# Mechanism beyond Markov models: How and why to use non-Markovian analysis of trajectory data

Ernesto Suárez

2017

# MSM in Science

Number of scientific publications per year

2016:  20,700

2015:  20,500

.

.

.

2000: 4,600

# MSM in Computational Chemistry

# Markov State Models (MSM)

Widely used to analyze and interpret molecular trajectories. The final goal is to infer long time behavior.

Main assumption: The Markov property

$$k_{ij}(\tau) = P\{X_{t+\tau} = j \,|\, X_t = i\}$$

Regular simulations are Markovian in their full continuous phase spaces. However any discrete partition of the phase space generates non-Markovian trajectories.

# Learning process in MSM



$$\begin{pmatrix} \hat{k}_{11} & \hat{k}_{12} & \hat{k}_{13} & \cdots \\ \hat{k}_{21} & \hat{k}_{22} & \hat{k}_{23} & \\ \hat{k}_{31} & \hat{k}_{32} & \hat{k}_{33} & \\ \vdots & & & \ddots \end{pmatrix}$$

$$\hat{k}_{ij}(\tau) = c_{ij} / c_i$$

$$\mathbf{K}^T \mathbf{p} = \mathbf{p}$$

Biased for kinetics

Voelz et al., J. Am. Chem. Soc., 2010, 132(5), pp 1526-1528

# Protein Models



Trp-cage  208μs

Chignolin  106μs

Villin  125μs

NTL9  1100μs

Shaw et al., Science 2011, 334(6055), pp. 517-520

# MSM Analysis: Standard Recipe

1. Divide the space in "Markovian" regions

2. Estimate parameters and <span style="color:red">select a lag time</span>

3. Analysis

# Estimating kinetic properties

## Mean First Passage Time (MFPT)

# Mean First Passage Time (MFPT)



$\tau$ = The lag-time >> integration time step δt

$$\mathrm{MFPT} = \frac{1}{k_{AB}}$$

NOT MSM

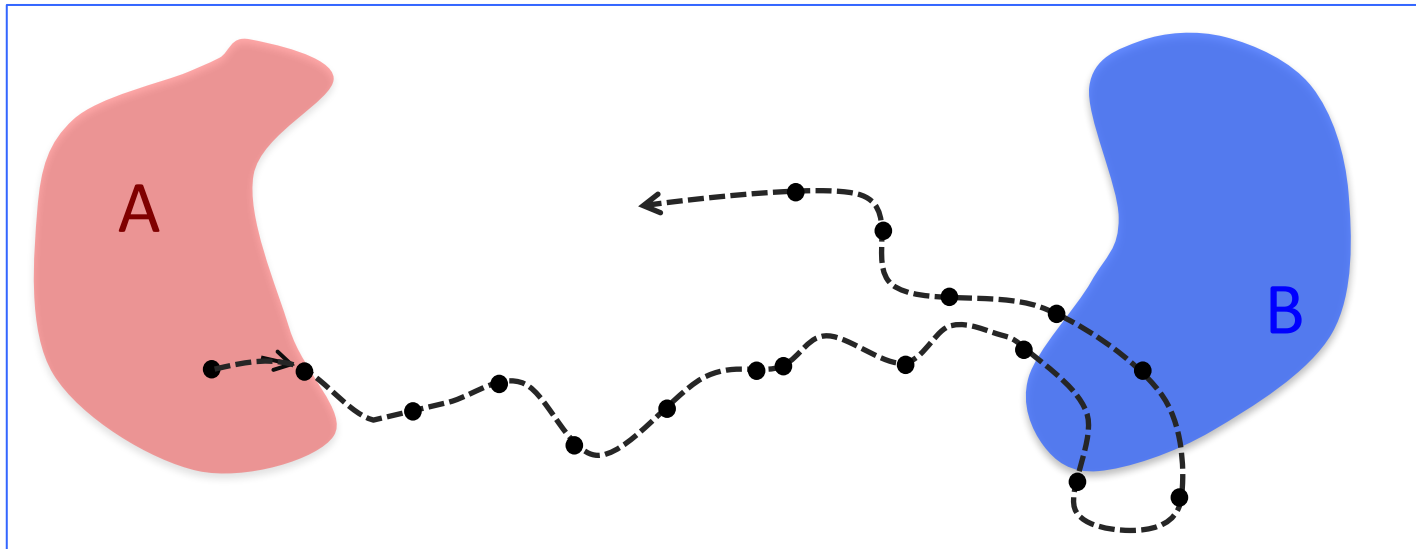# Mean First Passage Time (MFPT)



Direct average,
no MSM or any other
model assumption

$$\text{MFPT} = \frac{1}{N} \sum_i^N FPT_i$$

$$\sim \frac{\tau}{p^{eq}(B)}$$

MFPT
estimation
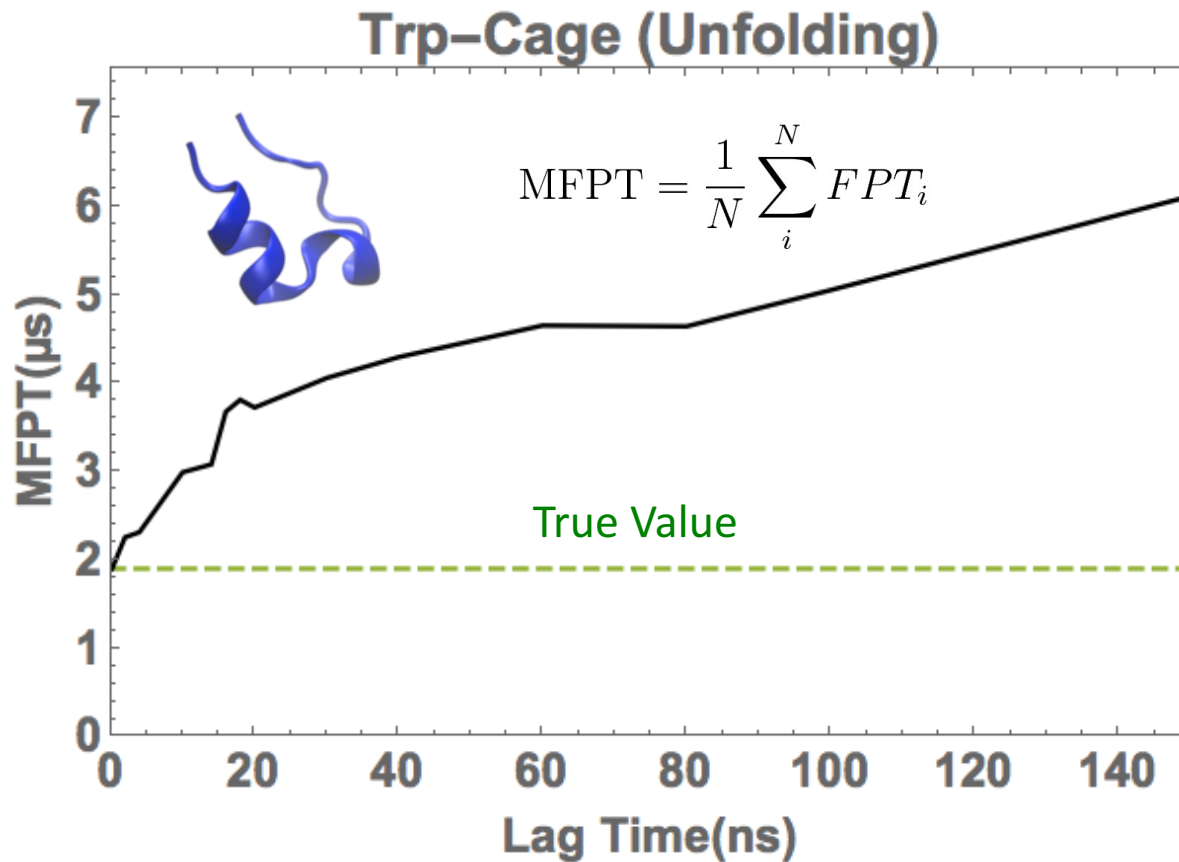
Exact Value

$\tau$

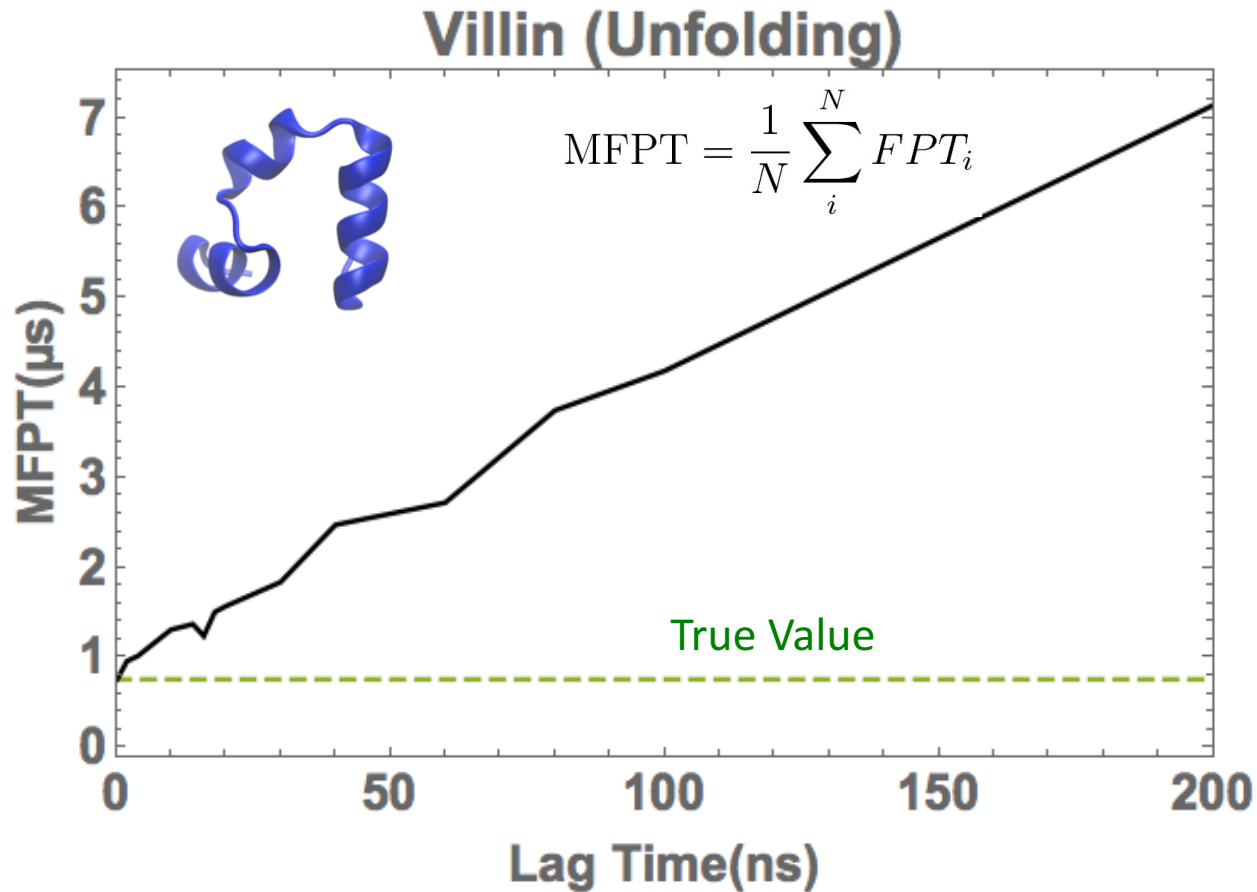Our estimation of the kinetic properties are lag-time dependent while the thermodynamic properties are the same for every lag-time.

# MFPT vs lag-time



Trp–Cage (Unfolding)

$$\mathrm{MFPT} = \frac{1}{N} \sum_{i}^{N} FPT_i$$

True Value

No MSM or any other model assumption

# MFPT vs lag-time



Villin (Unfolding)

$$\text{MFPT} = \frac{1}{N} \sum_i^N FPT_i$$

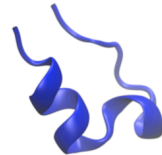True Value
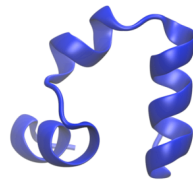
No MSM or any other model assumption
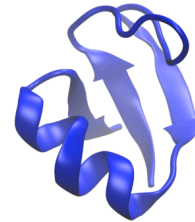
# Data choices

Long MD Simulations



Trp-cage 208μs

Chignolin 106μs

Villin 125μs

NTL9 1100μs

**Full data set**
- Markovian
- Non-Markovian

**Reduced data set (< 5%)**
- Non-Markovian

# Markov State Models

# Markov MFPT vs lag-time

**Chignolin**

### Chignolin (Folding)

Direct

$$\mathrm{MFPT} = \frac{1}{N}\sum_i^N FPT_i$$

True Value

MSM

MFPT(µs)

$\tau$ = Lag Time (ns)

### Chignolin (Unfolding)

MSM

Direct

True Value

$\tau$ = Lag Time (ns)

# Markov MFPT vs lag-time

**Trp-cage**



**Trp–Cage (Folding)**

Direct $\mathrm{MFPT} = \dfrac{1}{N}\sum_i^N FPT_i$

True Value

MSM

MFPT(µs)

Lag Time(ns)

**Trp–Cage (Unfolding)**

Direct

MSM

True Value

Lag Time(ns)

# Markov MFPT vs lag-time

**Villin**



**Villin (Folding)**

Direct

$$\text{MFPT} = \frac{1}{N} \sum_i^N FPT_i$$

MSM

True Value

MFPT(μs)

Lag Time(ns)

**Villin (Unfolding)**

Direct

MSM

True Value

Lag Time(ns)

# Markov MFPT vs lag-time

**NTL9**

# MSM Analysis

- Biased for kinetic properties
  - Discretization error ⬆MFTP
  - Markov error ⬇MFPT

# Non-Markovian Analysis
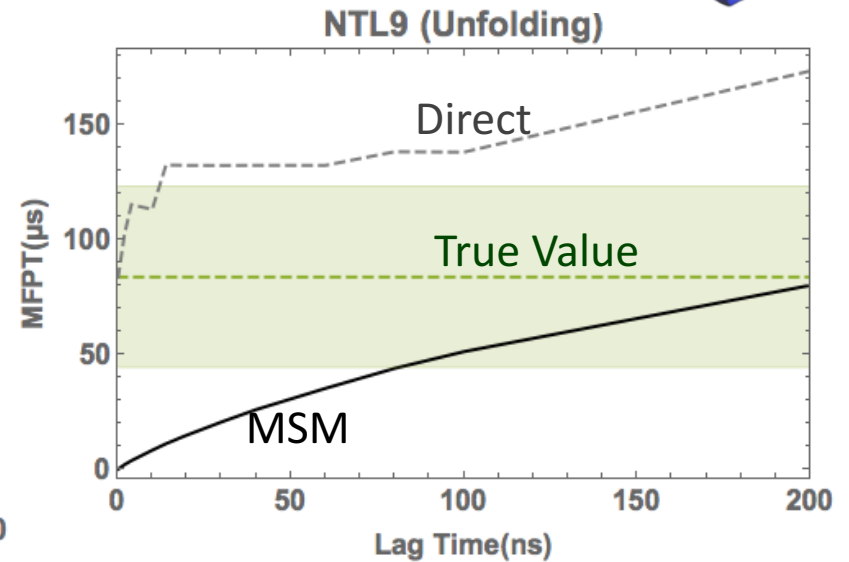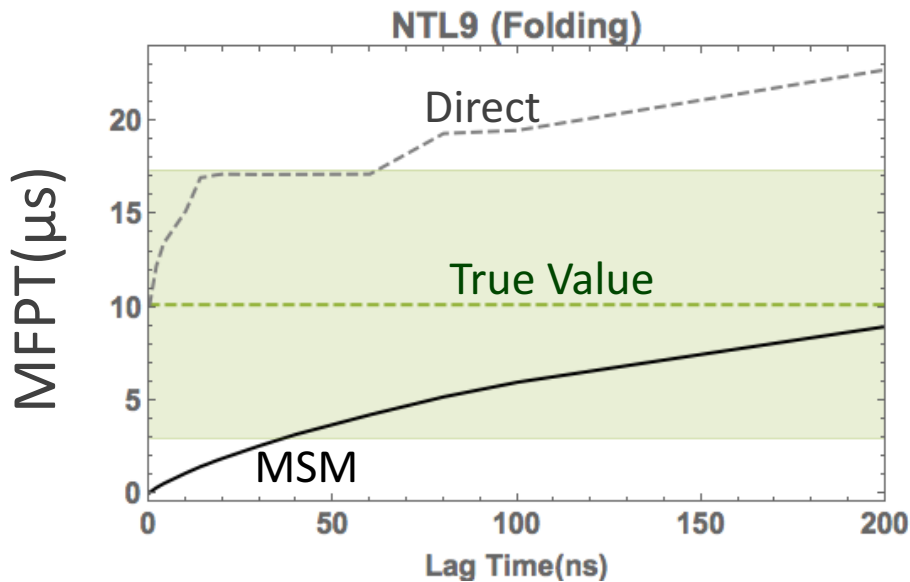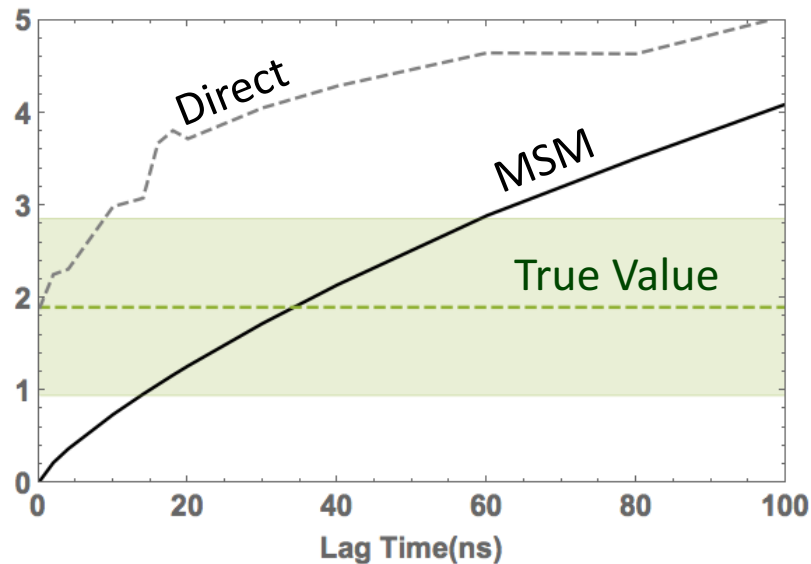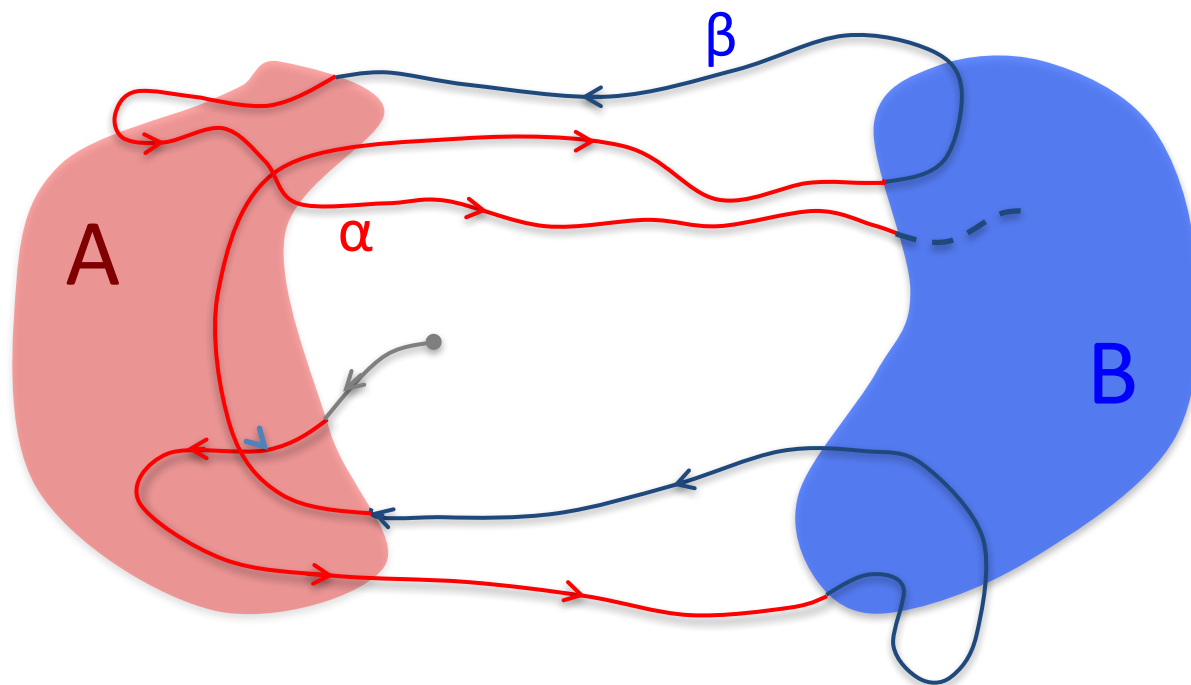
# Beyond Markov: Color

α = Last in A    β = Last in B

Suarez et al., J. Chem. Theory Comput., 2014, 10 (7), pp 2658–2667
Vanden-Eijnden et al., J. Chem. Phys., 2009, 131(4), pp 44120

# Beyond Markov: Color

$$k_{ij}(\tau) = P\{X_{t+\tau} = j | X_t = i\}$$

$$k_{ij}^{\mu\nu}(\tau) = P\{X_{t+\tau} = j, L_{t+\tau} = \nu | X_t = i, L_{t+\tau} = \mu\} \quad \mu, \nu = \alpha, \beta$$

Suarez et al., J. Chem. Theory Comput., 2014, 10 (7), pp 2658–2667
Vanden-Eijnden et al., J. Chem. Phys., 2009, 131(4), pp 44120

# Beyond Markov: Color

Unbiased kinetics

*2N x 2N*
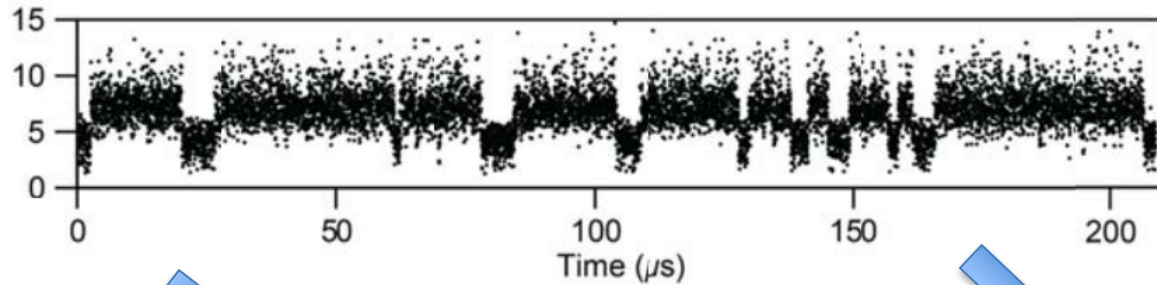
Biased kinetics

*N x N*

$$\begin{bmatrix} k_{11} & : & k_{12} & : & k_{13} \\ \cdots & & \cdots & & \cdots \\ k_{21} & : & k_{22} & : & k_{23} \\ \cdots & & \cdots & & \cdots \\ k_{31} & : & k_{32} & : & k_{33} \end{bmatrix} \Longrightarrow \begin{bmatrix} k_{11}^{\alpha\alpha} & 0 & : & k_{12}^{\alpha\alpha} & 0 & : & 0 & k_{13}^{\alpha\beta} \\ 0 & 0 & : & 0 & 0 & : & 0 & 0 \\ \cdots & & & \cdots & & & & \cdots \\ k_{21}^{\alpha\alpha} & 0 & : & k_{22}^{\alpha\alpha} & 0 & : & 0 & k_{23}^{\alpha\beta} \\ k_{21}^{\beta\alpha} & 0 & : & 0 & k_{22}^{\beta\beta} & : & 0 & k_{23}^{\beta\beta} \\ \cdots & & & \cdots & & & & \cdots \\ 0 & 0 & : & 0 & 0 & : & 0 & 0 \\ k_{31}^{\beta\alpha} & 0 & : & 0 & k_{32}^{\beta\beta} & : & 0 & k_{33}^{\beta\beta} \end{bmatrix}$$

Example with 3 bins. *A* is defined as bin 1 and *B* as bin 2

$$\mathcal{K}^T \mathbf{p}^\mu = \mathbf{p}^\mu \qquad p_i^{\mathrm{eq}} = p_i^\alpha + p_i^\beta$$

Suárez et al., J. Chem. Theory Comput., 2014, 10 (7), pp 2658–2667

# MSM vs Non-Markovian Analysis



MSM

$$\begin{bmatrix} k_{11} & \vdots & k & \vdots & k_{13} \\ \cdots & N \times N & \cdots \\ k_{21} & \vdots & k_{22} & \vdots & k_{23} \\ \cdots \\ k_{31} & \vdots & k_{32} & \vdots & k_{33} \end{bmatrix}$$
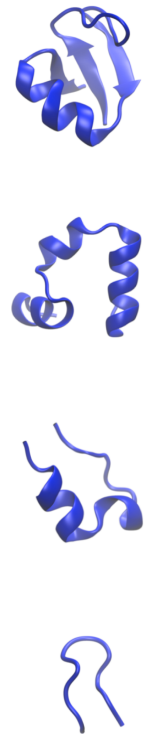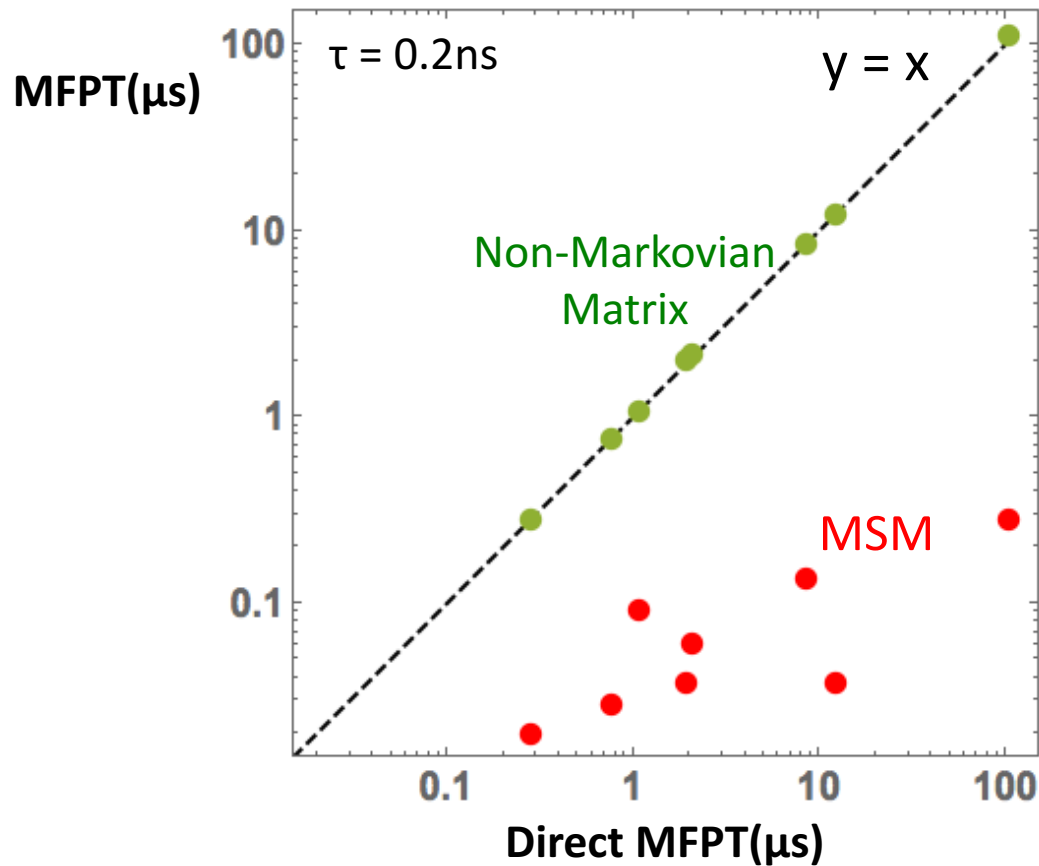
MFPT

Non-Markovian
Matrix

$$\begin{bmatrix} k_{11}^{\alpha\alpha} & 0 & \vdots & k_{12}^{\alpha\alpha} & 0 & \vdots & 0 & k_{13}^{\alpha\beta} \\ 0 & 0 & \vdots & 0 & 0 & \vdots & 0 & 0 \\ \cdots & 2N \times 2N & \cdots \\ k_{21}^{\alpha\alpha} & 0 & \vdots & k_{22} & 0 & \vdots & 0 & k_{23}^{\alpha\beta} \\ k_{21}^{\beta\alpha} & 0 & \vdots & 0 & k_{22}^{\beta\beta} & \vdots & 0 & k_{23}^{\beta\beta} \\ \cdots \\ 0 & 0 & \vdots & 0 & 0 & \vdots & 0 & 0 \\ k_{31}^{\beta\alpha} & 0 & \vdots & 0 & k_{32}^{\beta\beta} & \vdots & 0 & k_{33}^{\beta\beta} \end{bmatrix}$$

MFPT

# MSM vs Non-Markovian Analysis

No lag-time optimization



τ = 0.2ns

y = x

MFPT(μs)

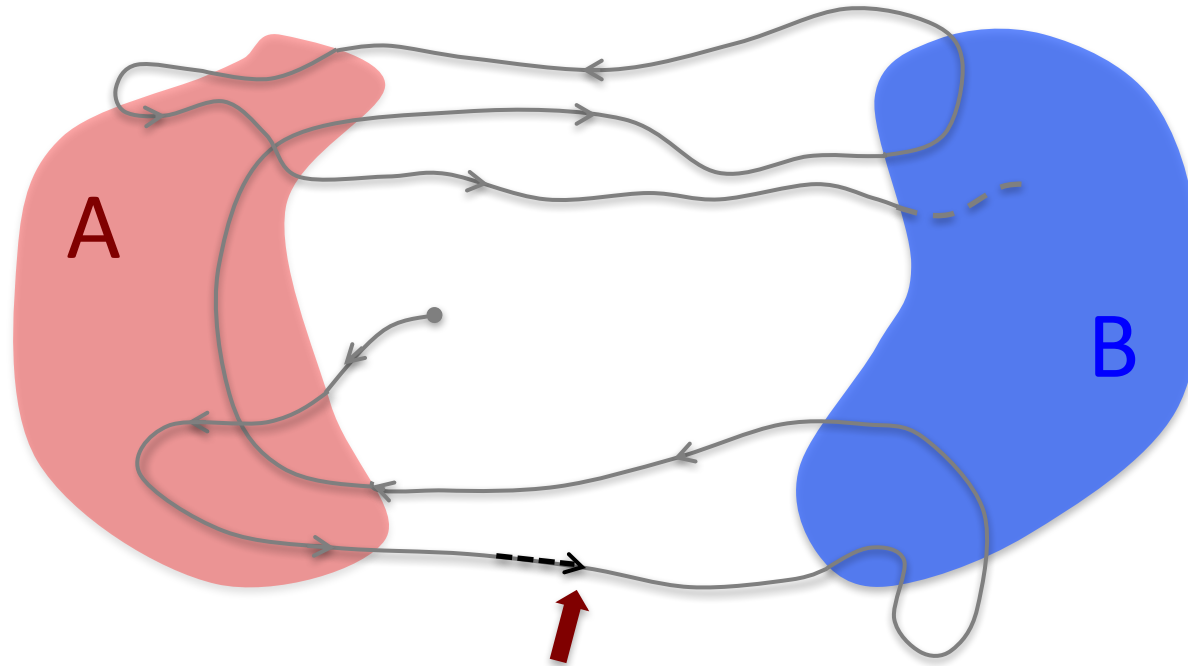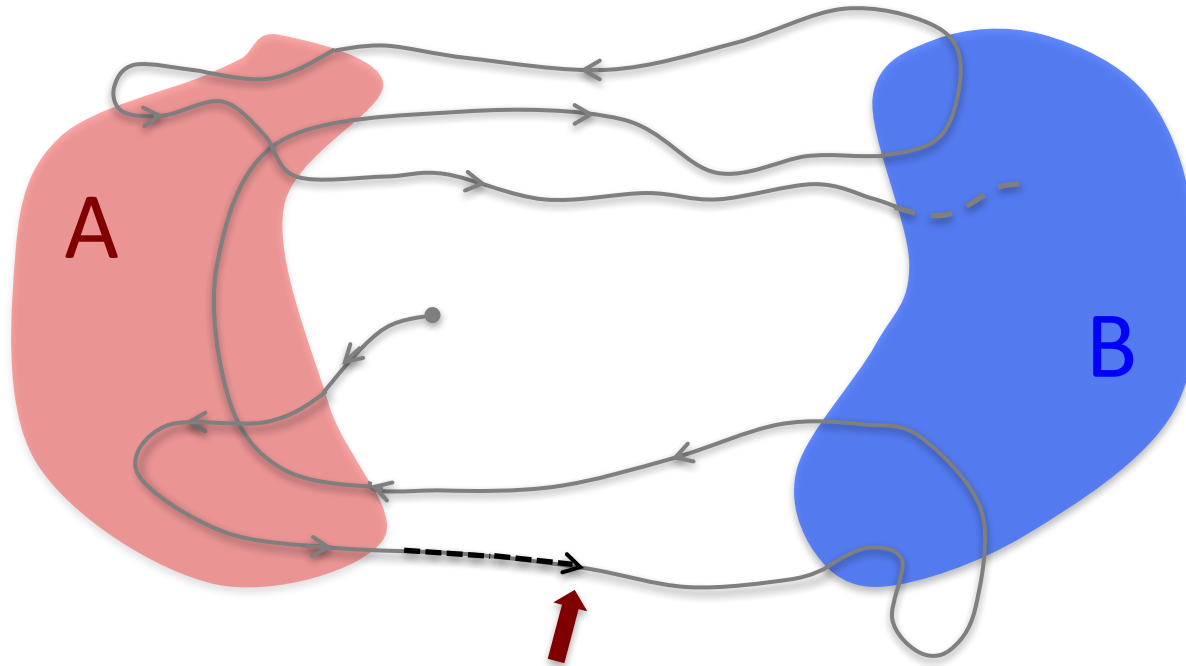Non-Markovian Matrix

MSM

Direct MFPT(μs)

# Non-Markovian Analysis

With sufficient history (color) information we get

- Unbiased thermodynamics (populations)
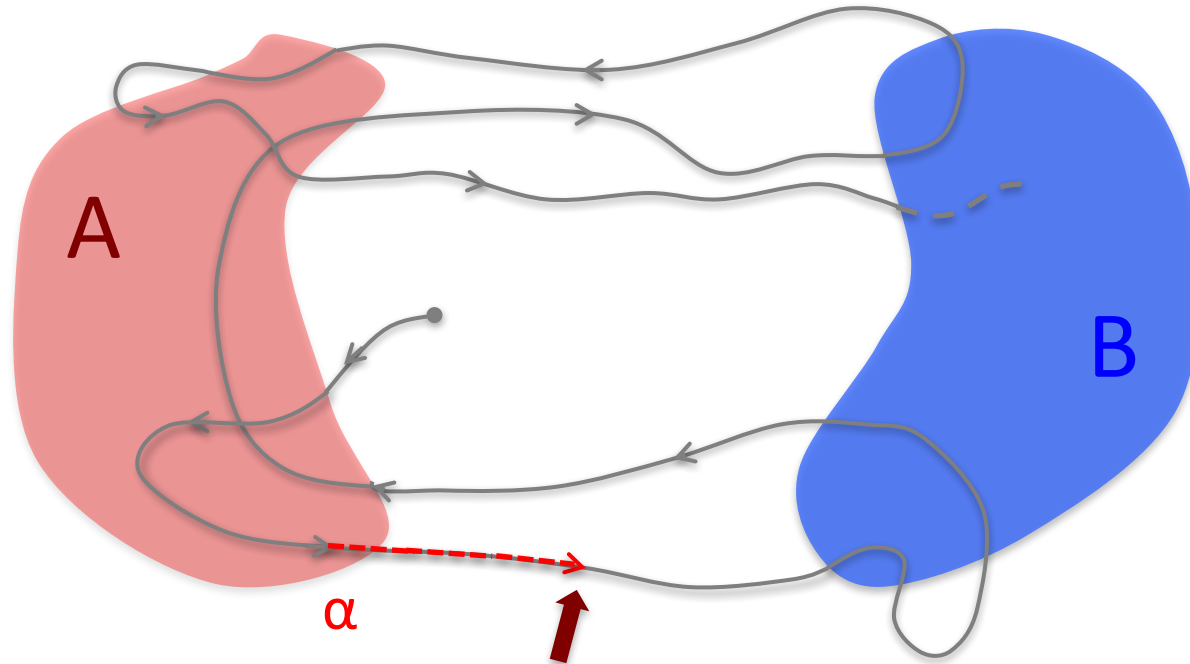
- Unbiased MFPT ($\tau \longrightarrow 0$)

Suarez et al., J. Chem. Theory Comput., 2014, 10 (7), pp 2658–2667

# Limited color/history info

# Limited color/history info

# Limited color/history info
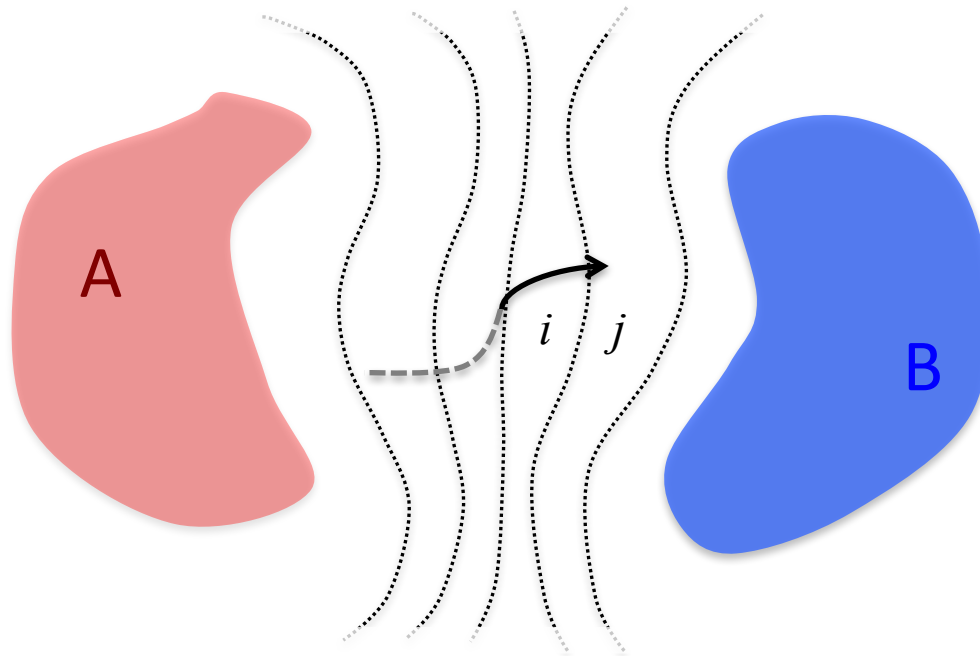
# Other non-Markovian Analyses

## Markov + Color

When examining a given time point of the trajectory for estimating a labeled rate, the α or β label are assigned if possible given the amount of history. Otherwise the label is assigned stochastically assuming a Markov behavior.
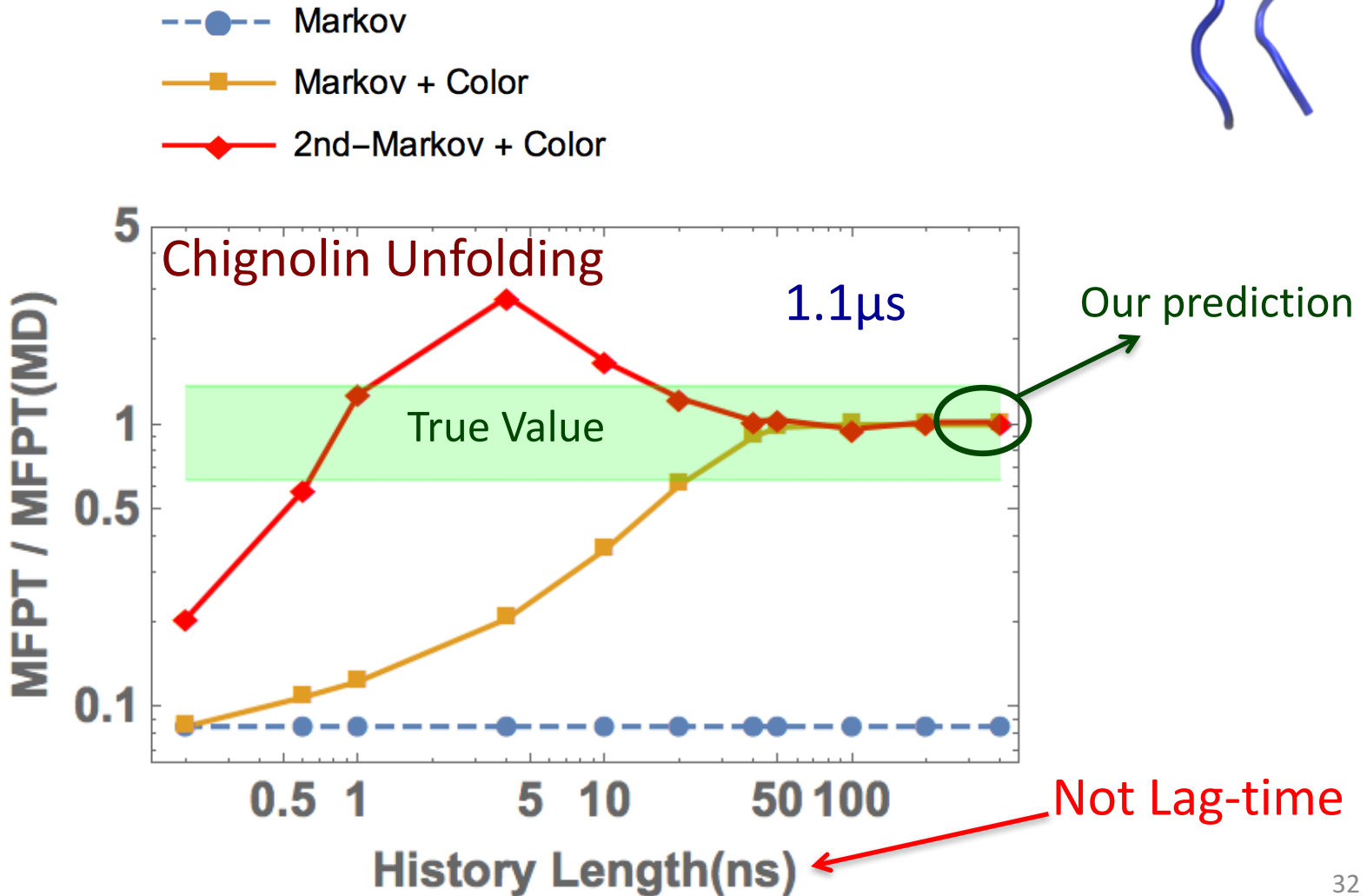
# Other non-Markovian Analyses

## 2nd Order Markov + Color

When examining a given time point of the trajectory for estimating a labeled rate, the α or β label are assigned if possible given the amount of history. Otherwise the label is assigned stochastically assuming a 2nd-Order Markov model.
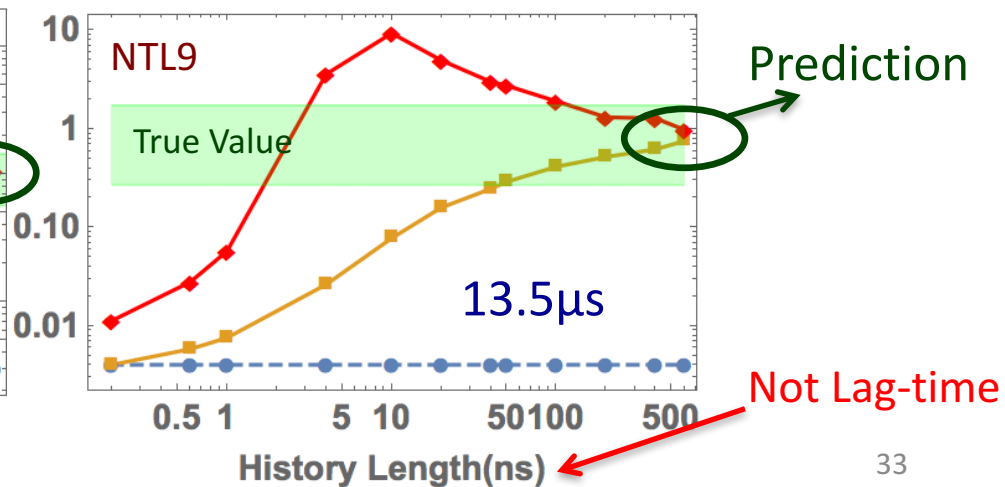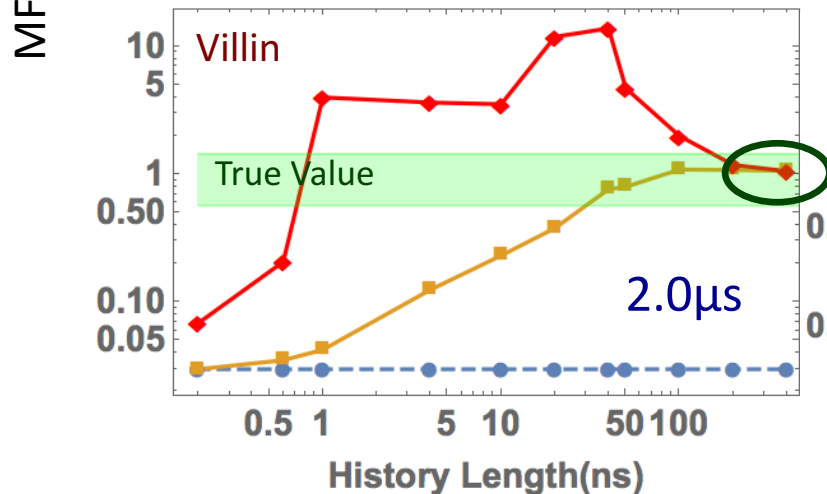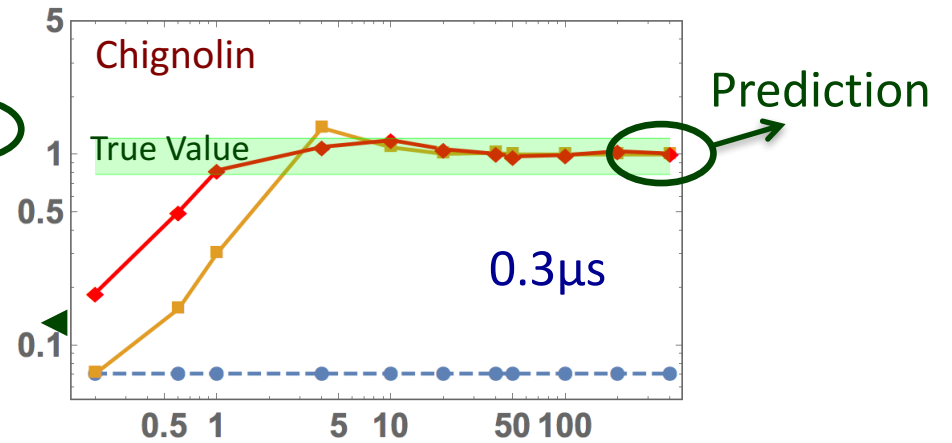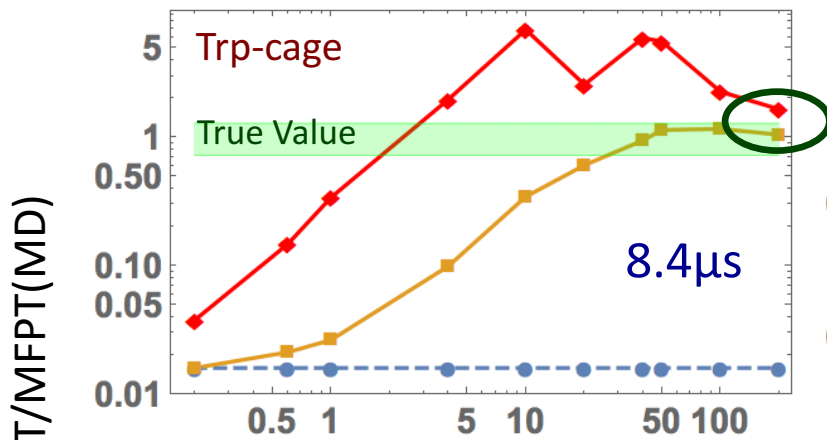
# Non-Markovian Analyses

# Non-Markovian Analyses (Folding)



Legend: ---●--- Markov ; —■— Markov + Color ; —◆— 2nd−Markov + Color

MSMBuilder States

Trp-cage — 8.4µs

Chignolin — 0.3µs — Prediction

Villin — 2.0µs

NTL9 — 13.5µs — Prediction

True Value
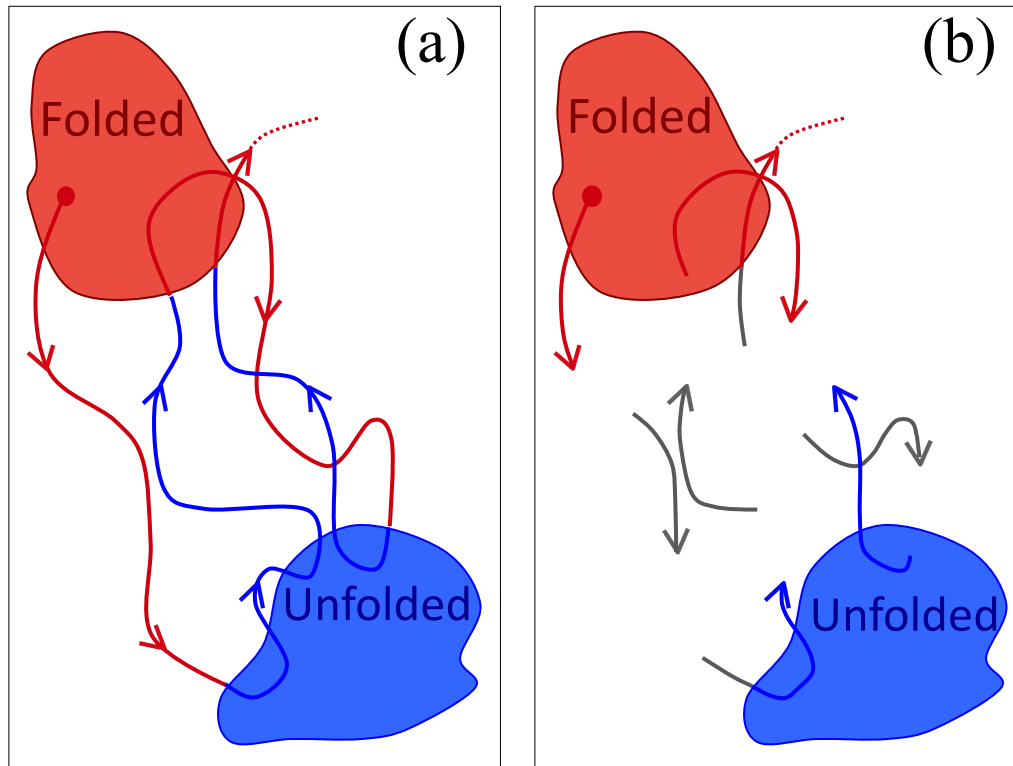
MFPT/MFPT(MD) vs History Length(ns)

Not Lag-time

# Reduced data set:

# Non-Markovian Analysis
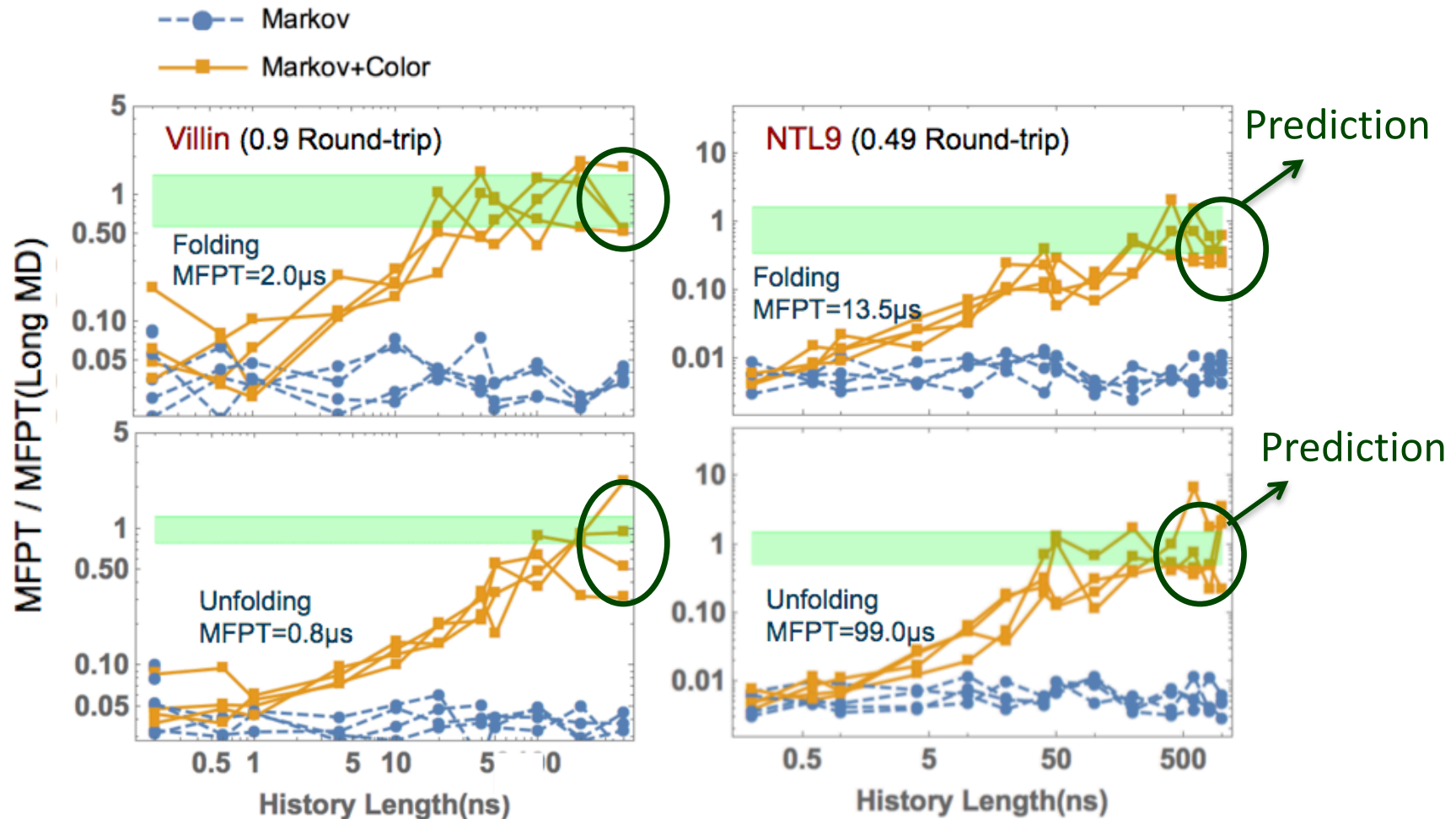
# Reducing the amount of Data (< 5%)

# Non-Markovian Analyses

**Reduced data, MSMBuilder States**

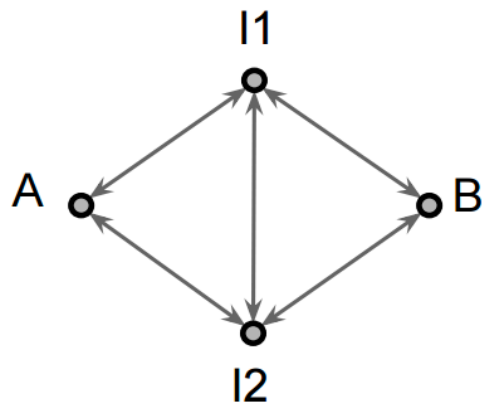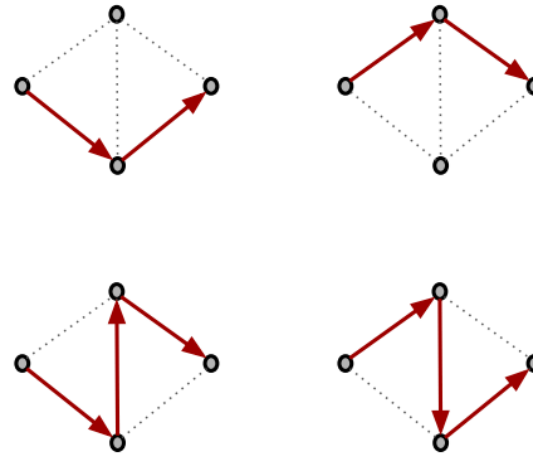# Non-Markovian Analyses

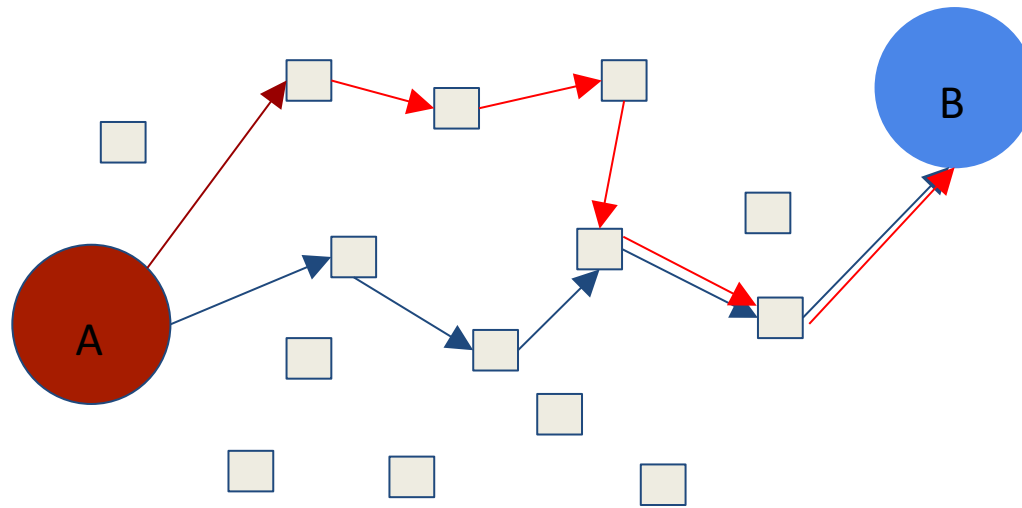**Reduced data, MSMBuilder States**
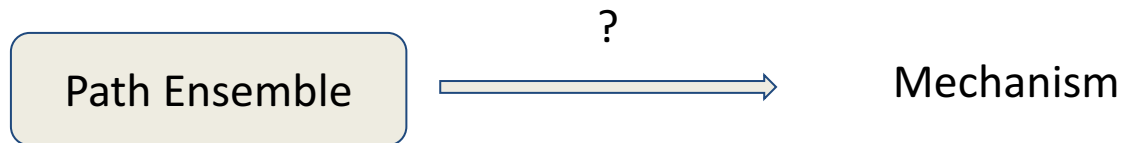
# Mechanism

# Markov vs Non-Markovian

# Mechanism

# Mechanism

# Mechanism

In practice we have…



Path Ensemble $\xrightarrow{\quad ? \quad}$ Mechanism

# Mechanism: Fundamental Sequence



Path Ensemble  ——?——>  Mechanism

# Mechanism: Fundamental Sequence

Defining the "backbone" of the path or *fundamental sequence* will allow us to divide the path ensemble in classes using an equivalence relation. Two paths that share the same fundamental sequence belong to the same class.
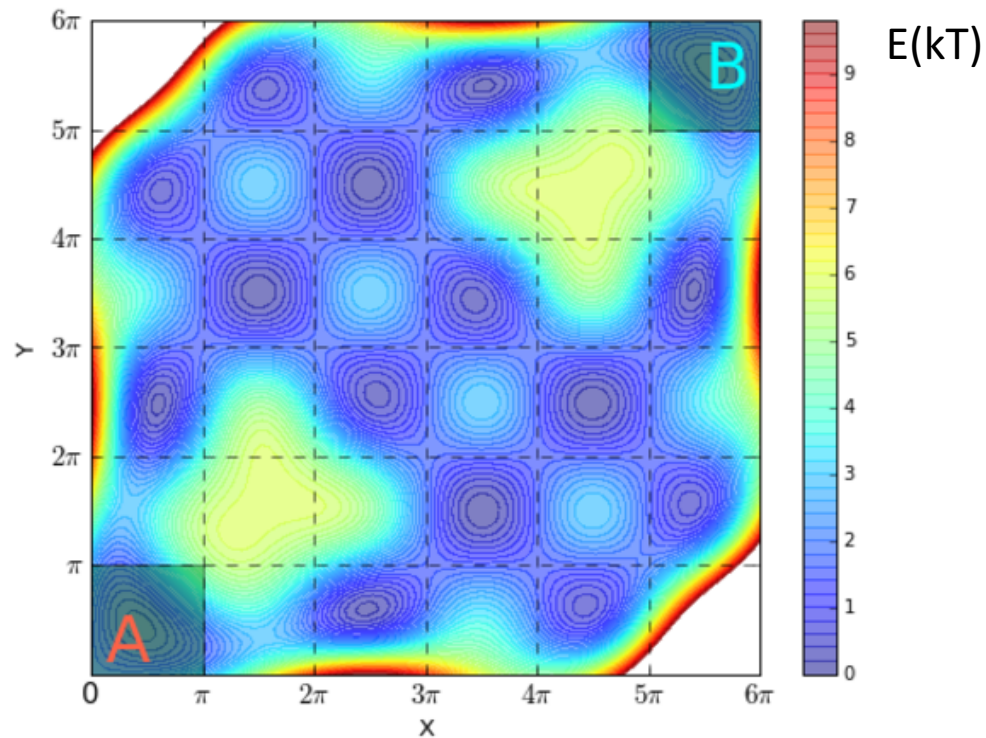
Def.

The fundamental sequence of a path is the most likely sequence that is consistent with the connectivity of the path. The likelihood is maximized in both directions.

$$\mathrm{FS}^* = \underset{\mathbf{q} \in \Gamma(G)}{\arg\min} \left\{ \sum_{i=1}^{|\mathbf{q}|-1} -\log(k_{q_i,q_{i+1}} k_{q_{i+1},q_i}) \right\}$$
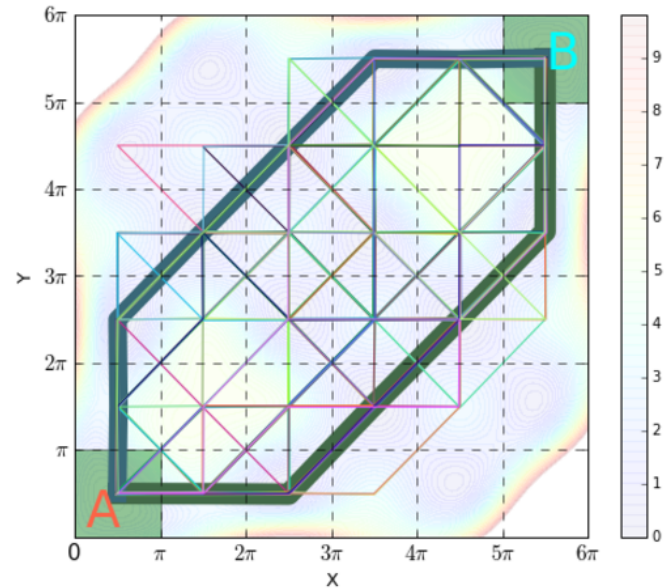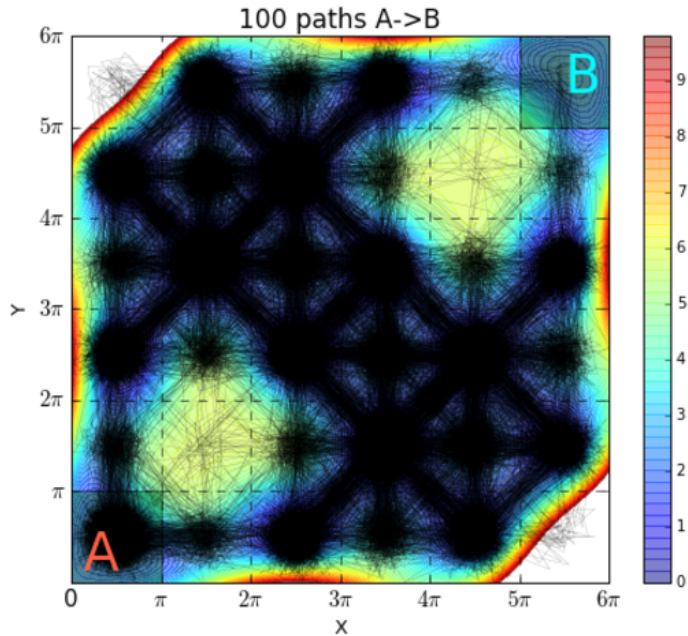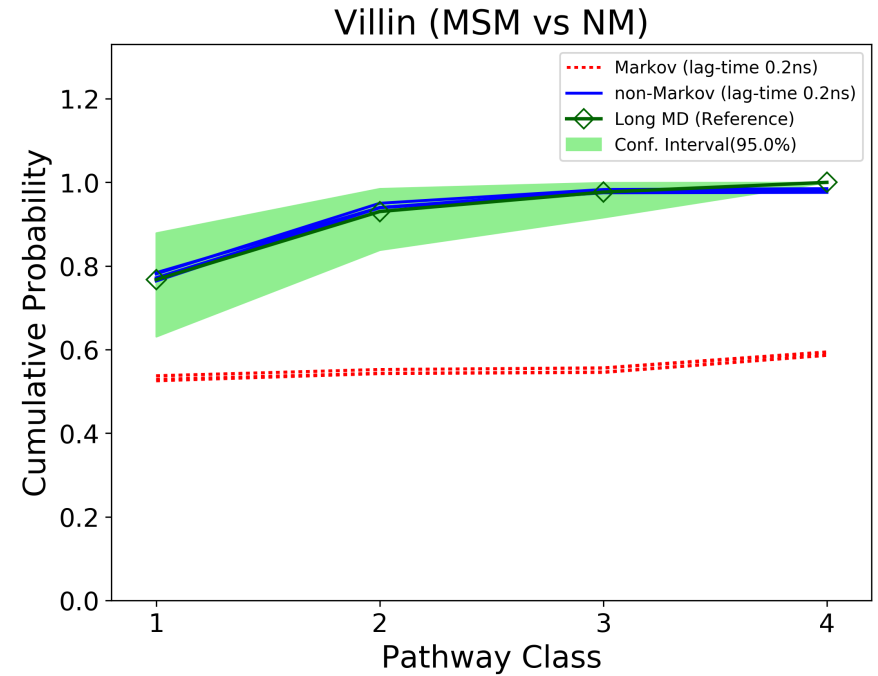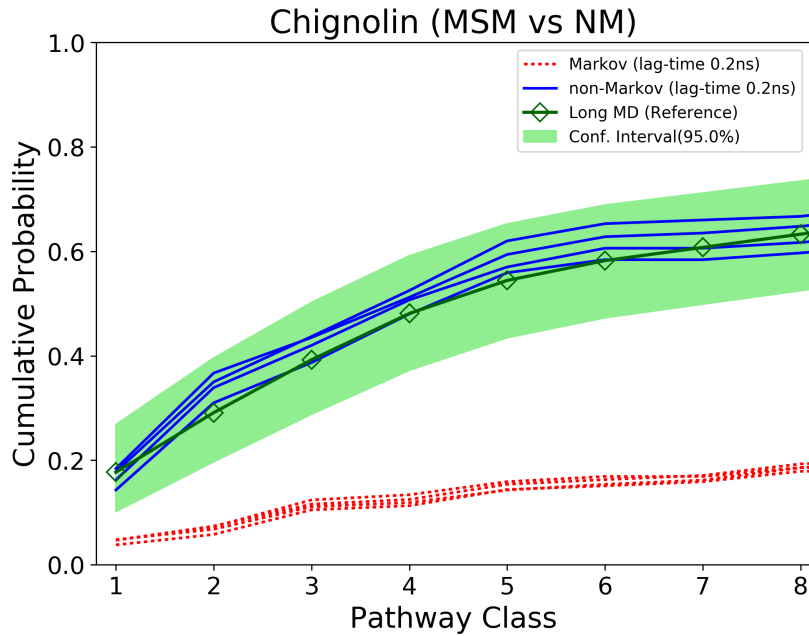
# Mechanism: Fundamental Sequence

Example: 2D toy model

# Mechanism: Fundamental Sequence
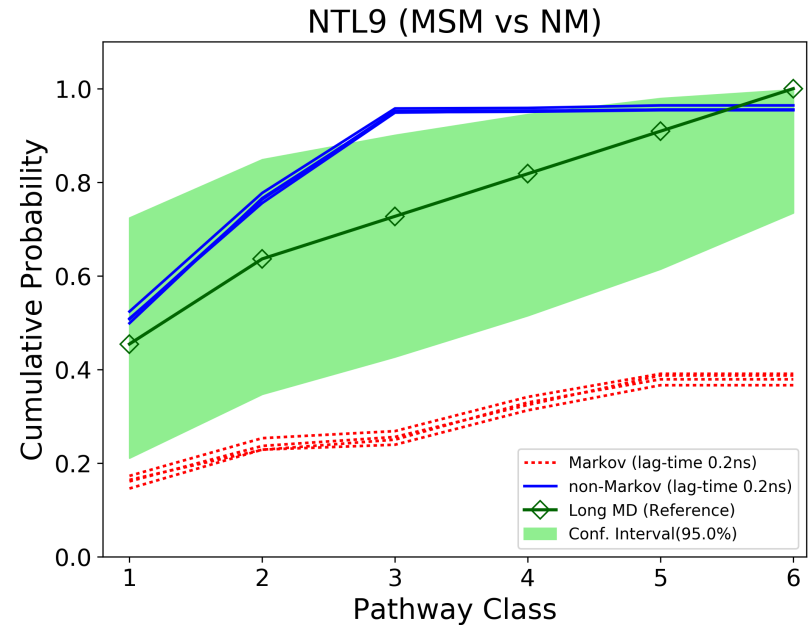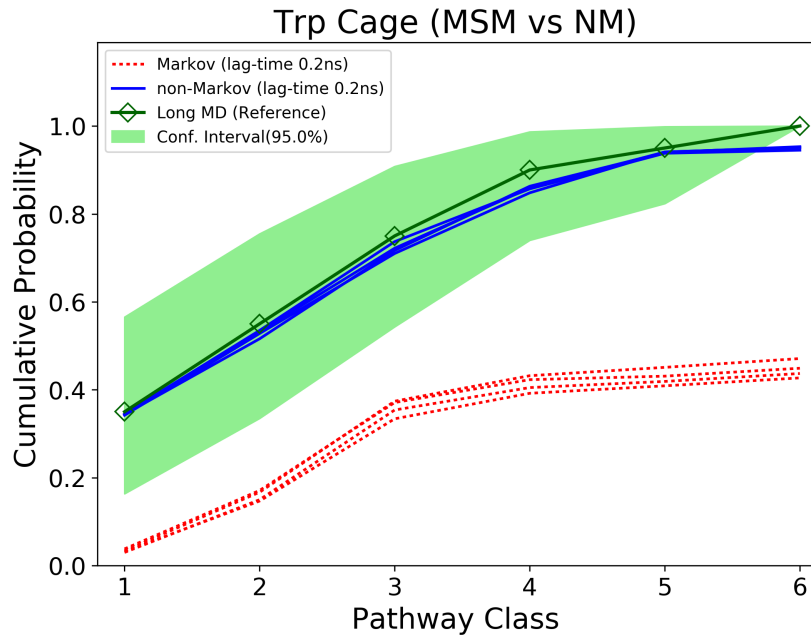
Example: 2D toy model

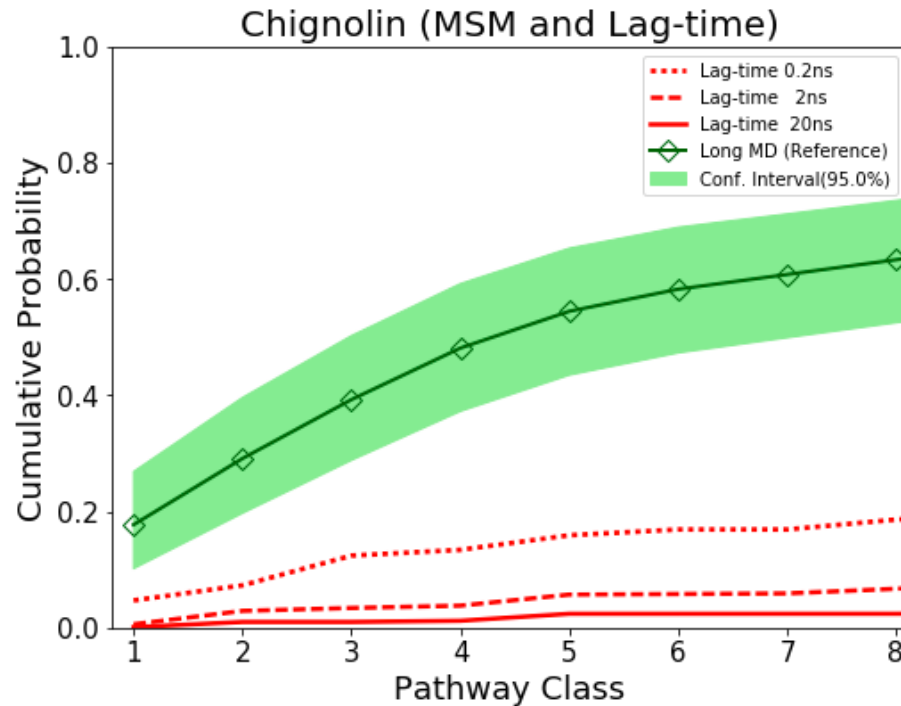# Mechanism: MSM vs NM

## Classification based on the FS

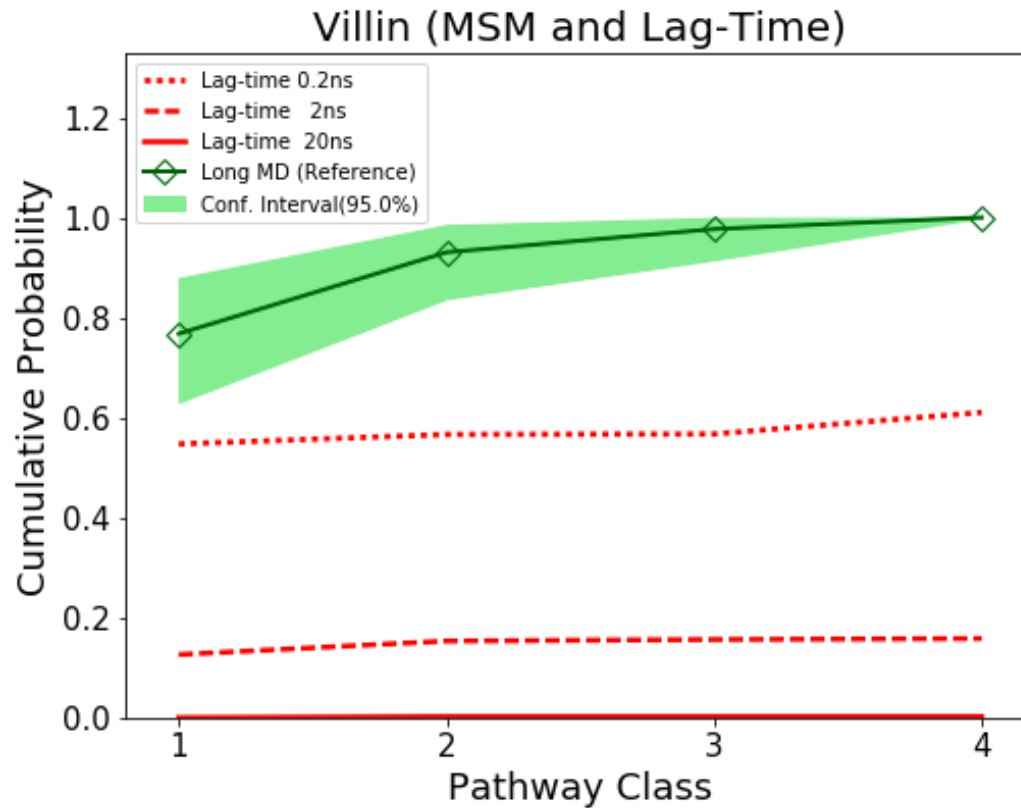# Mechanism: MSM vs NM

## Classification based on the FS

# MSM vs Lag-time

## Classification based on the FS

# MSM vs Lag-time

## Classification based on the FS

# Conclusions

- The inclusion of color information in the analysis allows us to obtain unbiased MFPTs even when the partition of the space in bins is not optimal.

- In a non-Markovian regime, even with a relatively small amount of history (available in most of the MD simulations), we can improve dramatically the estimation of the MFPTs with respect to regular Markov Models.

- We can drastically reduce the amount of data and still obtain reasonable results.

- If the history is taken in to account, there is no need of lag-time "optimization".

- The NM approach drastically outperforms MSM in the description of the mechanism/path ensemble.
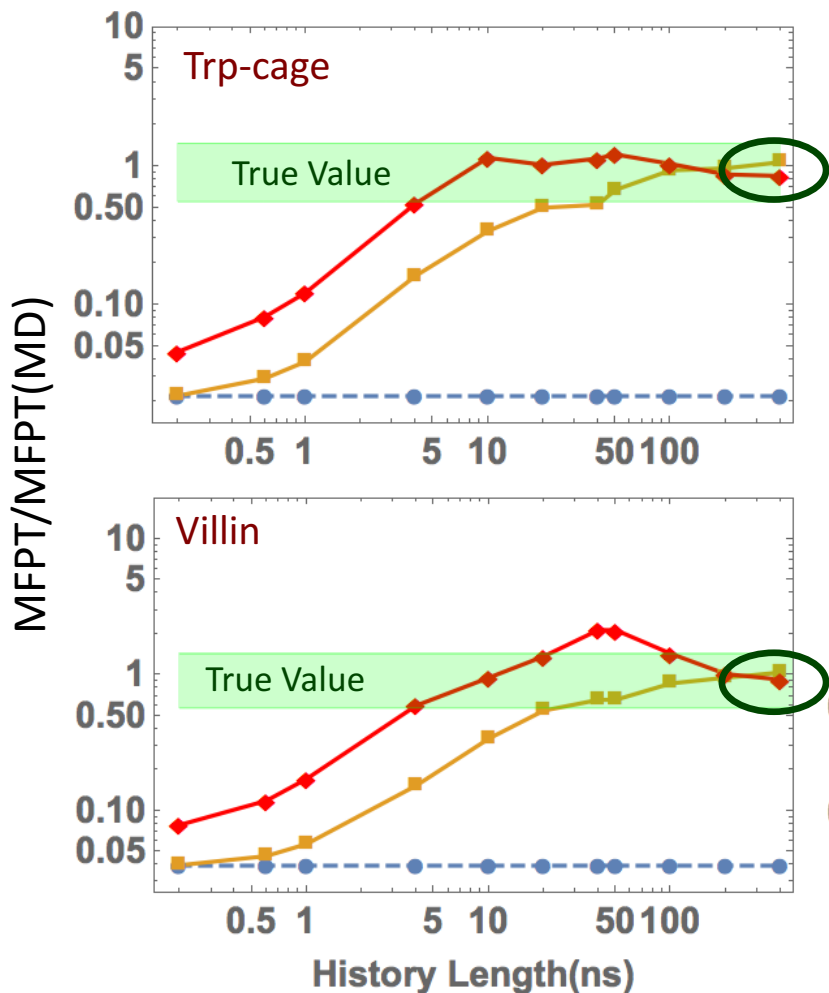
# Acknowledgment

- Joshua Adelman

- Daniel Zuckerman

- Justin Spiriti

- Rory Donovan

- Ramu Anandakrishnan
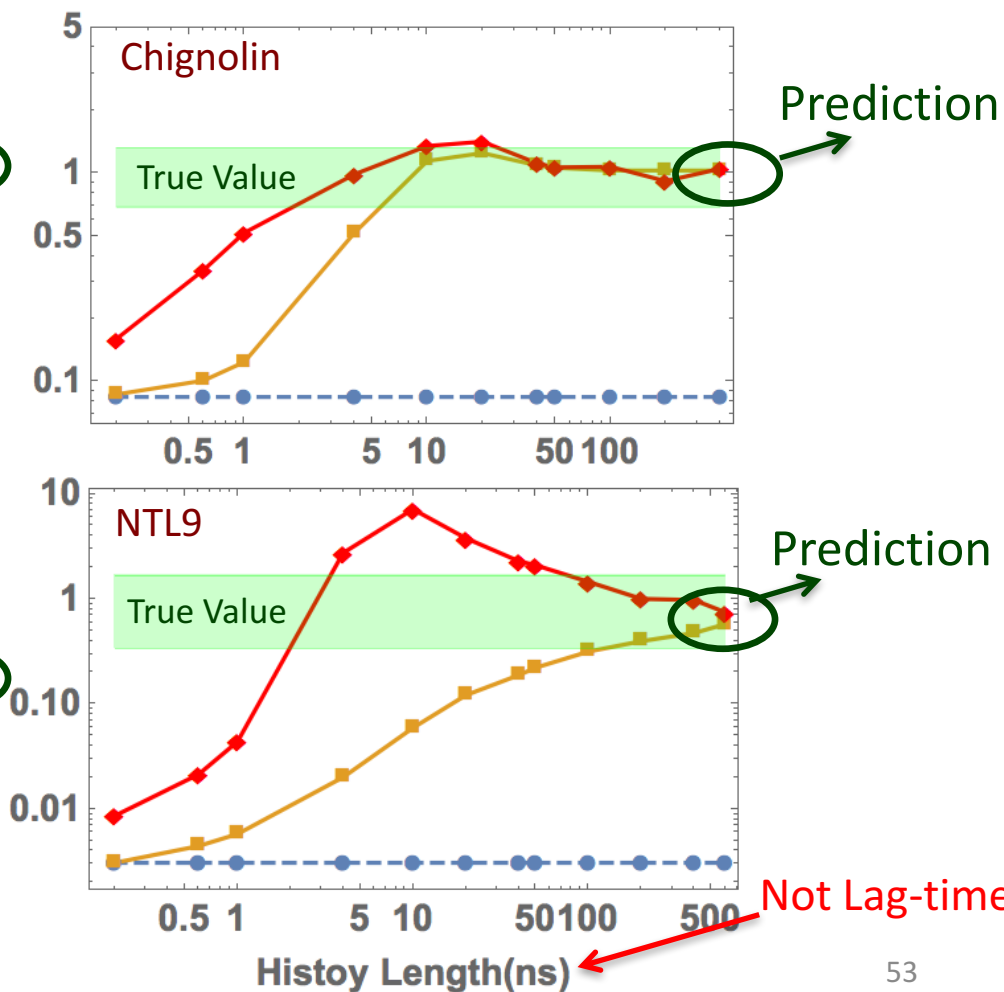
- Ariane Nunes

- Ian Welland

- Shaw group

- NIH

- NSF

# Supporting Information

# Non-Markovian Analyses(Folding)
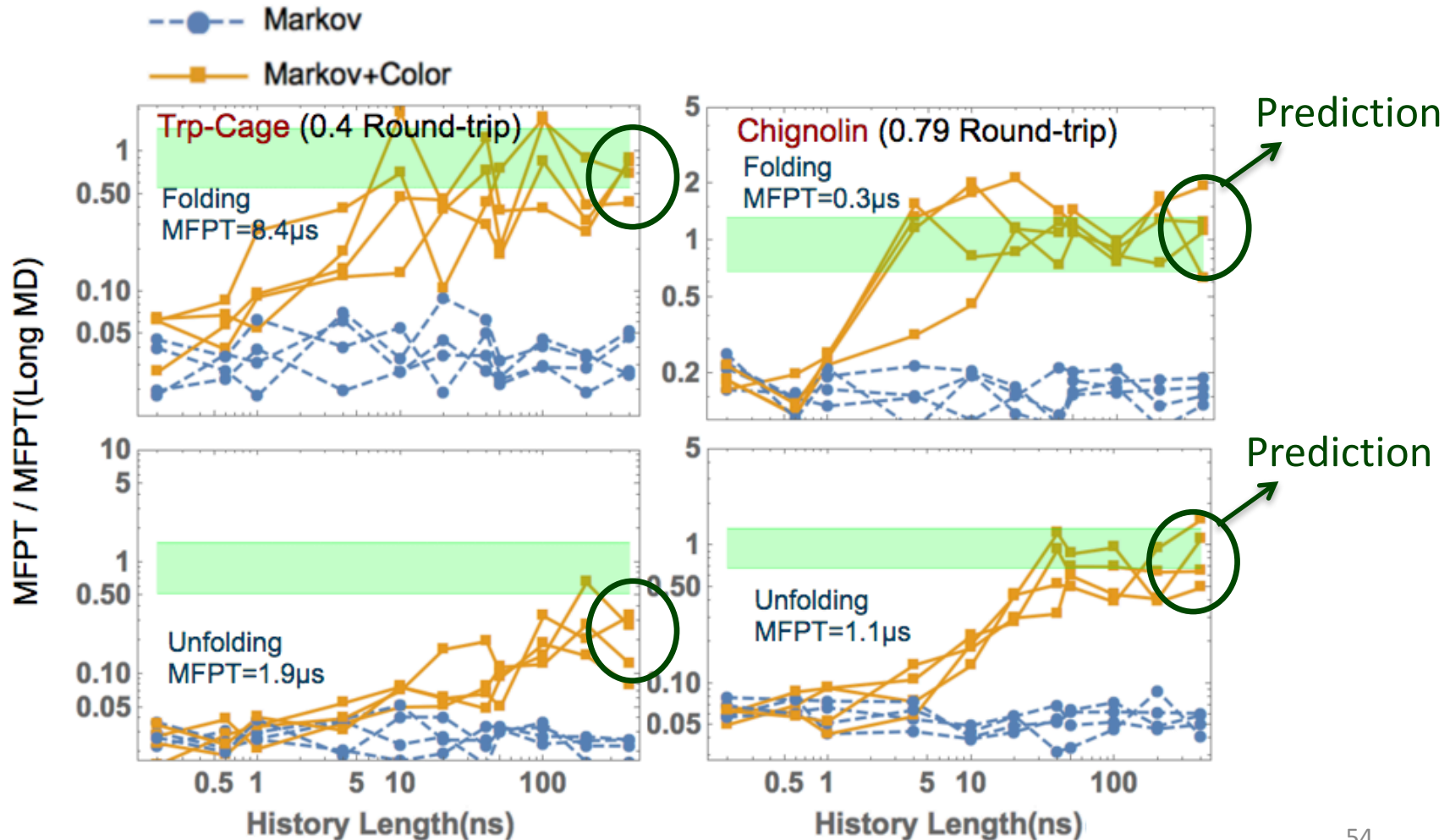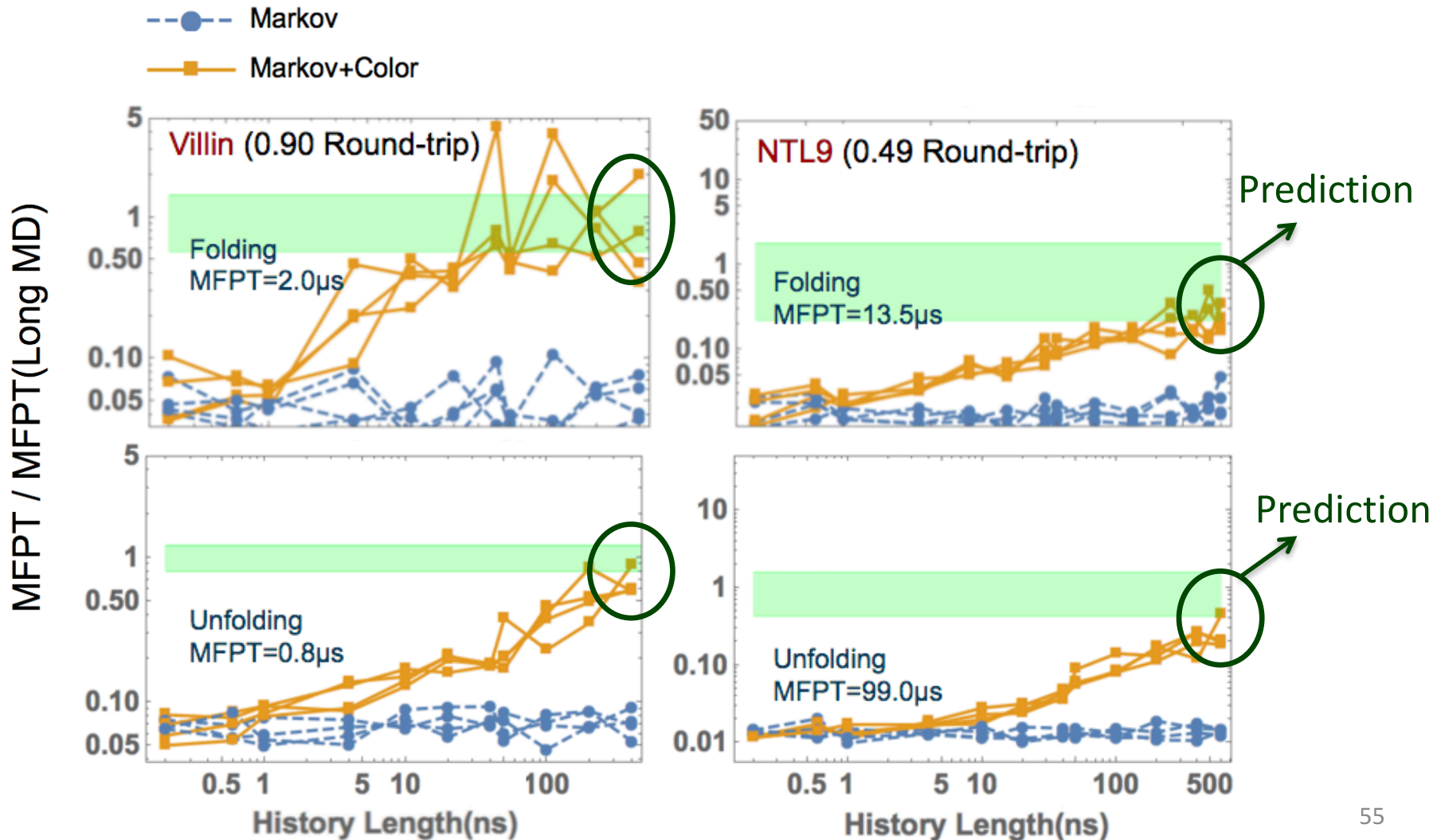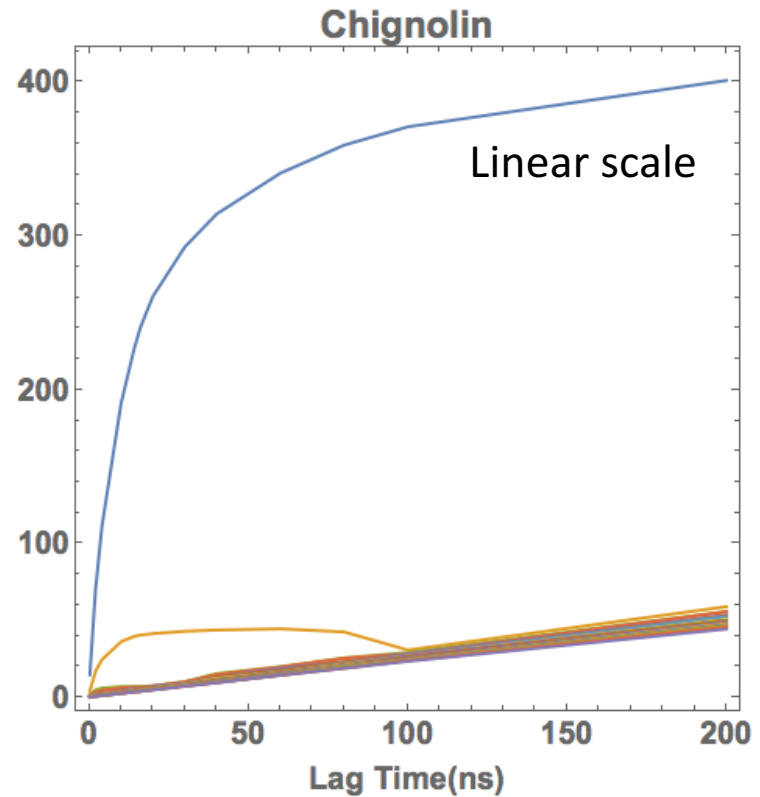
# Non-Markovian Analyses

**Reduced data, RMSD-based States**

# Non-Markovian Analyses
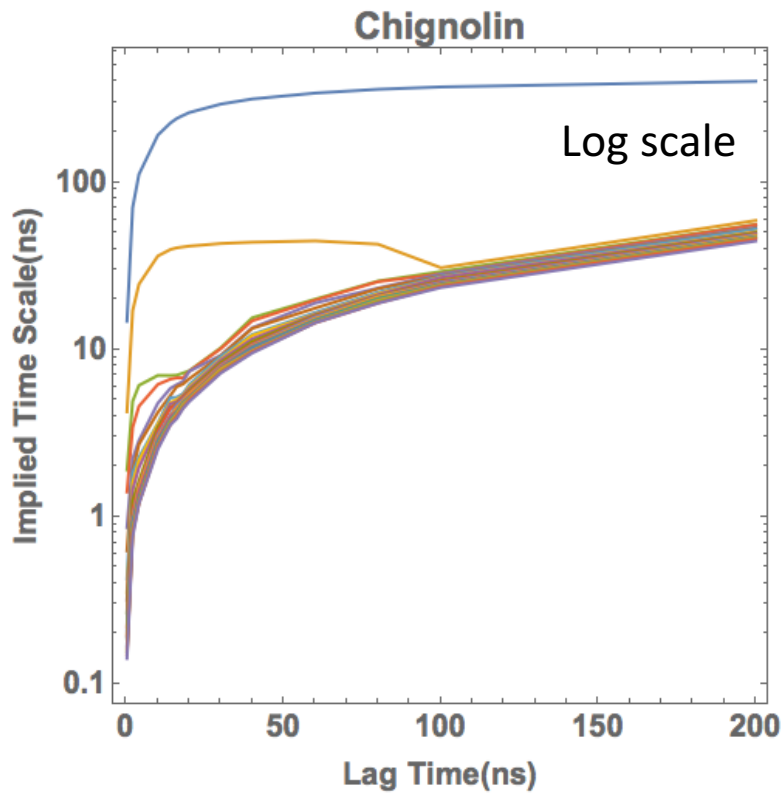
**Reduced data, RMSD-based States**

# MSM: Implied time scales

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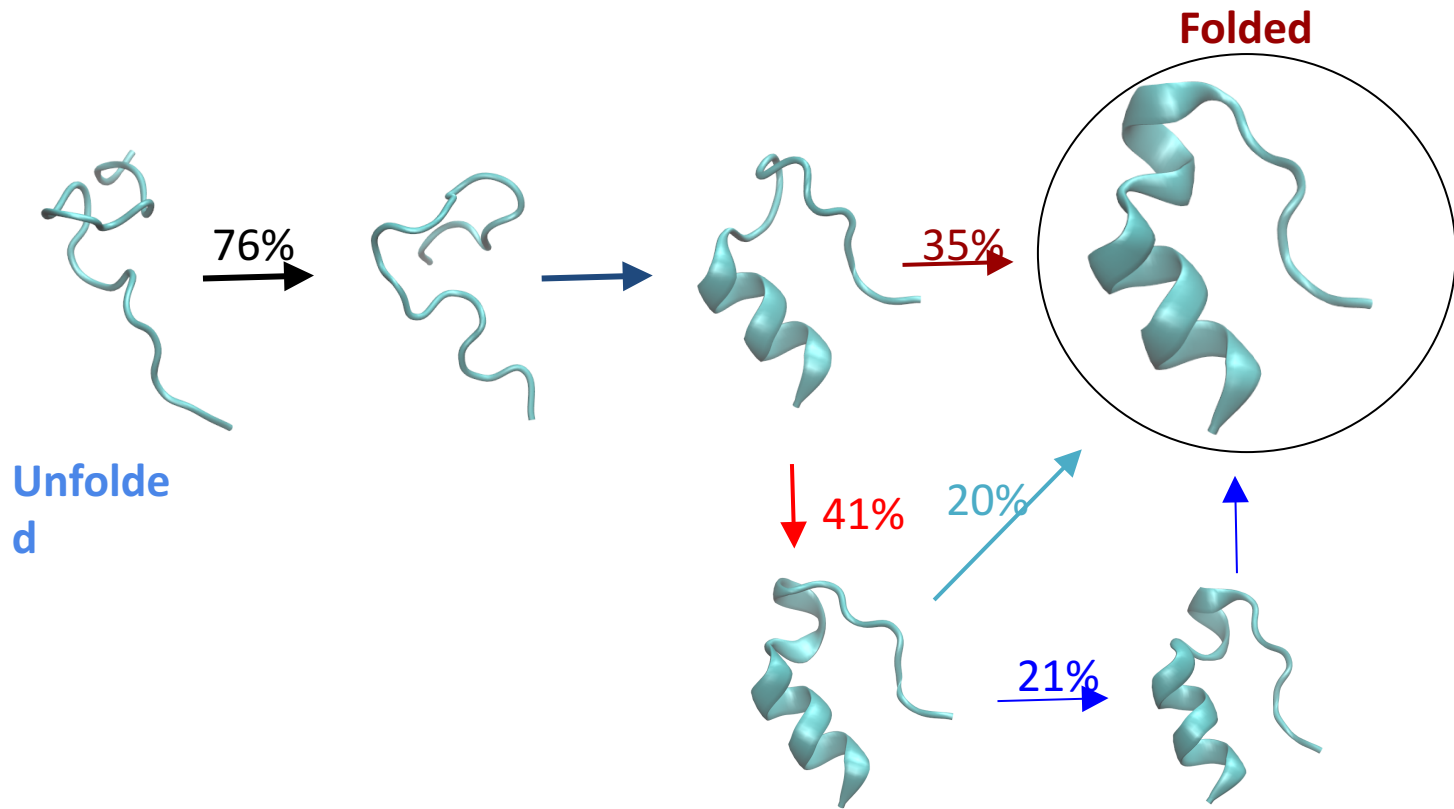


Log scale

Linear scale

# Beyond Markov: Color

$$\mathcal{K}_{ml} = P\{X_{t+\tau} = \left\lceil \frac{l}{2} \right\rceil, L_{t+\tau} = \nu | X_t = \left\lceil \frac{m}{2} \right\rceil, L_t = \mu\}, \ \mu, \nu = \begin{cases} \alpha, \alpha & \text{if } m, l \text{ are odd} \\ \alpha, \beta & \text{if only } m \text{ is odd} \\ \beta, \alpha & \text{if only } m \text{ is even} \\ \beta, \beta & \text{if } m, l \text{ are even} \end{cases}$$
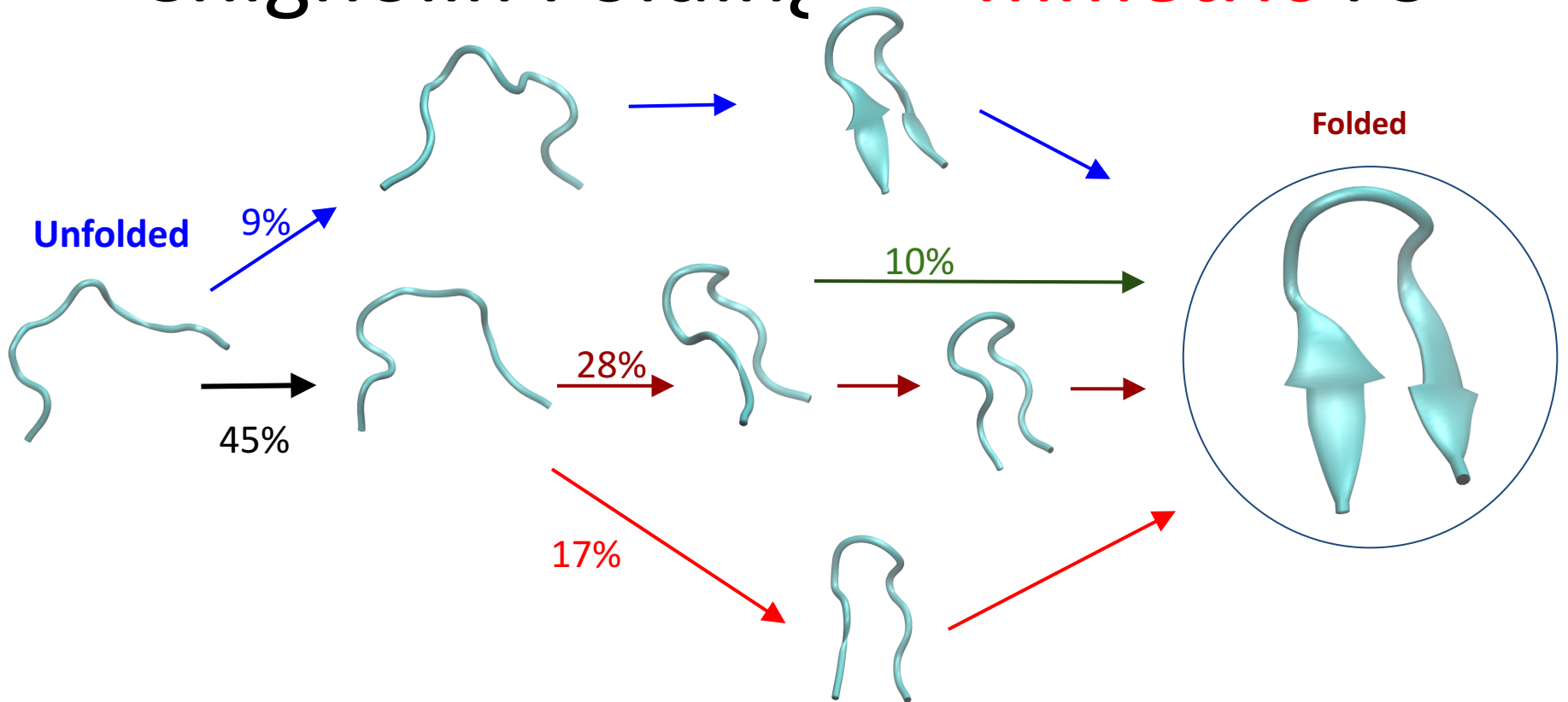
Suarez et al., J. Chem. Theory Comput., 2014, 10 (7), pp 2658–2667
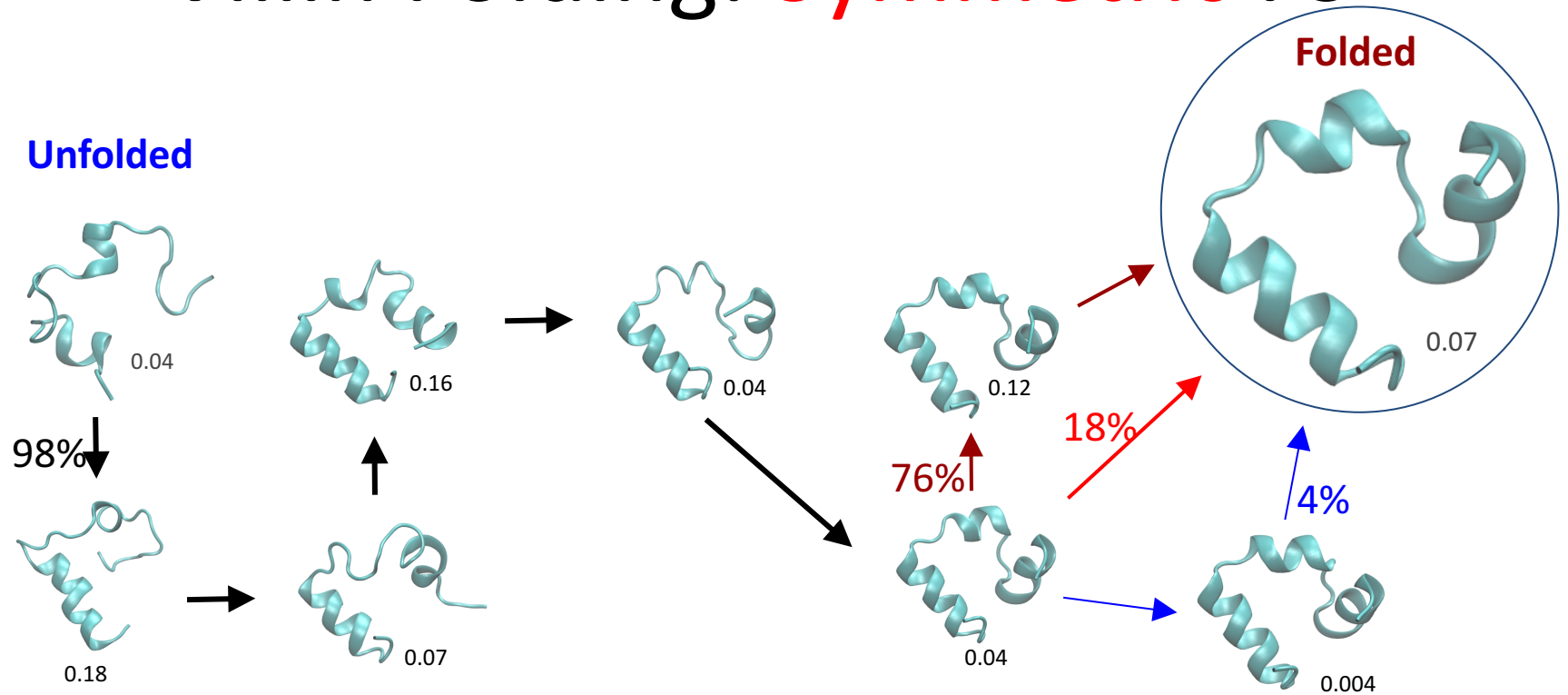Vanden-Eijnden et al., J. Chem. Phys., 2009, 131(4), pp 44120
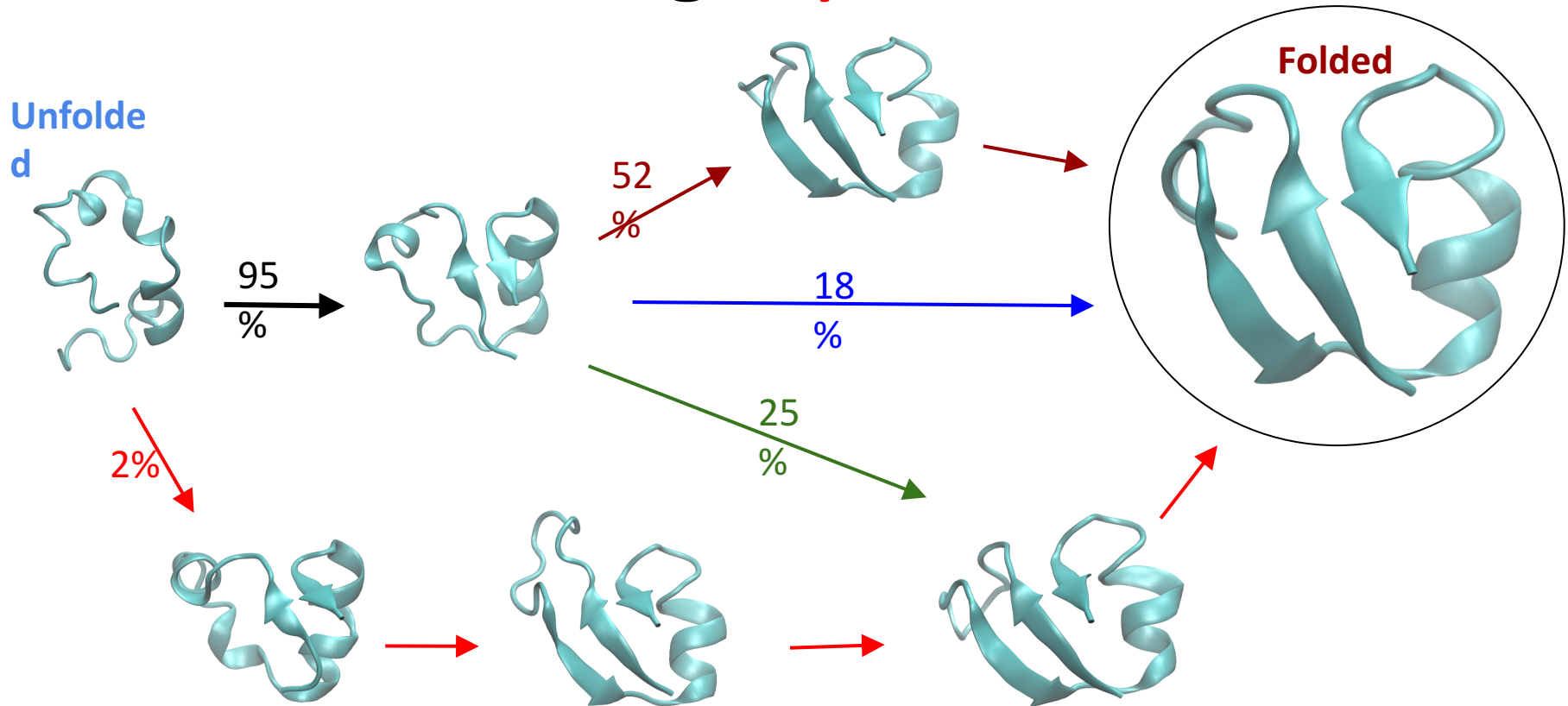
# Trp-cage Folding: Symmetric FS



**Folded**

76%   35%   41%   20%   21%

**Unfolded**

# Chignolin Folding: Symmetric FS



Folded

Unfolded

9%

45%

28%

10%

17%

# Villin Folding: Symmetric FS



**Unfolded**

**Folded**

0.04

0.16

0.04

0.12

0.07

98%

0.18

0.07

76%

18%

0.04

4%

0.004

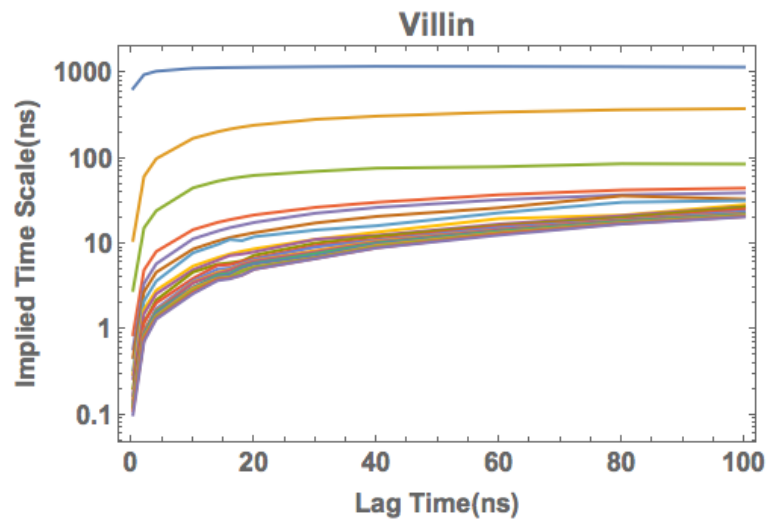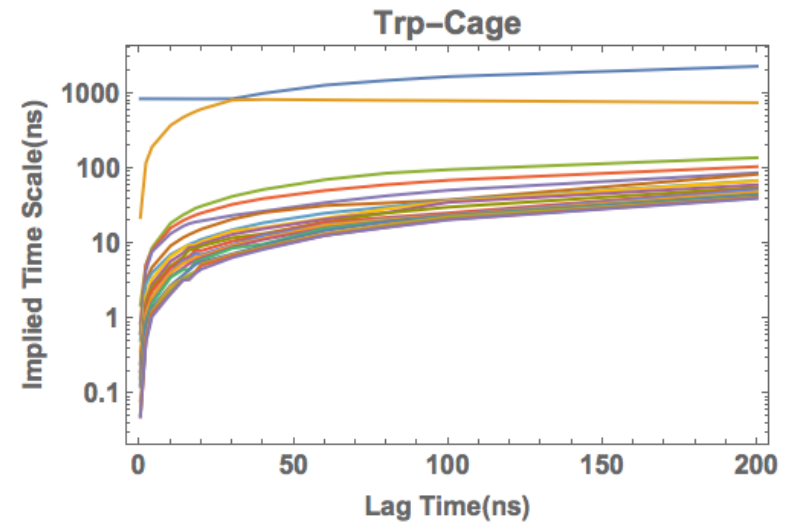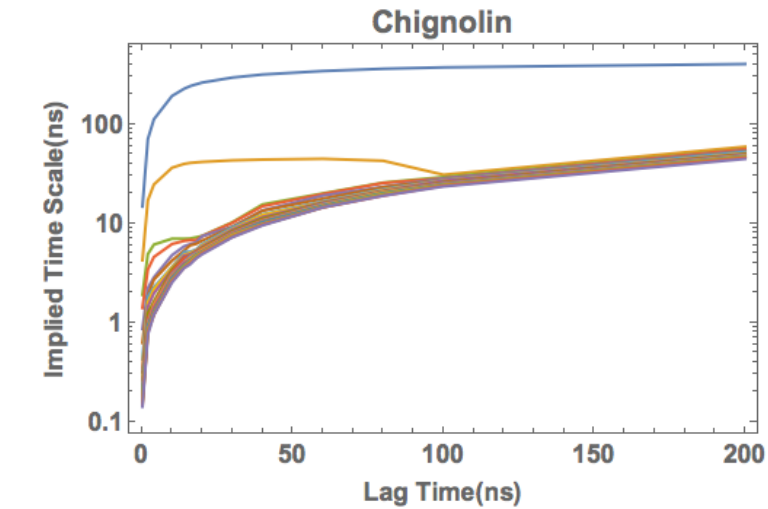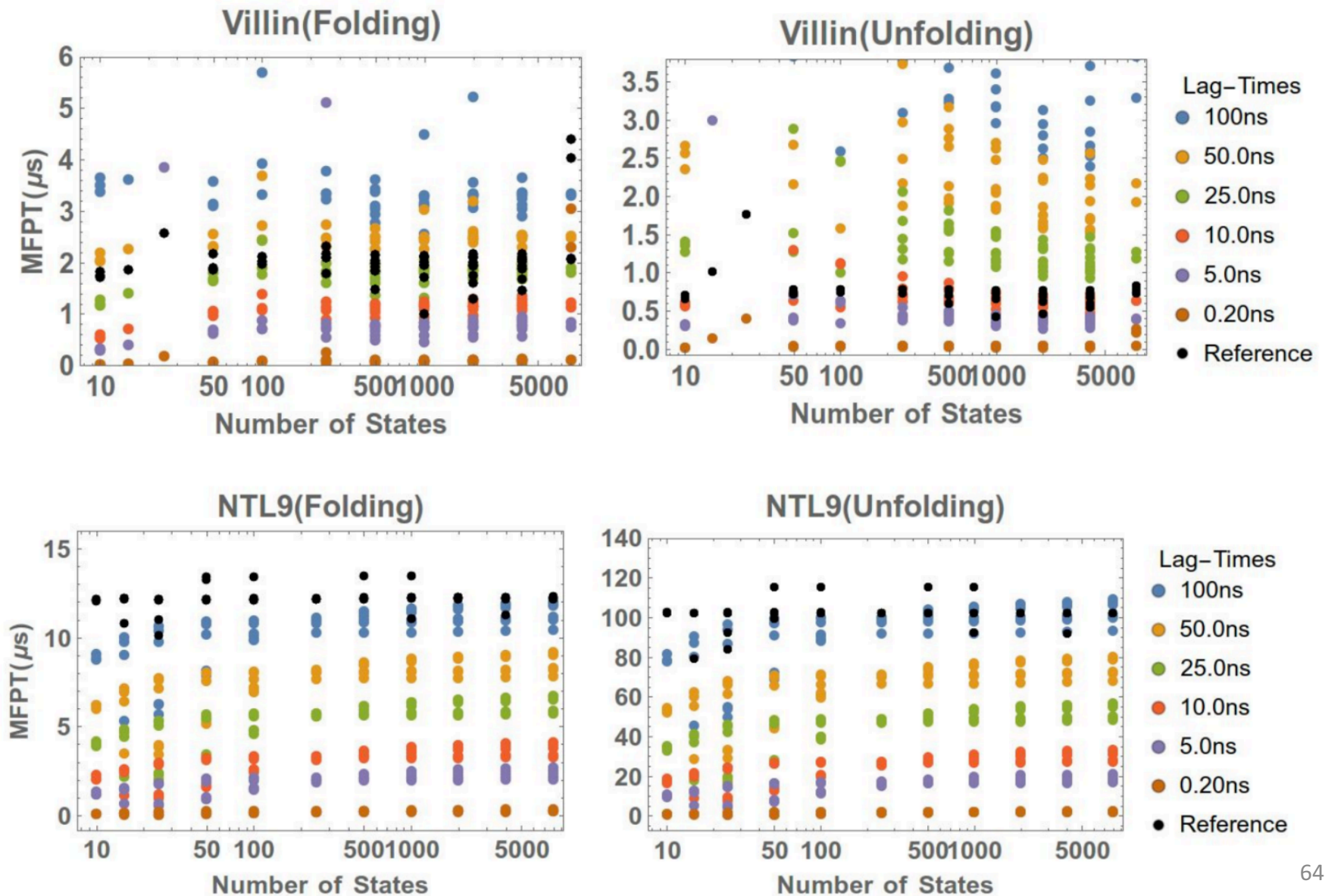# NTL9 Folding: Symmetric FS

# Protein models

Table 1: Protein models used for Markovian and non-Markovian analyses. For each system, the table shows the number of residues, the total simulation time used in the analysis and the state definitions based on heavy-atom RMSD with respect to the folded structure whose protein data bank code (PDB ID) is given in the last column.

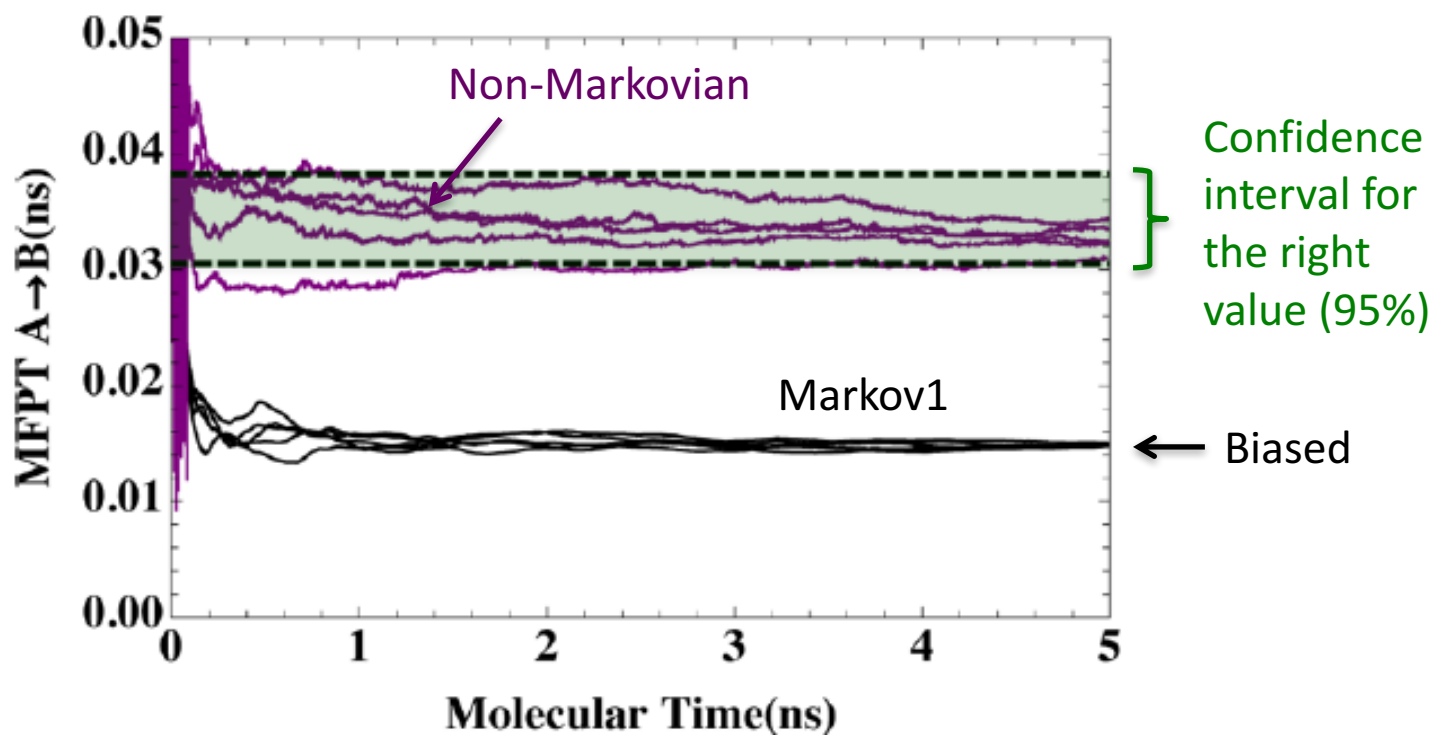| Protein | Num. Residues | Time($\mu s$) | RMSD (Folded) | RMSD (Unfolded) | Reference Structure (PDB ID) |
|---------|------|--------|--------------|----------------|--------------|
| Chignolin | 10 | 106 | $< 1.10$Å | $> 7.00$Å | 5AWL |
| Trp-cage | 20 | 208 | $< 1.75$Å | $> 10.0$Å | 2JOF |
| NTL9 | 39 | 1100 | $< 1.50$Å | $> 10.0$Å | 2HBA |
| Villin | 35 | 125 | $< 1.50$Å | $> 11.0$Å | 2F4K |

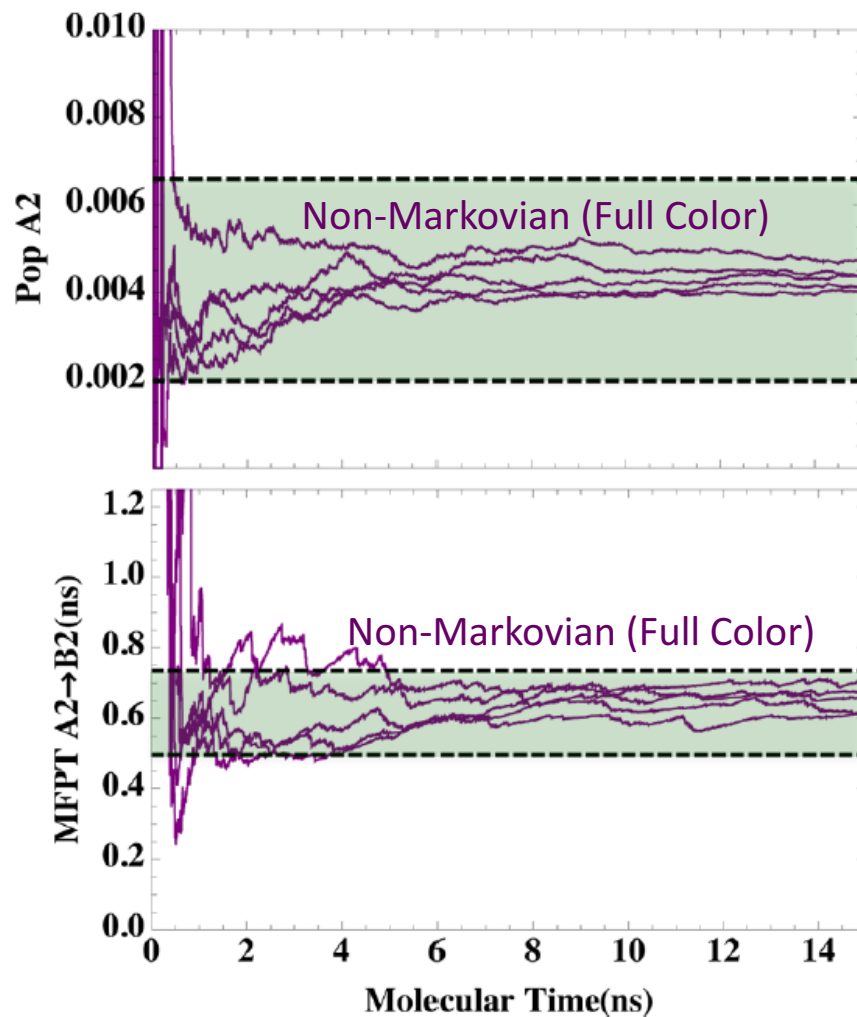# Implied time scales

# Markov State Models

# Example: Methane/Methane

Dissociation process, 5 independent WE simulations.

# Example: Ala4

5 independent WE simulations.

# Ala4

First passage time distribution