

Markov state models for molecular dynamics

Michael Maduabum
Adithya Ganesh
Axel Sly



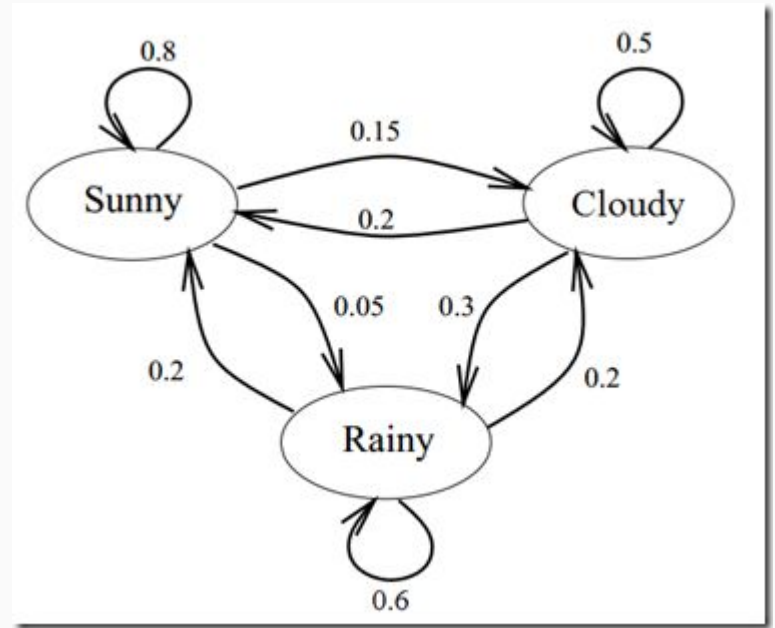
Everything you wanted to know about Markov State Models but were afraid to ask

Vijay Pande, Kyle Beauchamp, Gregory Bowman

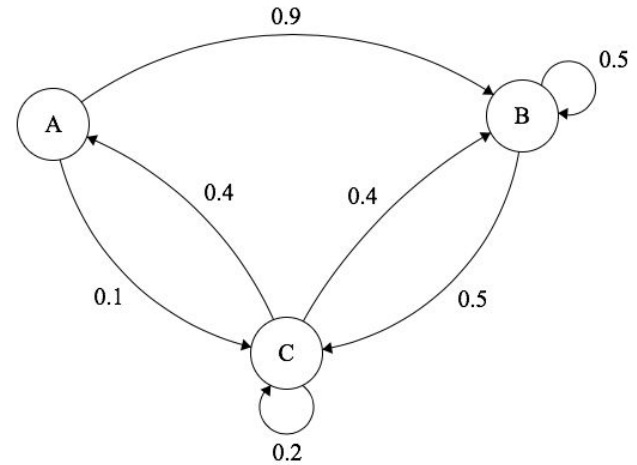
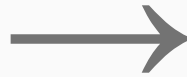
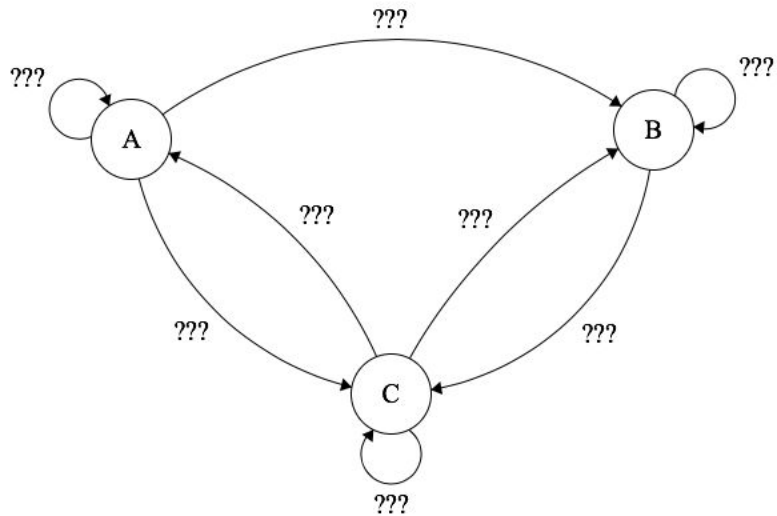


What is a Markov State Model?

- Way of modeling a random process
- Markov Assumption
 - Next state is determined only by current state

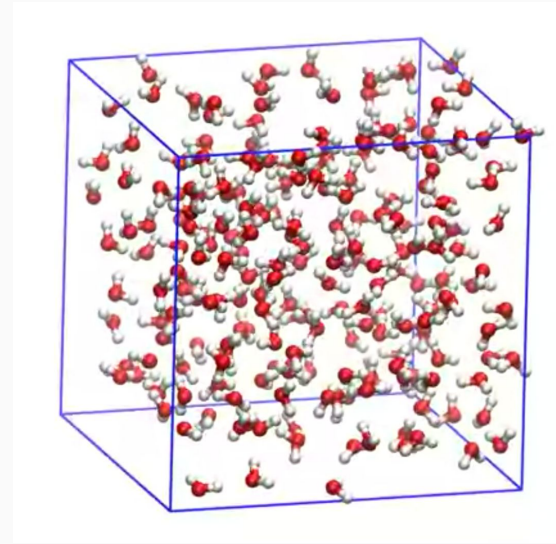


Transition probabilities



Applications

- Protein folding
- **Molecular Dynamics**
- Speech recognition
- Self-driving cars
- Gene prediction
- Machine translation
- Part-of-speech tagging



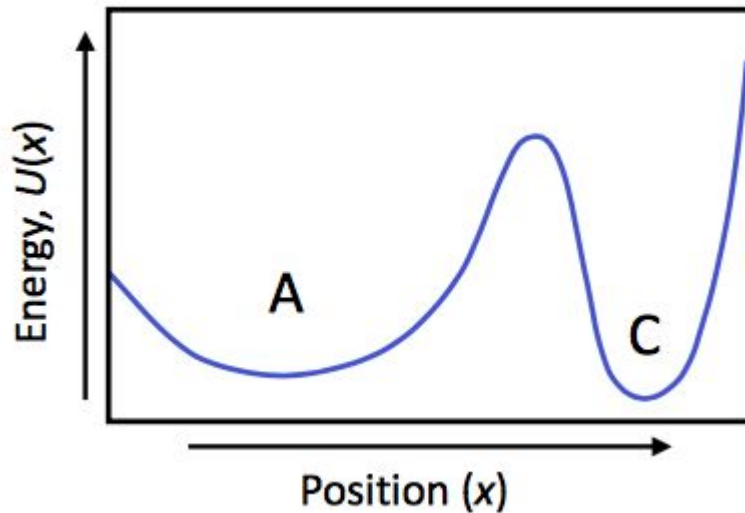
Challenges in Molecular Dynamics

- Computationally Intensive
 - Relevant simulation timescales milliseconds - seconds
- Force Fields are limited
- Analysis of simulation data
 - Trajectories



Building an MSM for MD

- Choose Markov states to correspond to kinetic macrostates
 - Similar conformations clustered to groups with a representative energy
- Build transitions matrix based on MD probabilities of state changes
 - Frequentist and Bayesian approaches



Molecular Dynamics with MSM's

- N states with transition probabilities between them
 - $N \sim$ thousands to millions
 - Large number of states gives more accurate representation
- Define states in a physically meaningful way
- Build the transition matrix efficiently
- Start with some initial data set
 - e.g. a molecular dynamics simulation

Adaptive Sampling

- Often there is not enough simulation data to immediately build an accurate Model
 - Need to perform more simulations to get novel data
- Use transition probability matrix to select macrostate to start simulation in
 - Statistical analysis to run simulations to generate the most data
- Significantly improves efficiency

Caveats

- Sampling is *still* a challenge
 - Long-timescale events are out of reach
- Limited by force fields
- State decomposition is difficult
 - Complex systems require more than just kinetic approaches
 - e.g. knots, topological effects
- Need new metrics for structural similarity

Questions?



Markov state models of biomolecular conformational dynamics

Chodera et al., 2014

Adithya Ganesh

Markov Chain Simulation

- Markov Chain simulation
 - <http://setosa.io/blog/2014/07/26/markov-chains/index.html>

Recent theoretical advances

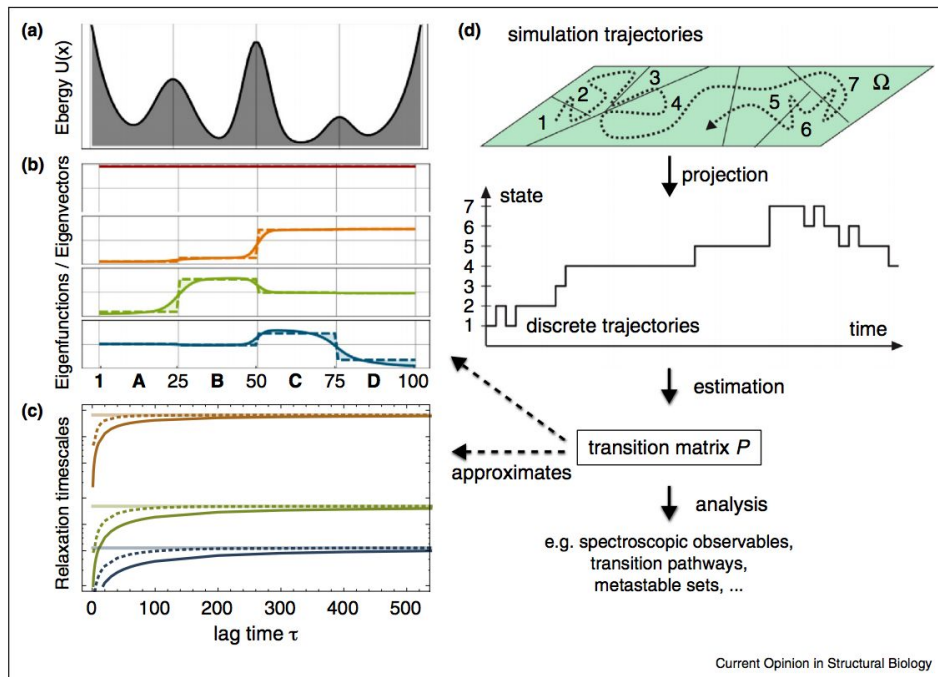
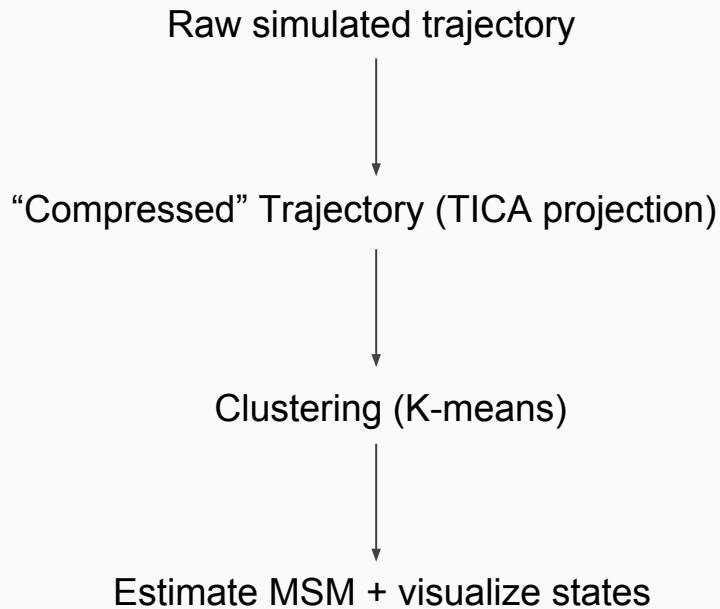
- Old approach
 - Maximizing “metastability” -- quantity related to *lifetime* of states
- Recent advances
 - Formulate analysis as an eigenvalue problem
- Key intuition
 - Eigenvector with eigenvalue 1 → steady state!

$$\boldsymbol{\pi}^T \mathbf{P} = \boldsymbol{\pi}^T .$$

Markov Chain Eigenvector Simulation

- **Key intuition**
 - Eigenvector (with $\lambda = 1$) \rightarrow steady state
- **Simulation**
 - <http://setosa.io/ev/eigenvectors-and-eigenvalues/>
- **PageRank**
 - *"The \$25,000,000,000 Eigenvector - The Linear Algebra Behind Google"*
 - <https://www.rose-hulman.edu/~bryan/googleFinalVersionFixed.pdf>
 - A very similar Markov chain problem
 - Obtain transition probabilities by observing which pages link to others
 - Eigenvectors give you the steady state "page ranks"

EMMA Python package -- demo



Questions?



HTMD: High-Throughput Molecular Dynamics for Molecular Discovery



Motivation

- Molecular dynamics simulations (MD) come with several limitations:
 - Data analysis
 - Reproducibility
 - Accuracy of force fields
 - Time sampling limitations
- Necessity of a more standardized methodology

High-Throughput Molecular Dynamics (HTMD)

- Programmable workspace for simulation based molecular discovery
- Python
- VMD

Integrated Platform

- Structure Manipulation
- System Building
- Molecular Simulation
- Adaptive Sampling
- Projecting and Clustering
- Markov State Models (MSMs)

2. METHODS

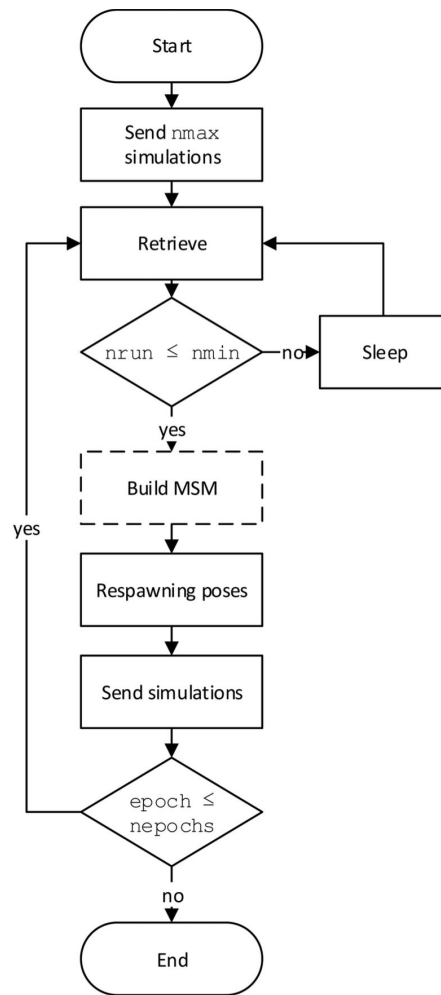
2.1. Integrated Platform for Molecular Modeling.

Listing 1: Manipulating molecular structures.

```
1 # Download the PDB structure of Barnase
2 mol = Molecule('2F4Y')
3 # Set the chain ID of Barnase to 'Y'
4 mol.set('chain', 'Y', sel='protein')
5 # Remove the carbon alpha of residue 15
6 mol.remove('residue 15 and name CA')
7 # Append Barstar to Barnase
8 mol.append(Molecule('2HXX'))
9 # Visualize in VMD or NGL
10 mol.view()
11 # Solvate in water box
12 solvate(mol, minmax=[[-50, -50, -50],
13                      [50, 50, 50]])
```

Adaptive Sampling

- Efficient and intelligent sampling of conformational space
- From previous simulation knowledge, it identifies conformational subspaces that are under-sampled



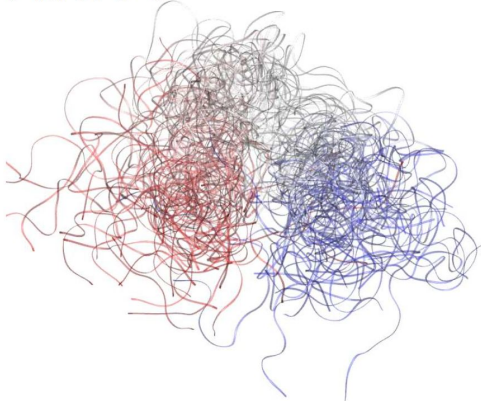
Protein Folding Analysis

```
1 # Collect and link simulation files
2 sims = simlist(glob('data/*/'),
  ↪ glob('input/*/structure.pdb'))
3 # Calculate protein contact maps
4 met = Metric(sims)
5 met.projection( MetricSelfDistance(
  ↪ 'protein and name CA',
  ↪ metric='contacts' ) )
6 data = met.project()
7 # Project on slowest 10 TICA dimensions
8 tica = TICA(data, 20)
9 dataTica = tica.project(10)
10 # Cluster data with kmeans
11 dataTica.cluster(
  ↪ MiniBatchKMeans(n_clusters=1000))
12 # Link a model with the data
13 model = Model(dataTica)
14 # Calculate the Markov model
15 model.markovModel(300, 4)
16 # Macrostates equilibrium distribution
17 model.eqDistribution()
18 # Visualize the macrostates in VMD
19 model.viewStates()
20 # Calculate the kinetics
21 kin = Kinetics(model, temperature=360)
22 kin.getRates()
```

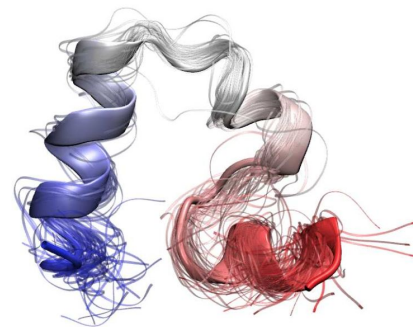
Protein Folding Analysis

The 4 macrostates produced by the Markov state model of Villin

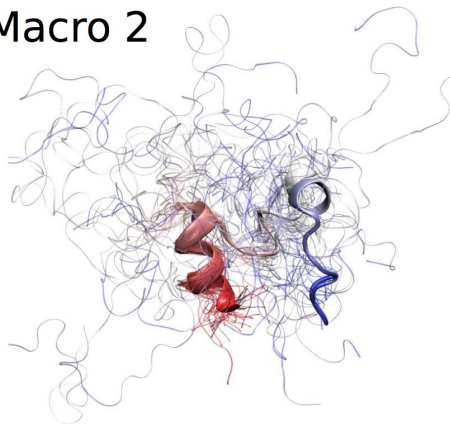
Macro 0



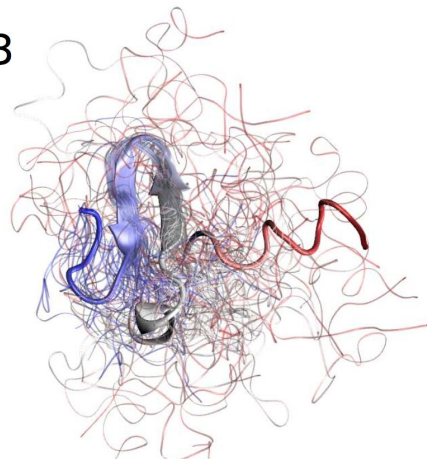
Macro 1



Macro 2



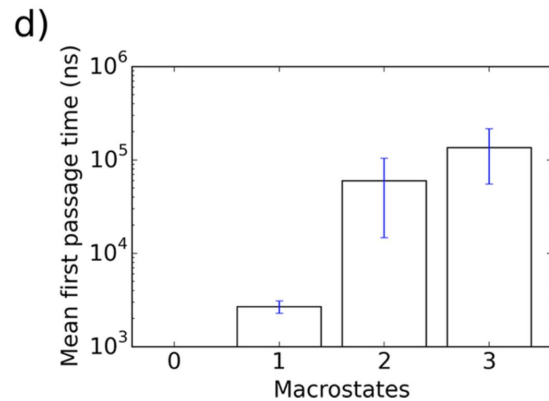
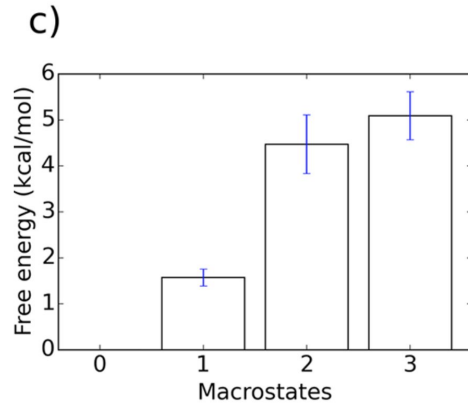
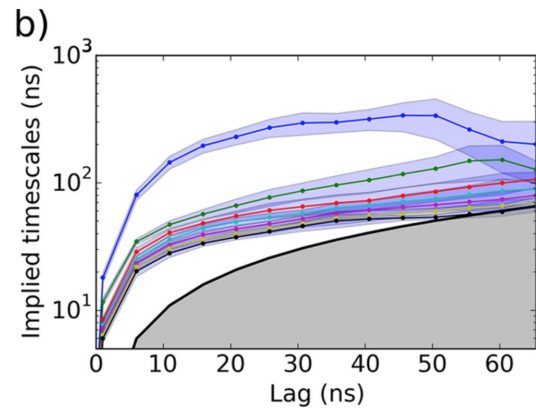
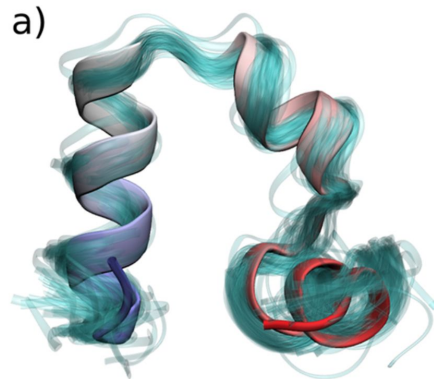
Macro 3



Questions?



Protein Folding Analysis



Adaptive Sampling Protein Folding

