*Connecting the dots:*

# A probabilistic model for biomolecular latent space trajectories
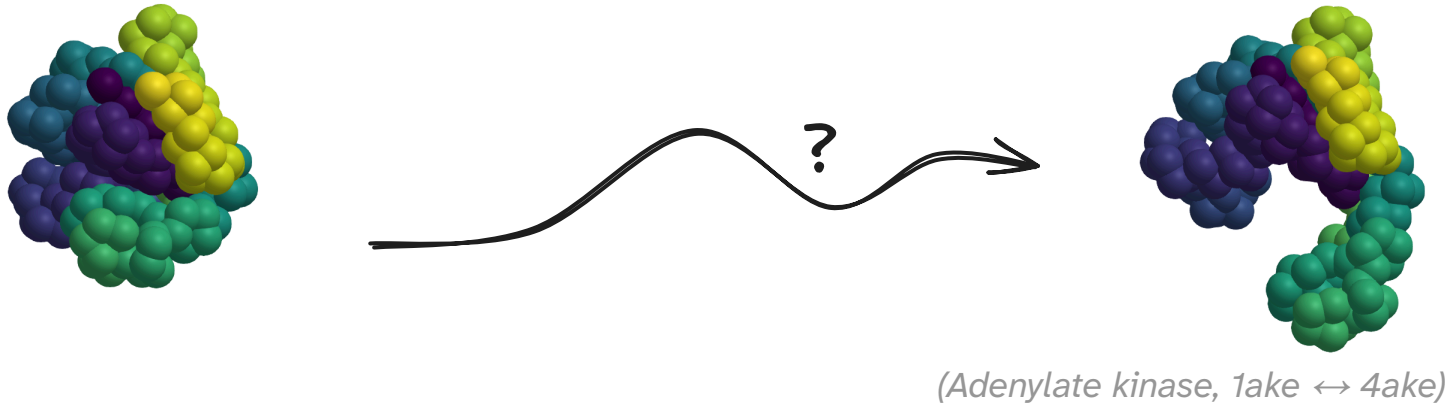
Mittelerde Meeting 2025

## Andreas Kröpelin

Microscopic Image Analysis Group

University Hospital Jena     Interactive Inference     CRC 1456

# Conformational Dynamics of Biomolecules

▶ want to understand *molecular machines*

▶ can observe individual conformations



*(Adenylate kinase, 1ake ↔ 4ake)*

▶ want to find *continuous dynamic* of conformational change
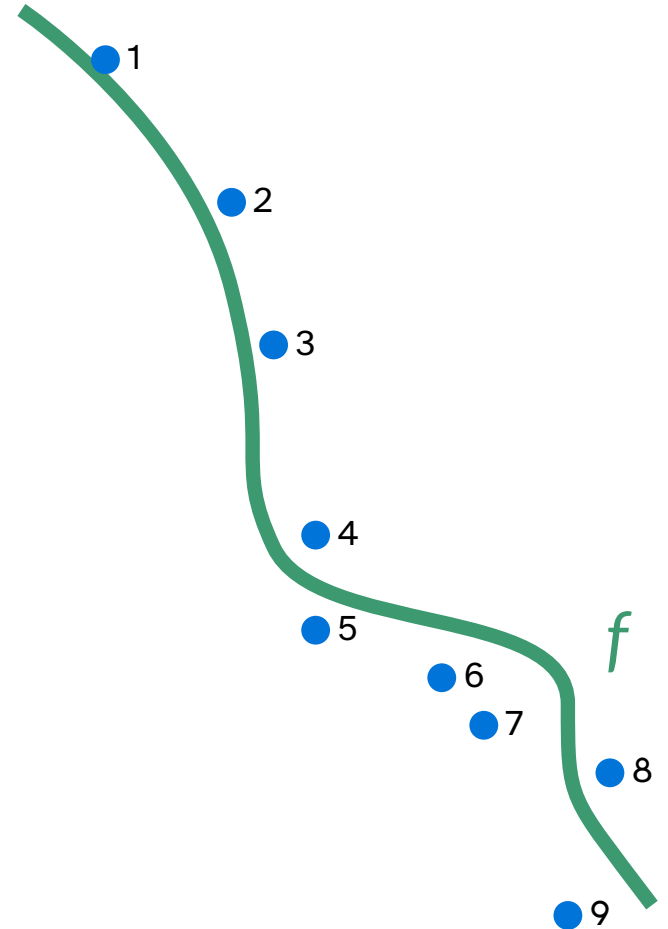(in some *latent space*)

# Mathematical abstraction

we have

observations $\mathcal{Y} = \{y_1, ..., y_m\} \subset \mathbb{R}^d$

and want to explain them with a

curve $f : [0, 1] \rightarrow \mathbb{R}^d$

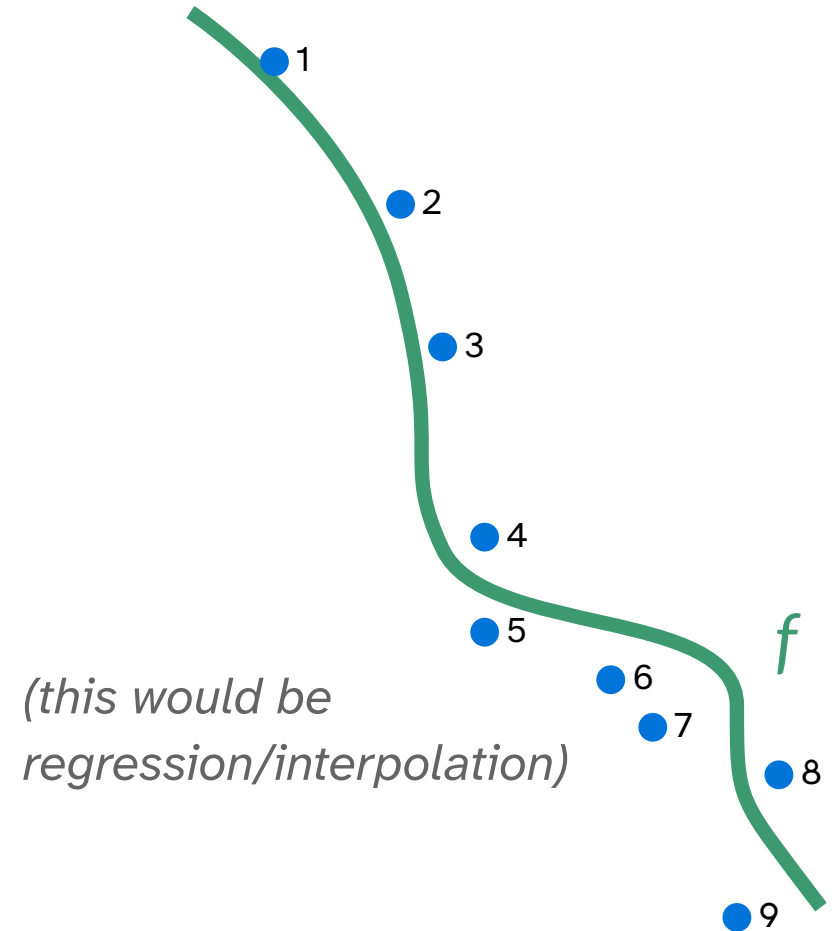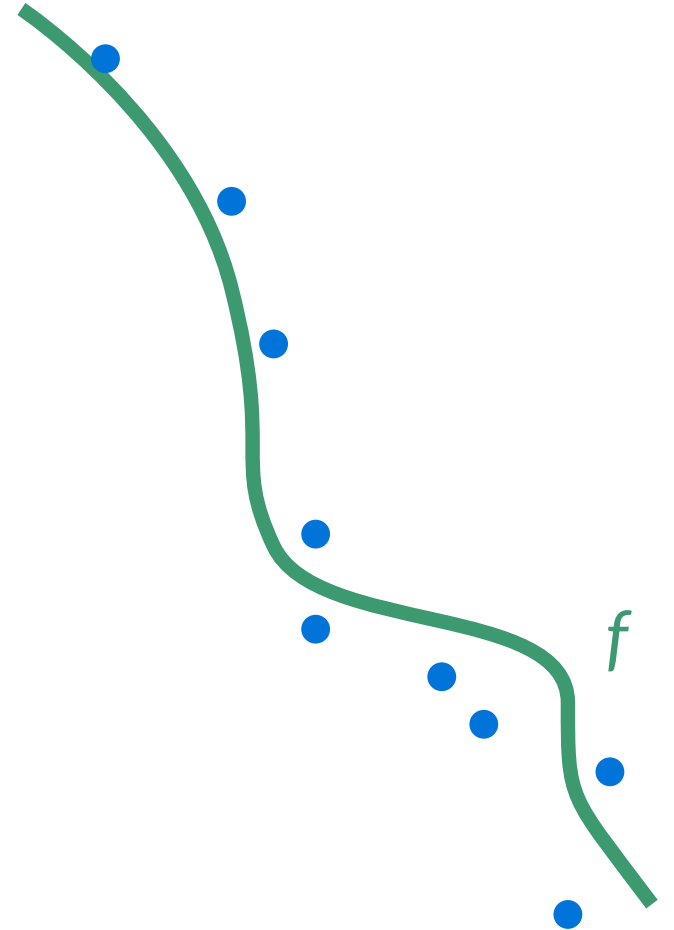*by the magic of* ✨ *Bayes* ✨

# Mathematical abstraction

we have

observations $\mathcal{Y} = \{y_1, ..., y_m\} \subset \mathbb{R}^d$

and want to explain them with a

curve $f : [0, 1] \rightarrow \mathbb{R}^d$

*by the magic of* ✨ *Bayes* ✨

*(this would be regression/interpolation)*

# Mathematical abstraction

we have

**unordered** observations $\mathcal{Y} = \{y_1, ..., y_m\} \subset \mathbb{R}^d$

and want to explain them with a
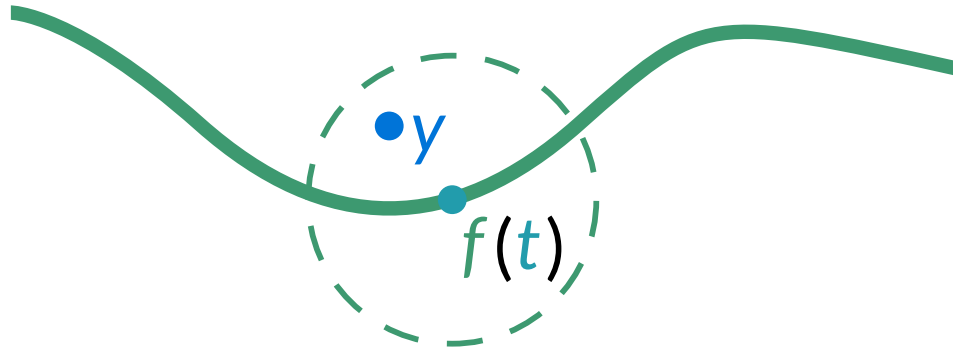
curve $f : [0, 1] \rightarrow \mathbb{R}^d$

*by the magic of* ✨ *Bayes* ✨

$f$

# Data generating process: Curve Mixture

Curve mixture $\mathcal{M}$ consists of curve $f$ and observation noise $\sigma$. Observation $y$ is generated from $\mathcal{M}$ like this:

1. draw $t \sim \mathcal{U}[0, 1]$

2. draw $y \sim \mathcal{N}\left(f(t), \sigma^2 I\right)$



**likelihood:** $\quad p(y \mid \mathcal{M}) = \dfrac{1}{Z(\sigma)} \displaystyle\int_0^1 \exp\left(-\dfrac{\|y - f(t)\|^2}{2\sigma^2}\right) \mathrm{d}t$

# Smooth curves

We pefer ⟨curve⟩ over ⟨curve⟩

... corresponding to preferring a low **bending energy**.

$$\textbf{prior:} \quad p(\mathcal{M}) = \exp\left(-\tau\left(\int_0^1 \|f''(t)\|^2 \, \mathrm{d}t\right) \Big/ \left(\int_0^1 \|f'(t)\| \, \mathrm{d}t\right)^2\right)$$

# Posterior: Likelihood and Prior combined

*Bayes' theorem:*

$$\underbrace{p(\mathcal{M} \mid \mathcal{Y})}_{\text{posterior}} = \prod_{y \in \mathcal{Y}} \underbrace{p(y \mid \mathcal{M})}_{\text{likelihood}} \cdot \underbrace{p(\mathcal{M})}_{\text{prior}}$$

**posterior:** How much do we believe $\mathcal{M}$ given $\mathcal{Y}$?

**likelihood:** from the data generating process

**prior:** for the smoothness

# Discrete approximation

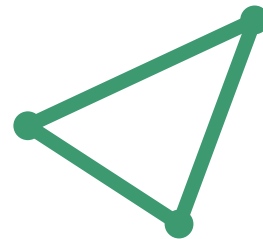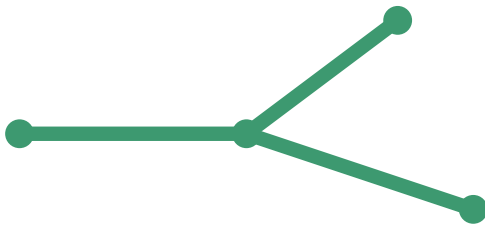Instead of arbitrary curves $f$, consider concatenations of **line segments**:



Model $\mathcal{M}$ now consists of
▶ observation noise $\sigma > 0$
▶ nodes $x_1, ..., x_n \in \mathbb{R}^d$
▶ connectivity information (which two nodes form a segment?)

# Properties of discrete model

👍 ▸ clear what the parameters are: $x_1, ..., x_n \in \mathbb{R}^d$ and $\sigma \in \mathbb{R}$

▸ can compute integrals in likelihood and prior exactly

▸ $p\left( \text{───} \mid \mathcal{Y} \right) = p\left( \text{───} \mid \mathcal{Y} \right)$
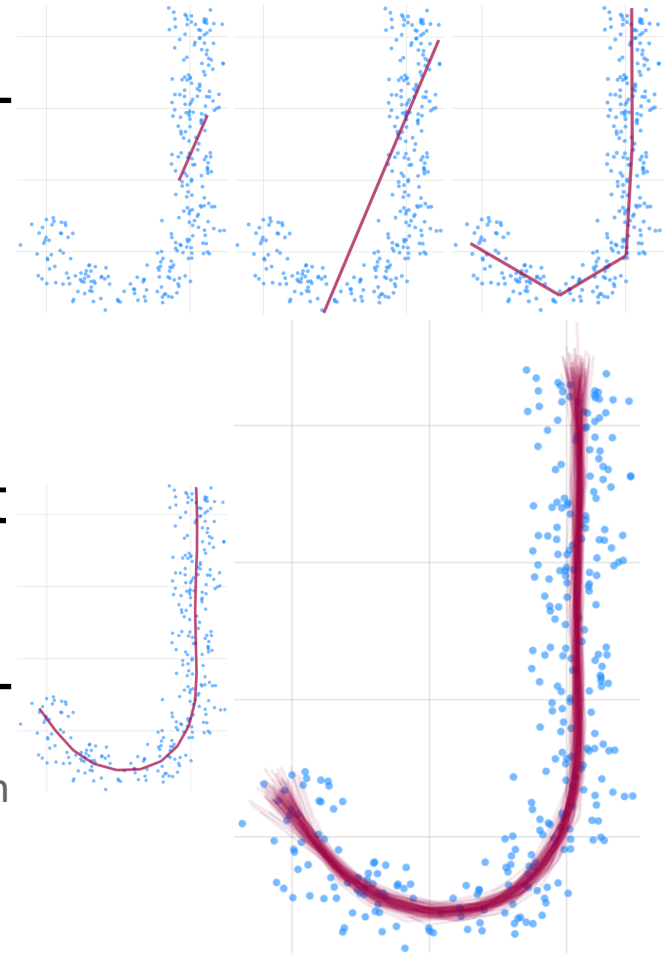
▸ easily generalizable to arbitrary topologies *(mixture of curve mixtures)*

👎 ▸ have to think about number of segments
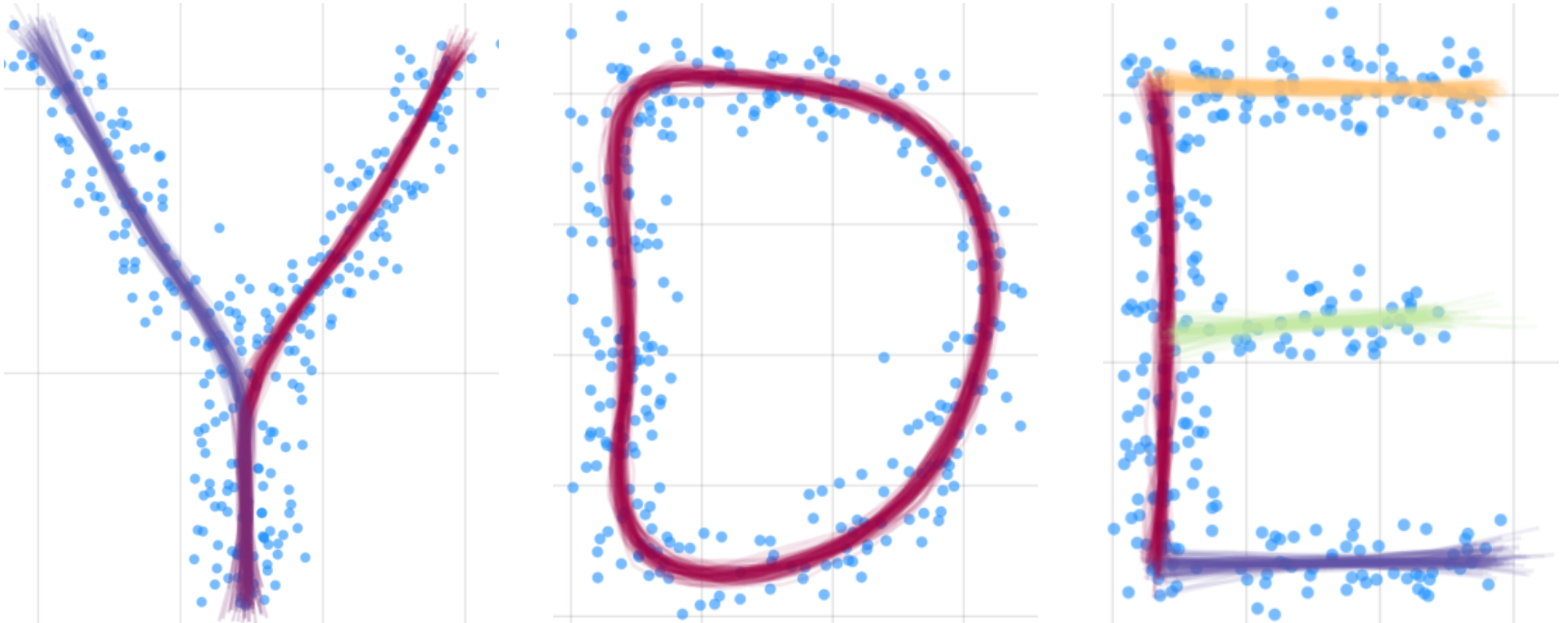
# Estimating a Curve Mixture

1  $\mathcal{M}$ ← single segment
2  **until satisfied**
3  ⌐ maximize posterior density (ADAM)
4  ⌐ finegrain $\mathcal{M}$, i.e. split every segment
5  sample $\mathcal{M}$ from posterior (NUTS)

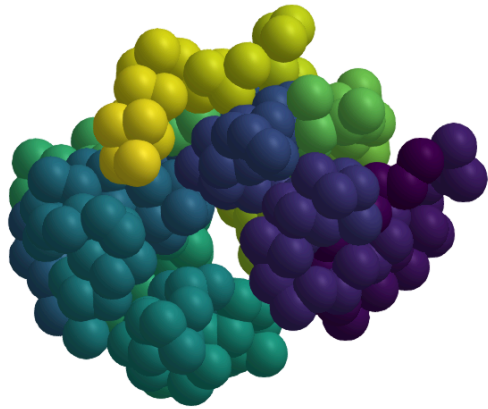(obtain gradient for ADAM and NUTS via automatic differentiation with `Mooncake.jl`)

# More letters!

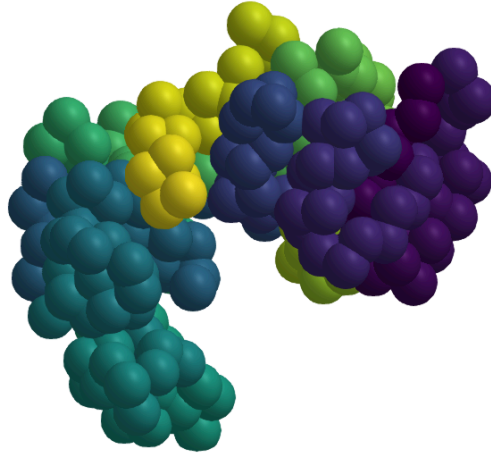(showcasing other topologies with a generalized model)
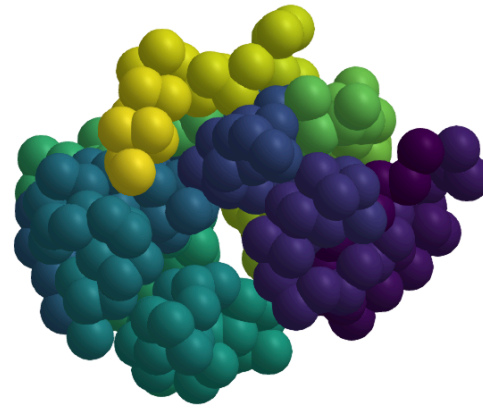
# Demo: Human Serum Transferrin (1a8e)

84 AA-chains from PDB with sequence similarity ≥ 90 % to 1a8e
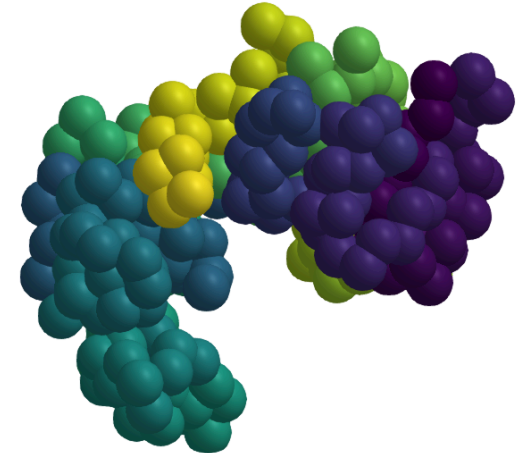


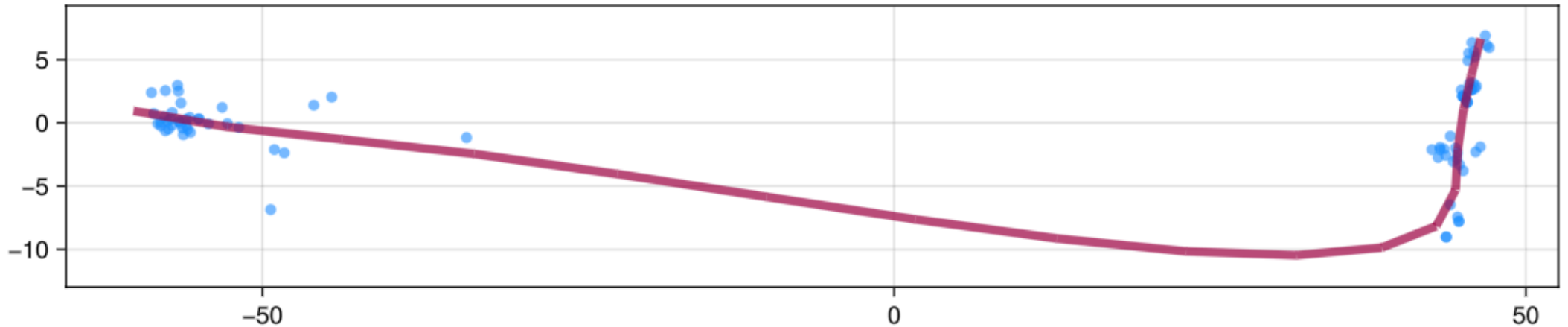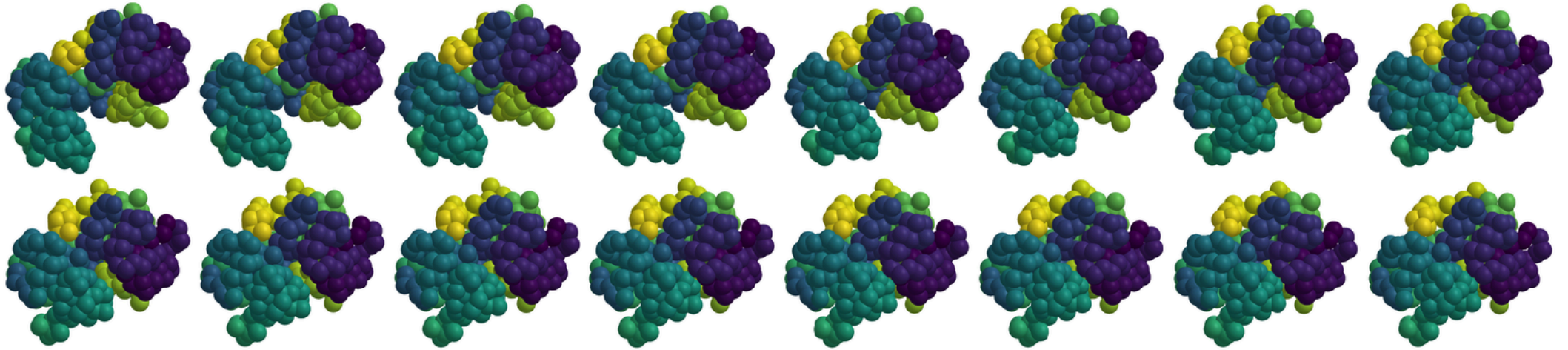1a8e                2h52                2o7u                5dhy

# Demo: Human Serum Transferrin (1a8e)

Curve Mixture in 2D latent space (PCA)
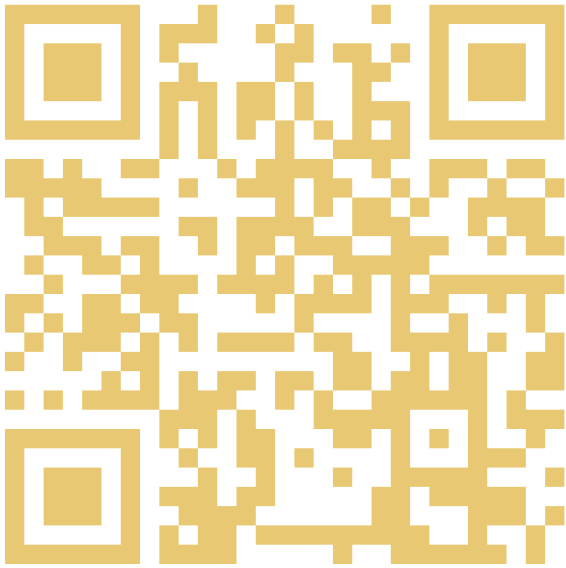
# Demo: Human Serum Transferrin (1a8e)

Traverse curve and project back into "structure space"

# **Summary and next steps**

▶ curve mixtures are a nice tool to find trajectories in data points

   ⤳   predictions biologically meaningful?

   ⤳   design user interface

▶ works reasonably well for PCA-embedded protein structures

   ⤳   other latent embeddings?

   ⤳   apply to RNA-seq data

▶ implementation in Julia let me focus on domain problem and performance (gradient, sampling, and optimization "for free")

# Let's discuss!

✉ andreas.kroepelin@uni-jena.de

🌐 a.ndreas.dev

🐘 @andreask@bayes.club

⇧slides and code on
GitHub