

# Markov State Model Review

H.Koh  
2024.11.14.

# Overview

## Non – Markovian Dynamical Approaches

Molecular  
Dynamics



Markov State  
Model



qMSM



IGME

- + Complementary to experimental approaches
- + Provide atomistic details of protein dynamics
- Required timescale often too large

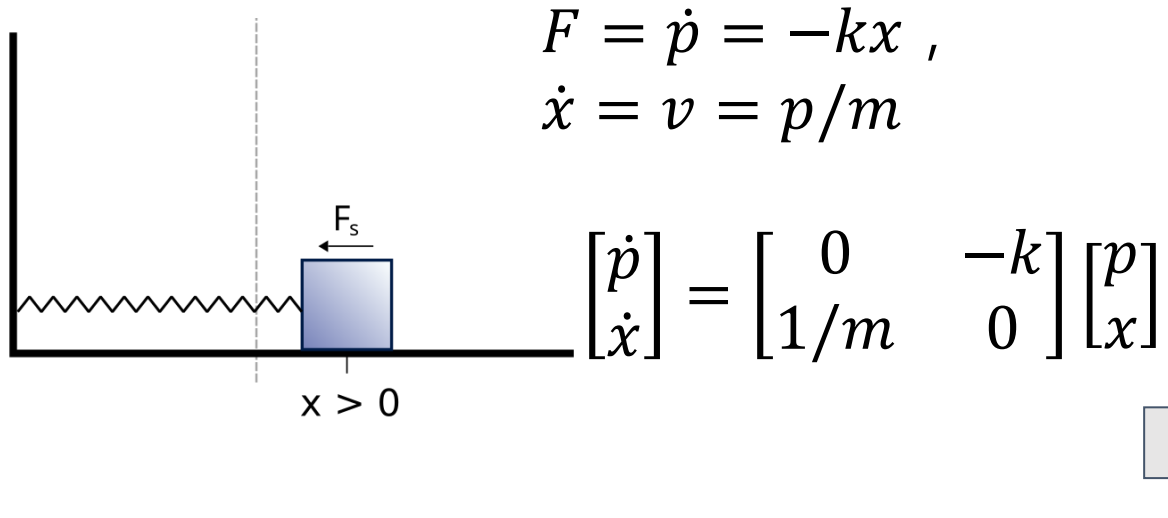
- + Long timescale dynamics
- + A large ensemble of short MD
- Hard to meet Memoryless conditions

- + Introduce Memory Kernel  $K$
- + Generalized Master Equation
- Issue of robustness
- Numerical instabilities in  $K$

- + Integrate Memory Kernel first
- + Fit analytic results with MD
- + Enhanced numerical stability

# Markov State Model

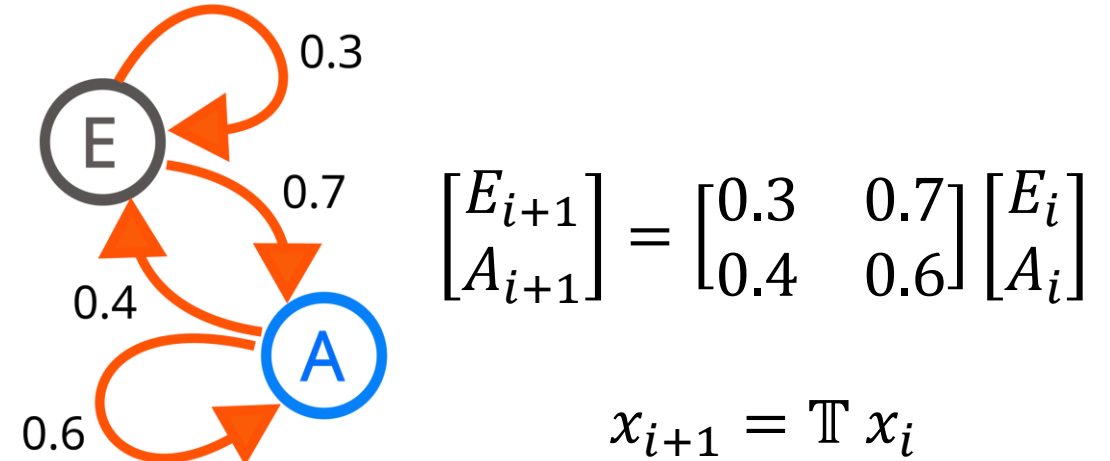
Wikipedia : Harmonic Oscillator



## Newtonian Physics

- once you know position and velocity (momentum), every dynamics is determined
- does not depend on previous history
- generalization : Liouville equation ?

Wikipedia : Markov Chain



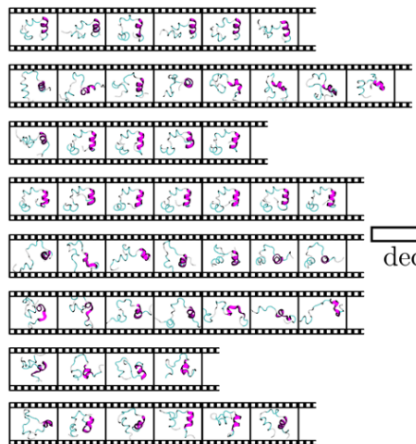
## Markov State Model ( $x(t + \tau) = \mathbb{T}(\tau) x(t)$ )

- $n \times n$  transition probability matrix  $\mathbb{T}(\tau)$
- entire configuration space -> divided into  $n$  states ( $n \times 1$  column vector  $x(t)$  )
- lag time  $\tau$
- Memoryless condition (Markovian)  
:  $\mathbb{T}$  does not depend on previous history (  $x(t)$  )

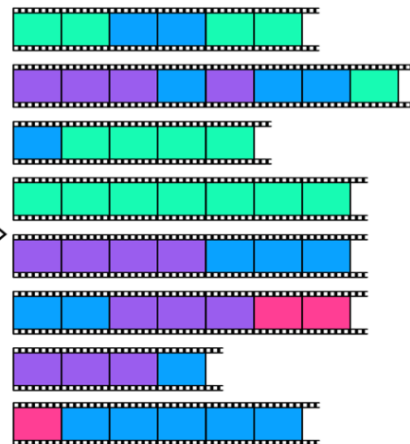
# Markov State Model

J. Am. Chem. Soc. 2018,  
140, 2386–2396

(a) simulation data



(b) state assignments



state  
decomposition

(c) observed transition  
counts matrix

	Green	Blue	Purple	Pink
Green	11	1	0	0
Blue	3	9	2	0
Purple	0	4	9	1
Pink	0	1	0	1

$A$

(d) reversible transition  
probability matrix

	Green	Blue	Purple	Pink
Green	0.85	0.15	0	0
Blue	0.14	0.62	0.21	0.03
Purple	0	0.24	0.72	0.04
Pink	0	0.25	0.25	0.5

$T$

$= \text{normalized } (A + A^T)/2$

If we can choose

- $n$  states that **capture the dynamics** of the system
- long enough  $\tau$  to be **Markovian**
- short enough  $\tau$  to **resolve the system dynamics**

then MSM can predict long-timescale dynamics based on a large ensemble of short MD simulations.

However, MSM has **limitations**

- **dynamic relaxation time**  $< \tau$
- $\tau < \text{length of individual MD simulations}$  available to estimate transition probabilities
- **large  $n$**   $\rightarrow$  **fast relaxation time**, but **hinder** the comprehension of biological mechanisms (Over fitting?)

Solution : **Non-Markovian Dynamics Approaches?**

# Generalized Master Equation (GME) Method

$$\dot{T}(t) = \dot{T}(0)T(t) - \int_0^t d\tau \mathcal{K}(\tau)T(t-\tau)$$

$\mathcal{K}(\tau)$  : memory kernel

$\tau_K$  : memory kernel relaxation time

$$\mathcal{K}(t \geq \tau_K) \approx 0$$

$$\frac{\dot{T}(t)}{T(t)} = \frac{\dot{T}(0)}{T(0)} - \int_0^t d\tau \mathcal{K}(\tau) \frac{T(t-\tau)}{T(t)}, \quad T(0) = \mathbb{I}_{n \times n}$$

- If  $\mathcal{K}(\tau) = 0$  (No memory),

$$\frac{\dot{T}(t)}{T(t)} = \frac{\dot{T}(0)}{T(0)} \rightarrow \frac{d \ln T(t)}{dt} = C \rightarrow T(t) = \exp(Ct)$$

- If  $\mathcal{K}(\tau) \neq 0$ ,

$$\ddot{T}(t) = \dot{T}(0)\dot{T}(t) - \mathcal{K}(t)T(0) = \dot{T}(0)\dot{T}(t) - \mathcal{K}(t)$$

$$\mathcal{K}(t) = \dot{T}(0)\dot{T}(t) - \ddot{T}(t)$$

Memory kernel  $\mathcal{K}(t)$  is something related with the second derivative of function (matrix)  $T(t)$  ?

# Generalized Master Equation (GME) Method

$$\dot{T}(t) = T(t)\dot{T}(0) - \int_0^{\min[\tau_K, t]} T(t-s)K(s)ds,$$

**Brute-force method** utilizing above equation

→ **quasi MSM (qMSM)**

- Theoretically rigorous
- Issue in ensuring the **Robustness**
- **Numerical instability** in the computation of  $K(s)$

Integrate until  $K(s)$  is fully decayed.



$$M_n = \int_0^{\tau_K} K(s)s^n ds$$

$$T(t)^{-1} \frac{d}{dt} T(t) = \dot{T}(0) - M_0 - \sum_{n=1}^{\infty} \left[ \frac{(-1)^n}{n!} T(t)^{-1} \frac{d^n}{dt^n} T(t) \right] M_n$$

**Integration of Memory Kernel  $K(s)$**   
(avoid numerical instability)

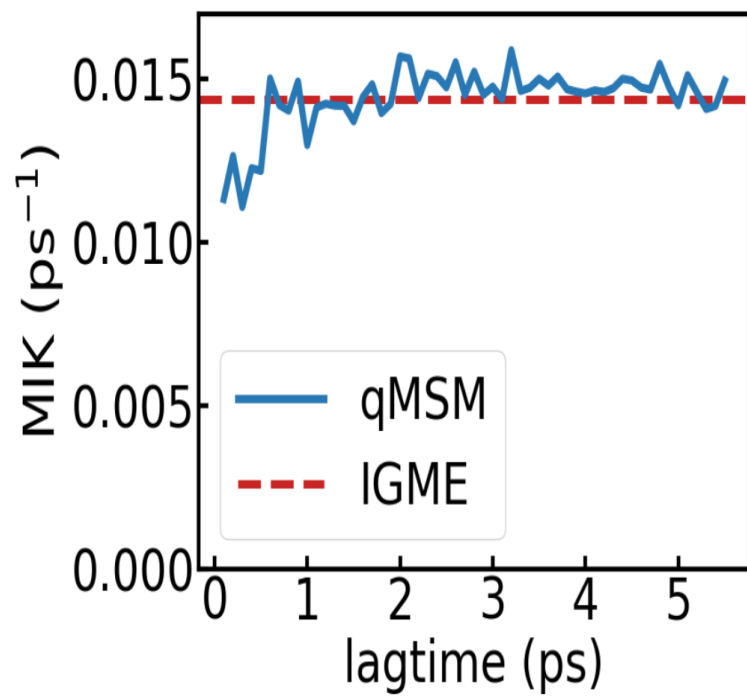
→ **Integrative GME (IGME)**

- Derive **Analytic Solution** for the GME
- **Fitting to MD** simulations
- Determine **hyperparameters**

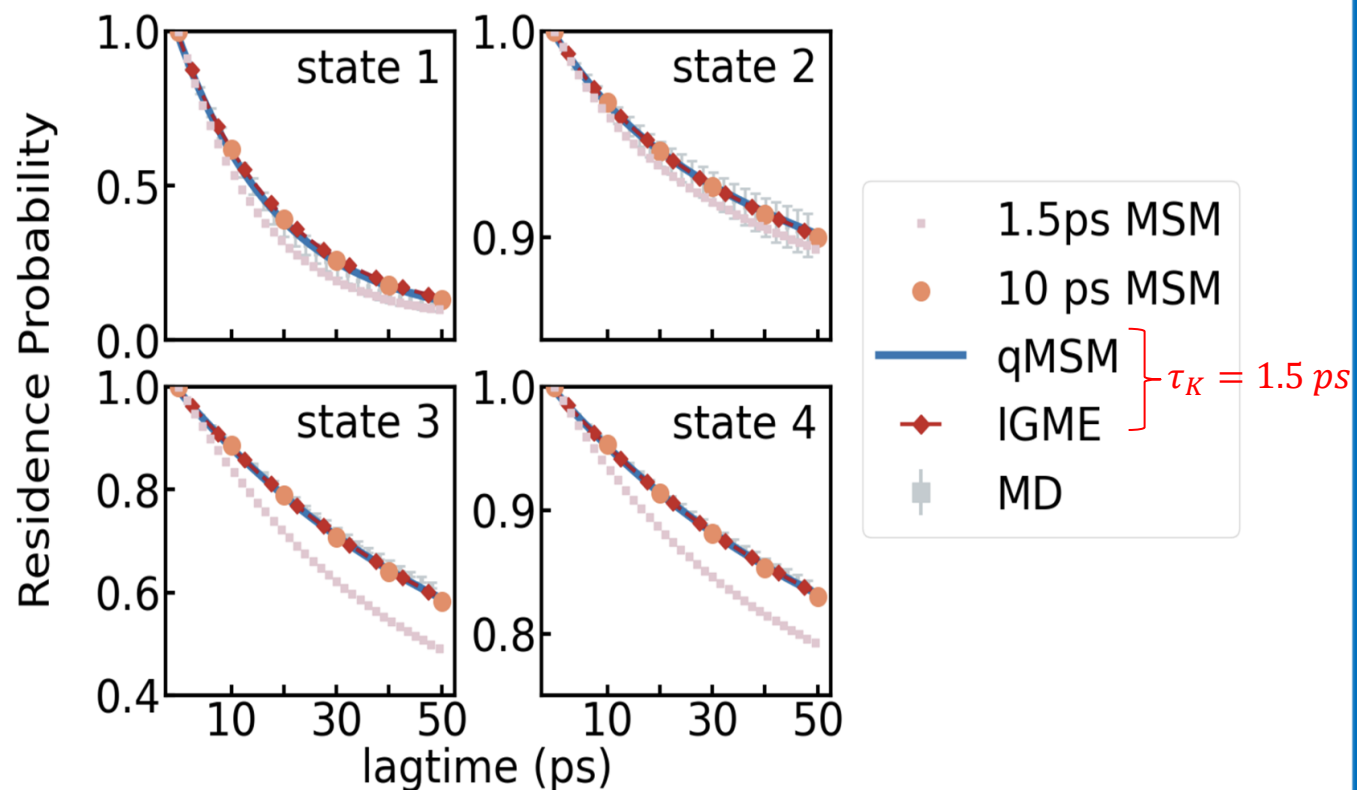
# Generalized Master Equation (GME) Method

## ❖ Alanine Dipeptide Tutorial

### Mean Integrated K (kernel)



### Comparison b/w different approaches



# Code Review

## ❖ GME Tutorials - Python 3.11.9

- Repository: [GME\\_tutorials on GitHub](#)

## ❖ msmbuilder2022

- Repository: [msmbuilder2022 on GitHub](#)

## ❖ Installation Command:

```
pip install  
git+https://github.com/msmbuilder/msmbuilder2022.git
```

## ❖ Common Installation Issue

- Error: `distutils.errors.DistutilsPlatformError: Microsoft Visual C++ 14.0 or greater is required`
- Solution:
  1. Refer to [this guide](#).
  2. Download [Visual C++ Build Tools](#).
  3. Select only Desktop development with C++ during installation.
  4. Restart your computer after installation.

## ❖ Issues with Python 3.12

- Fastcluster Compatibility
  - On Windows with Python 3.12, `msmbuilder2022` installation may trigger a `fastcluster` error.
- Numpy Compatibility
  - Python 3.12 includes a numpy version higher than `1.23.5`, causing errors in `igme.top_outputs` in `villin_headpiece_tutorial.ipynb`.
  - Solution:  
`pip install numpy==1.23.5`



MD + Time-  
resolved  
Spectroscopy



Software  
Engineering  
Web-app

What should I  
do now?