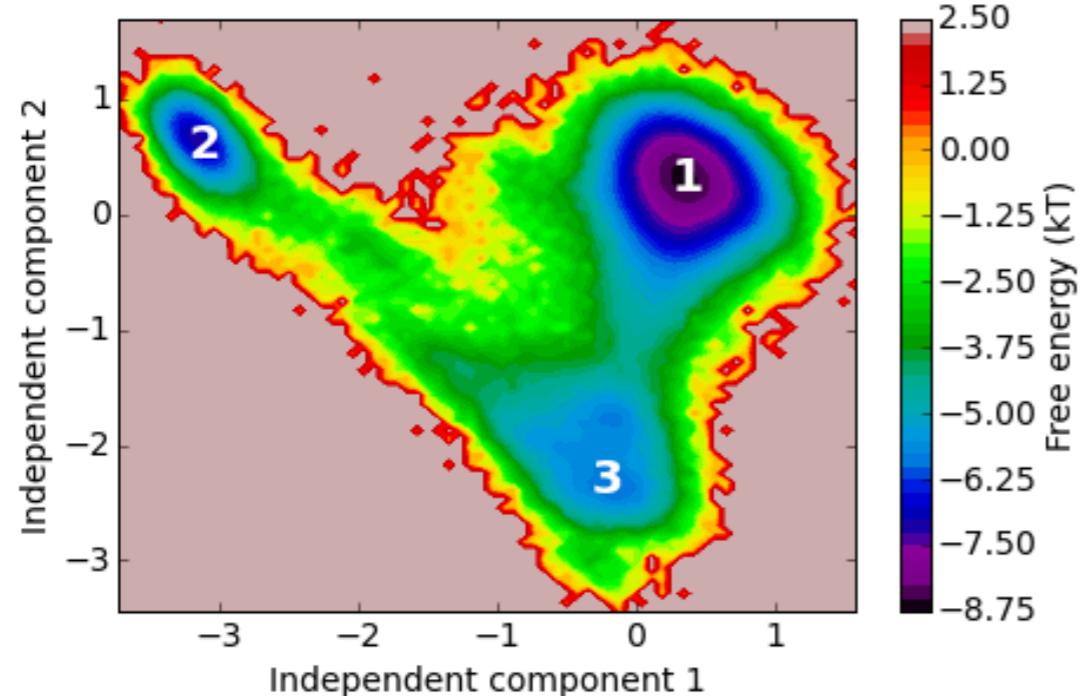
Reduction of $3N$ coordinates in order to make data more manageable/quantitative:

- ▶ clustering of relevant states
- ▶ Minimize information loss with the clustering

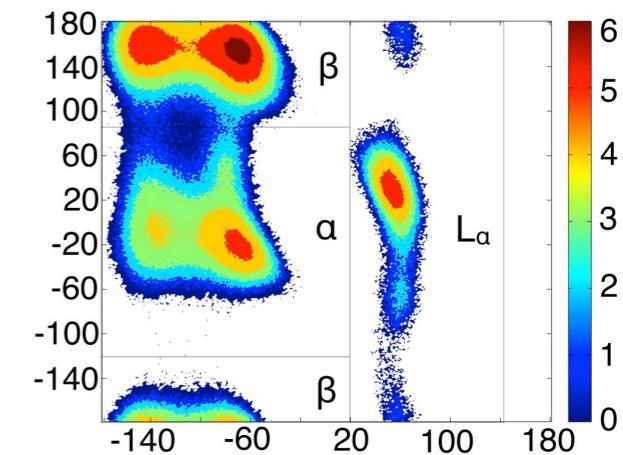
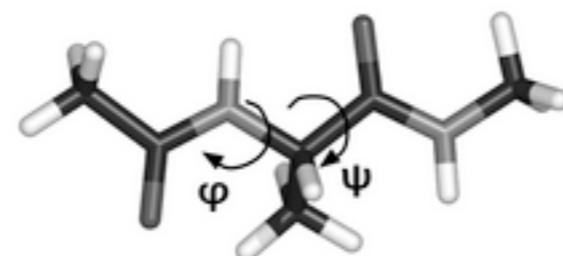
Dimensionality reduction with TICA

Complexity reduction and feature selection

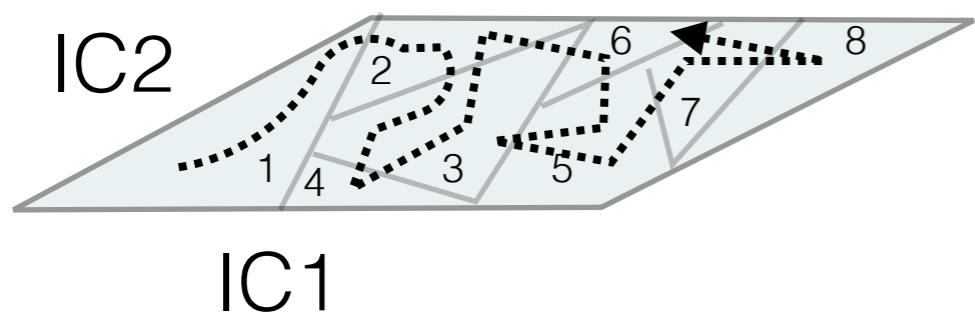
No point in clustering on $3N$ dimensions, there just aren't enough data points



Dimensionality reduction using TICA
(time lagged independent component analysis) or direct feature selection

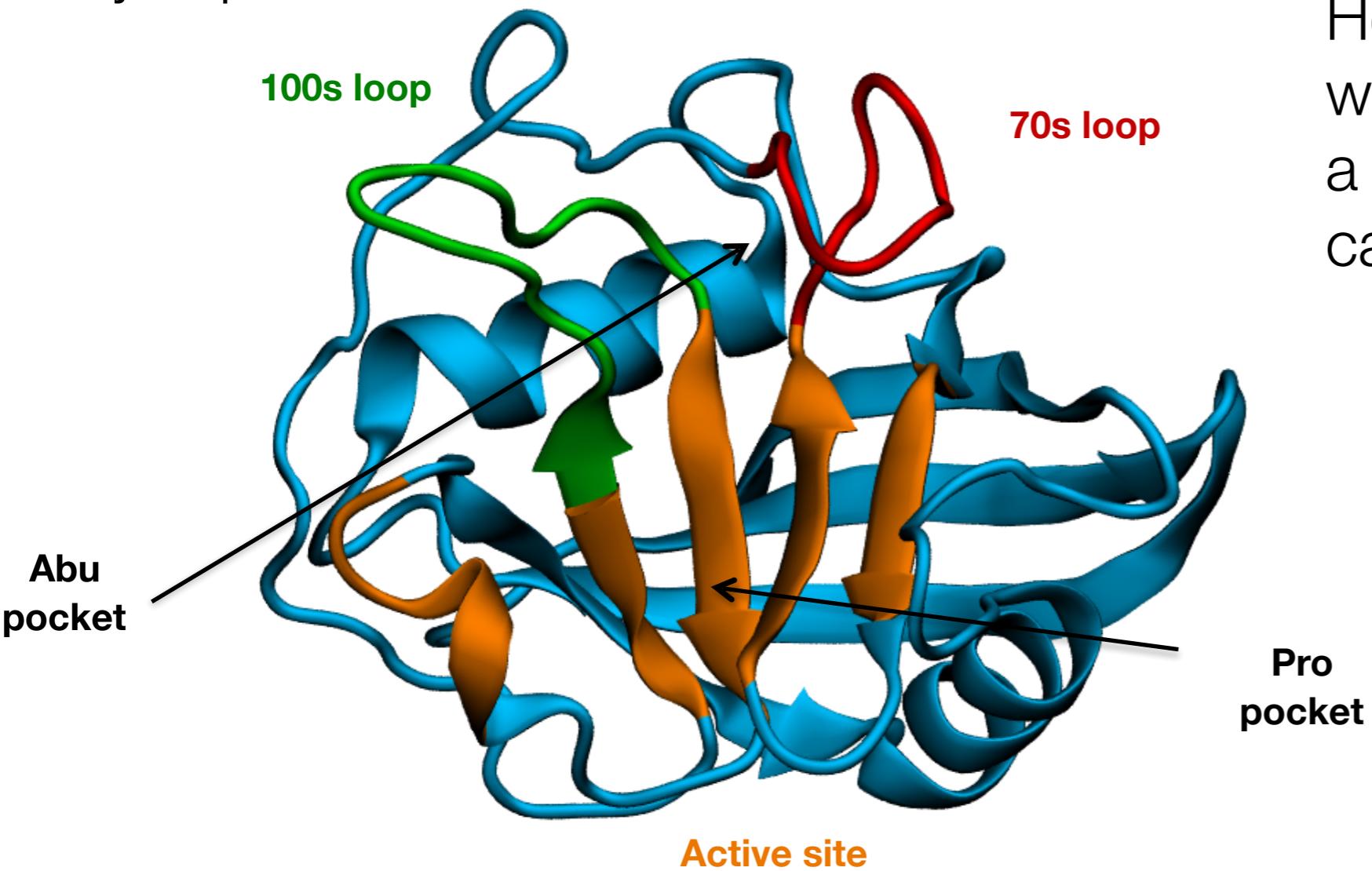


Cluster with your favourite clustering algorithm! (e.g. k-means)



Cyclophilin A – a heavily studied enzyme

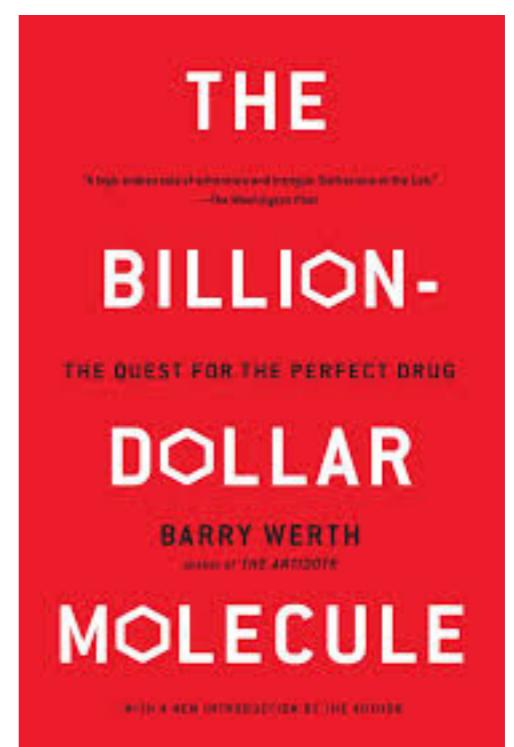
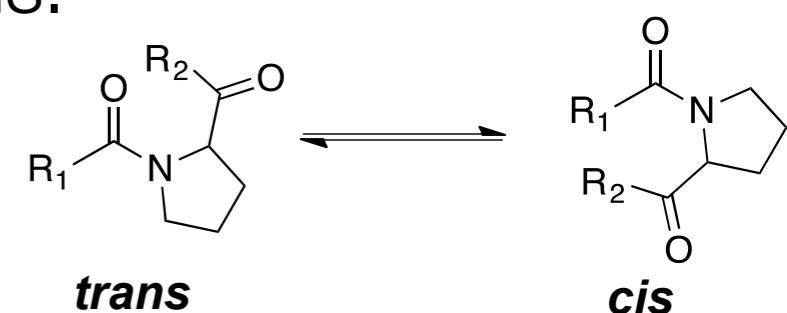
Cyclophilin A:



- 17 different isoforms with different functionality
- no isoform specific drug

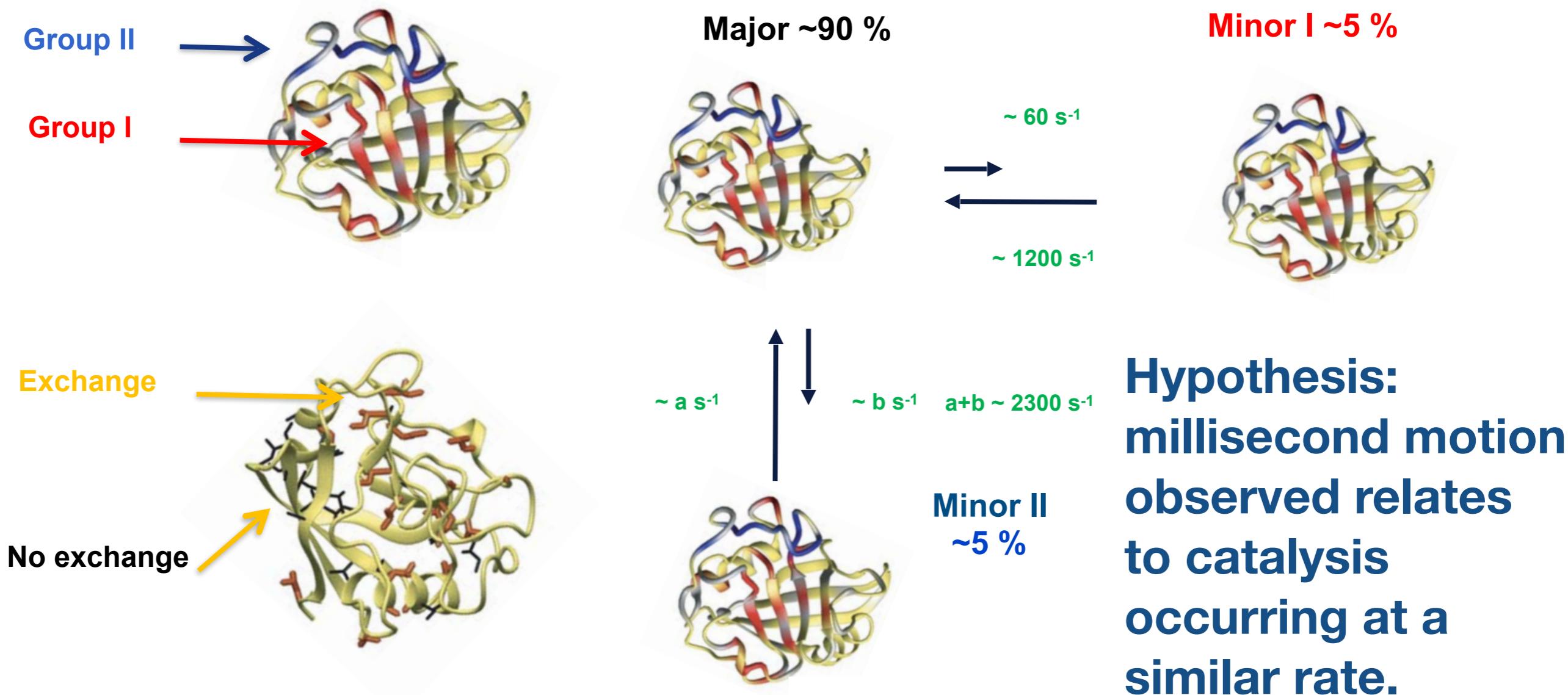
Cyclophilin A:

Hepatitis B, C and HI virus will ‘hijack’ CypA to use as a Chaperone in new virus capsids.



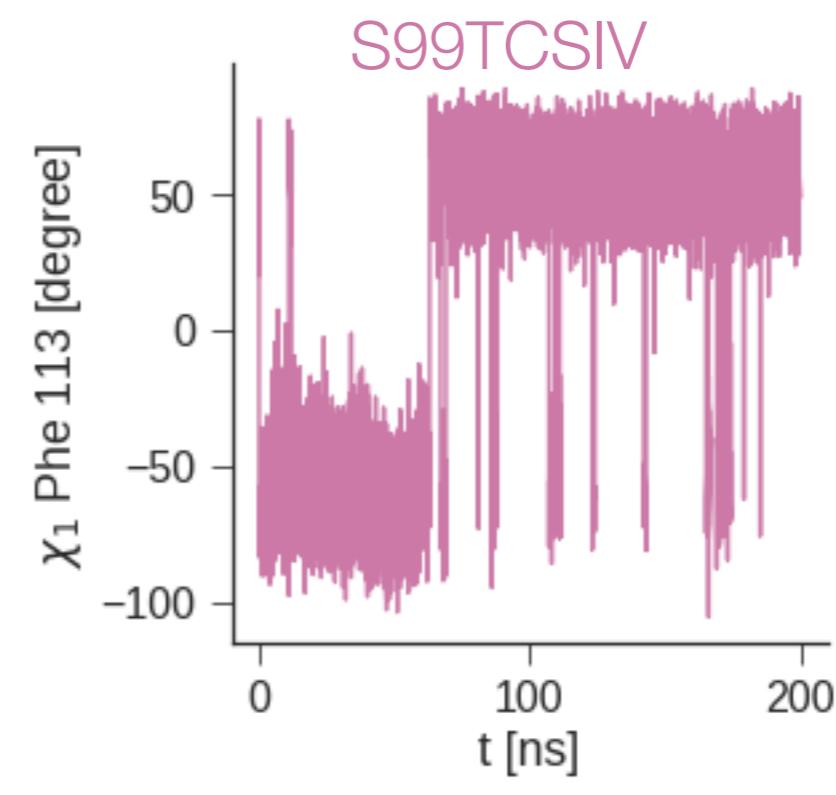
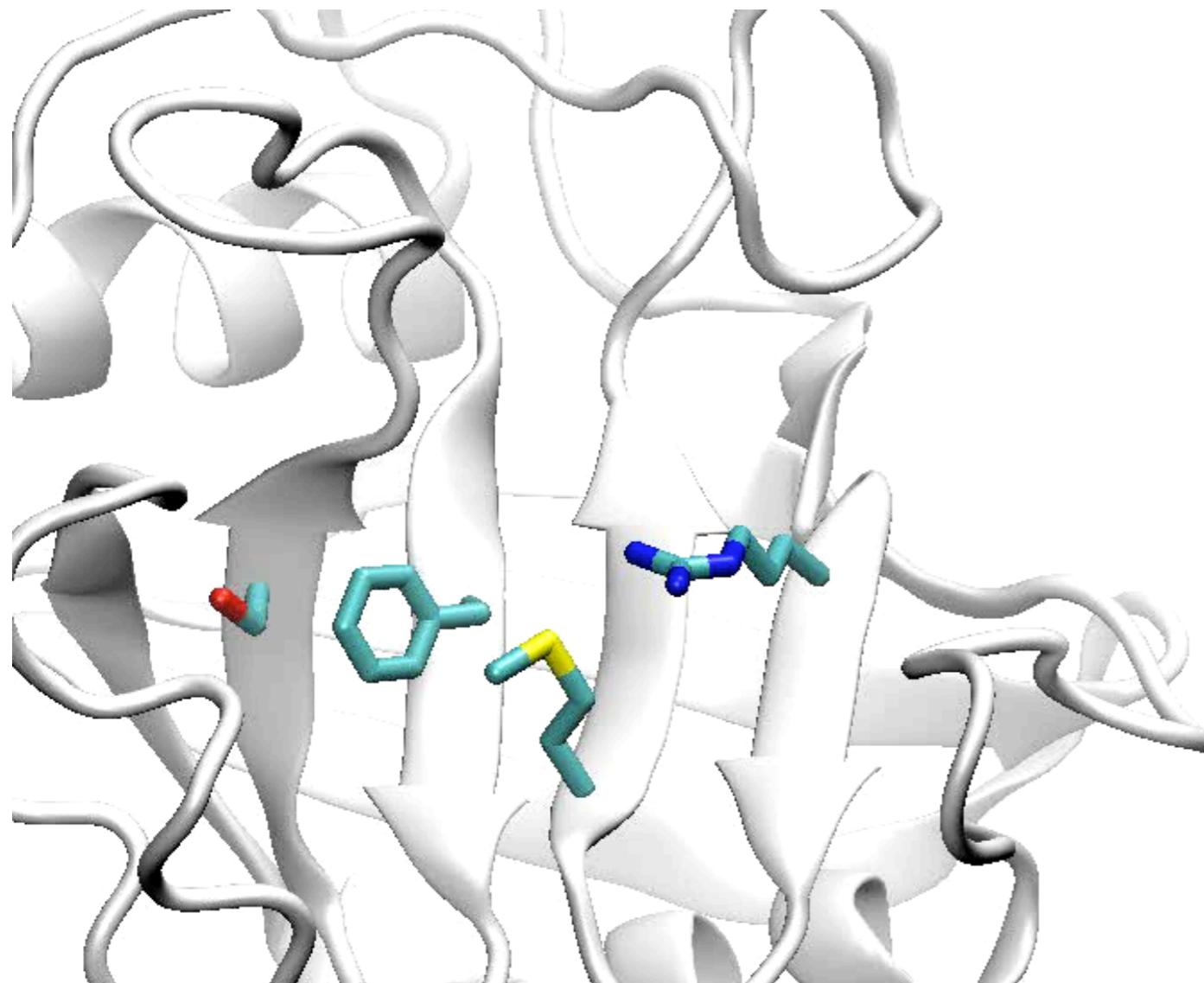
Cyclophilin A – a heavily studied enzyme

à NMR relaxation measurements on apo CypA show existence of two slow processes:



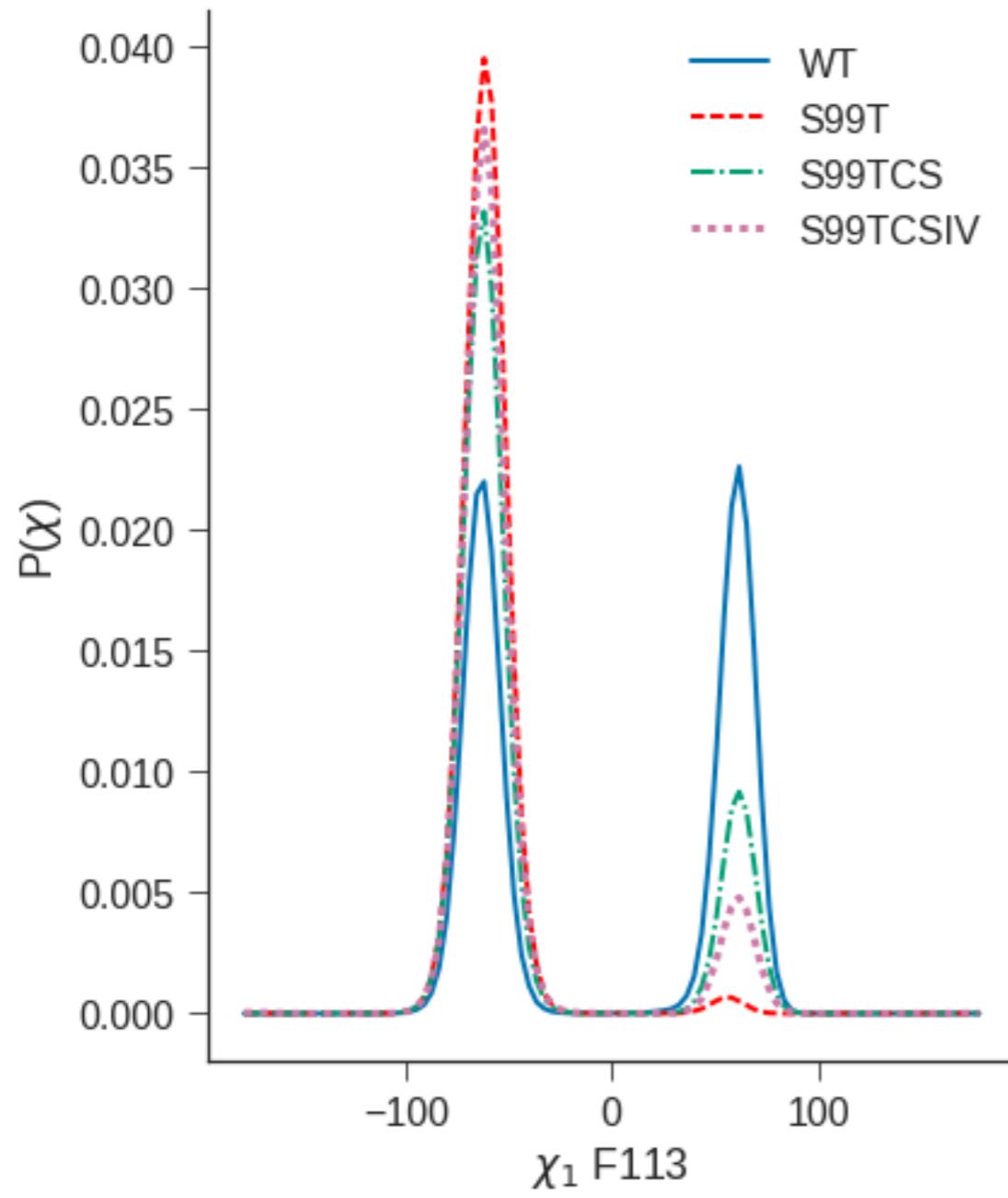
Hypothesis:
millisecond motion observed relates to catalysis occurring at a similar rate.

Dimensionality reduction with TICA



Dimensionality reduction with TICA

Stationary properties from MSM:

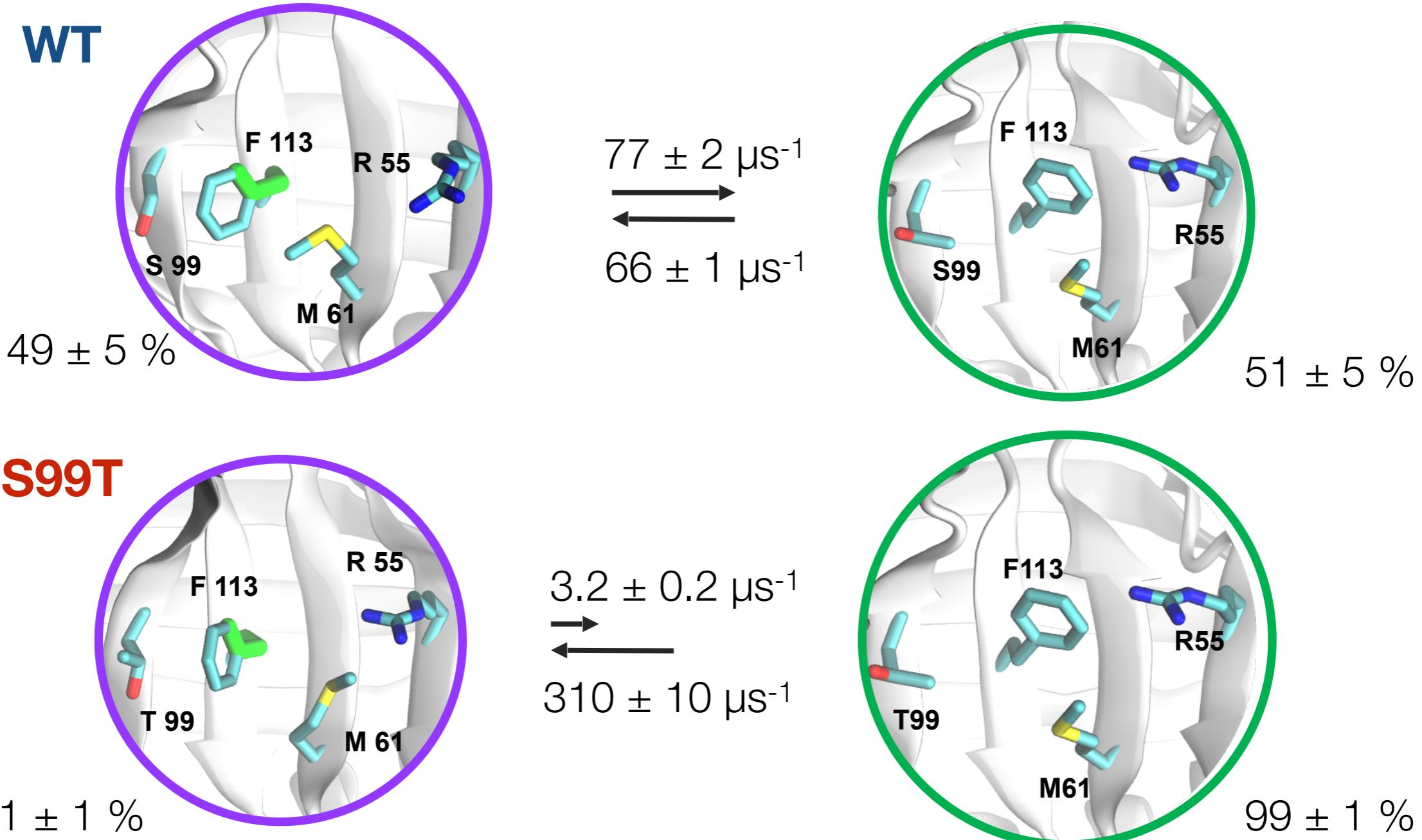


MSM on side chain dihedrals of residues that form part of group 1 residues and include Phe 113, χ_1 χ_3 Arg55, χ_1 χ_2 Met61, Ser/Thr99, Cys115, Ile97, Leu98, Gln 63



	MSM		X-ray occupancy	
Pop χ_1 :	-60°	+60°	-60°	+60°
WT:	0.51	0.49	0.37	0.63
S99T:	0.99	0.01	1	
S99TCS:	0.8	0.2	0.79	0.21
S99TCSIV:	0.9	0.1	1	

Dimensionality reduction with TICA



Outlook

- Markov models are away to extract equilibrium and dynamic properties from a timeseries.
- In biomolecular simulations this information can be used to guide experiments

Areas of ongoing reasearch

- Use ML to build multiple models and find heuristic for predicting the best model
- Difficult to find a unique metric to measure ‘best’ model
- Feature selection and dimensionality reduction
- Large spread of timescales, makes it difficult to focus on the process you are trying to understand (biomolecules)

Questions?

