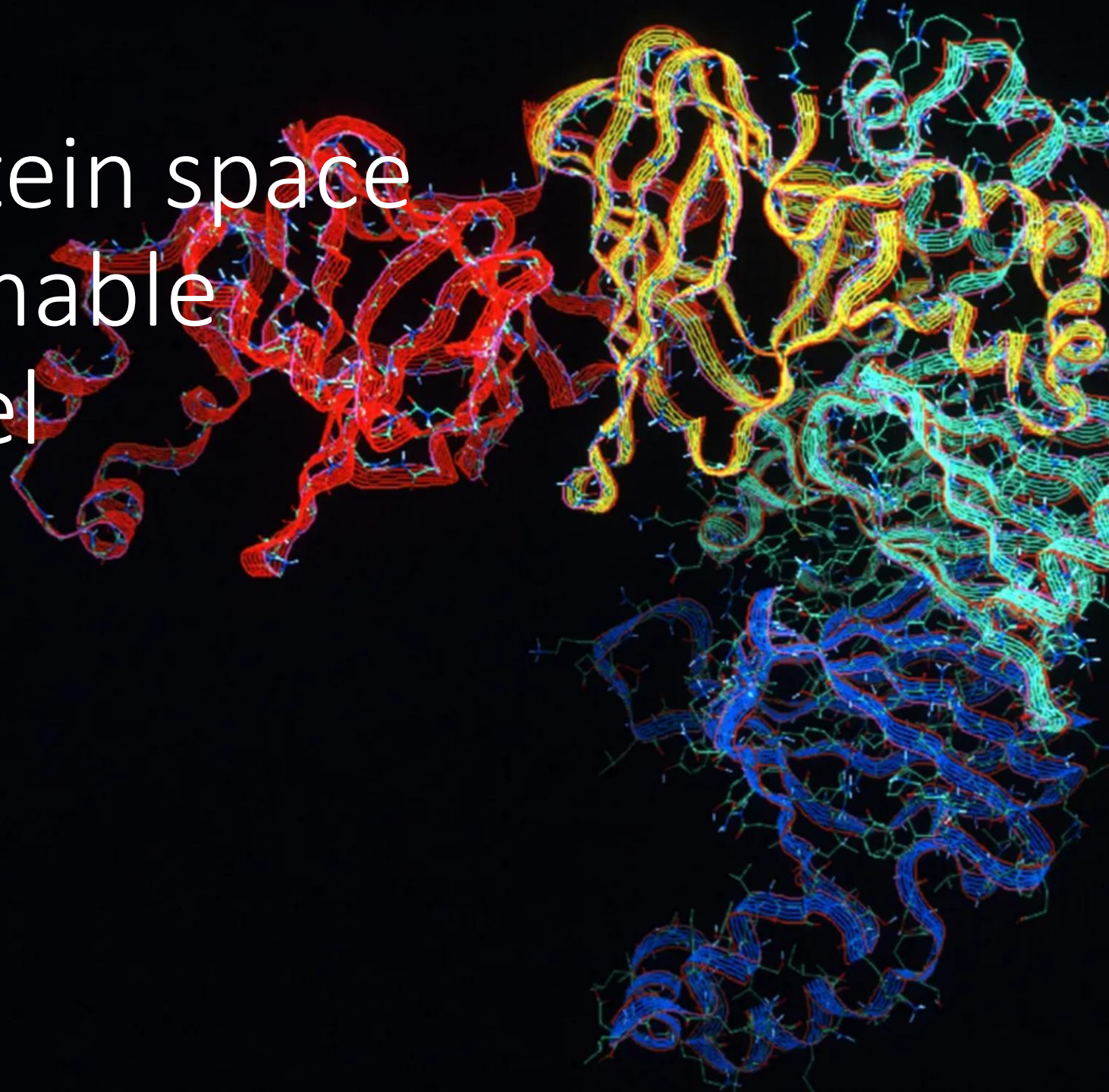# Illuminating protein space with a programmable generative model
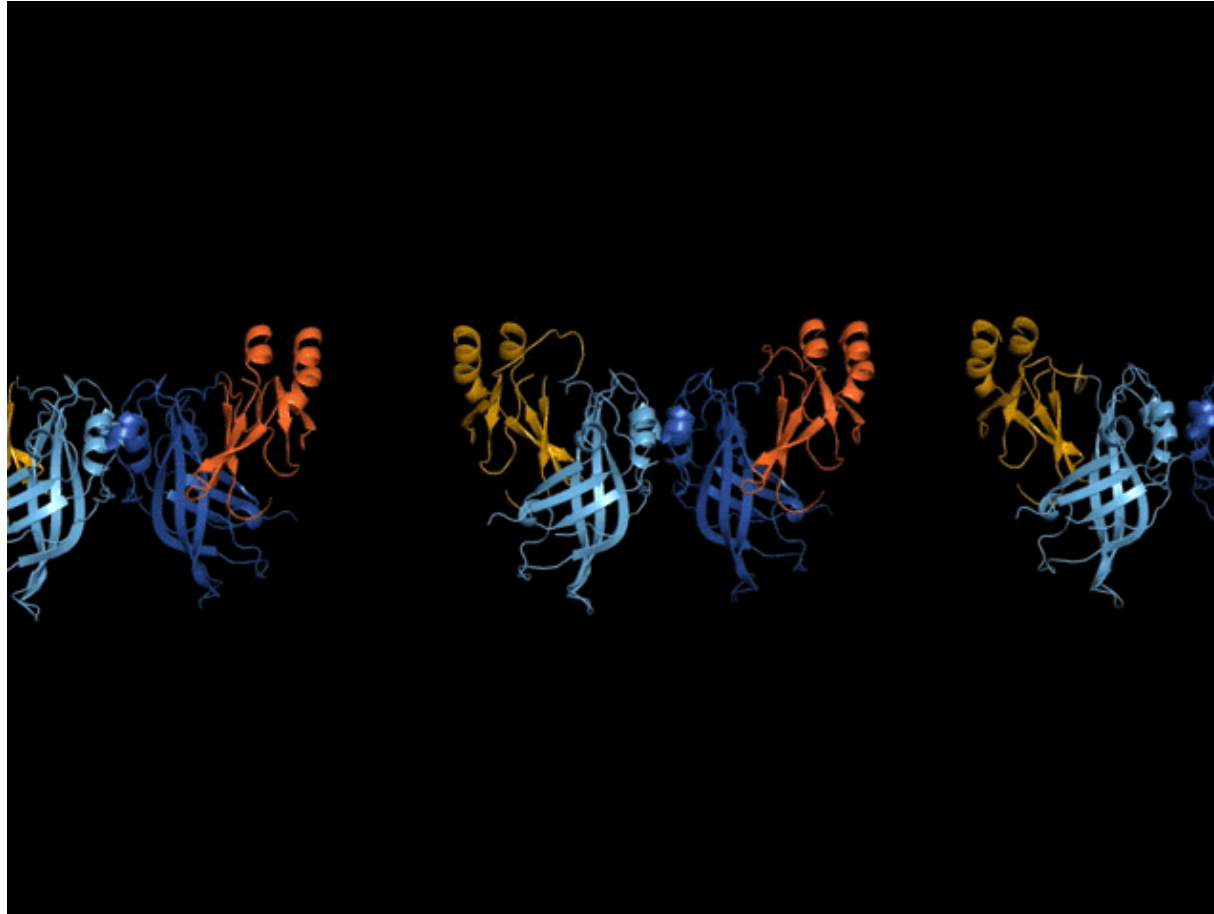
J Ingraham, M Baranov, Z Costello, V Frappier,
A Ismail, S Tie, W Wang, V Xue, F Obermeyer,
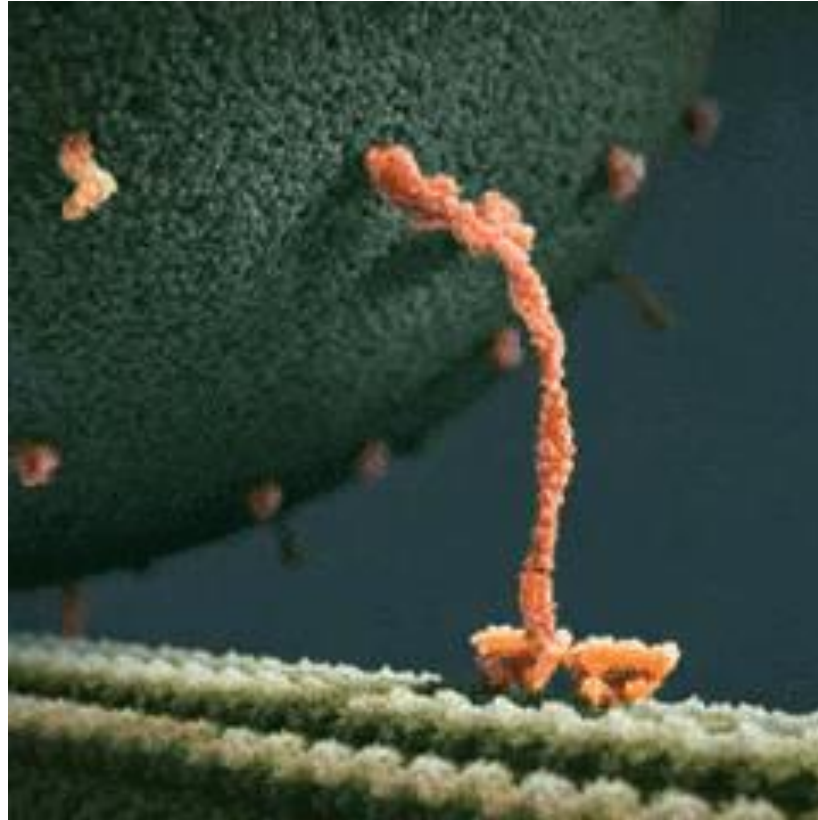A Beam, G Grigoryan

December 1, 2022

Presented by Meret Ackermann

# Proteins – the chief actors in cells

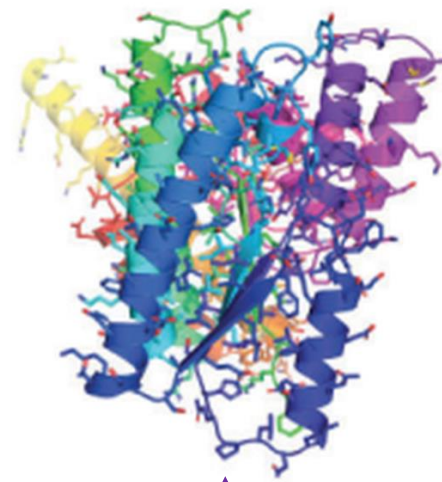# Structure – the key component of function

# Protein Structure Prediction



Sequence

Folding

AlphaFold
RosettaFold
…

# Protein Sequence prediction



Structure

Folding

Protein MPNN
GVP-GNN

…

# Structure first – Protein Design

Backbone

Sequence

Side Chains

Experimental characterization

# Generative diffusion process

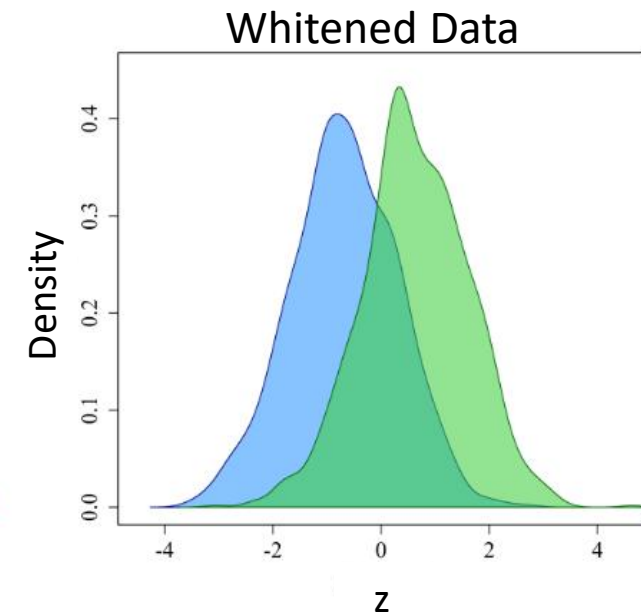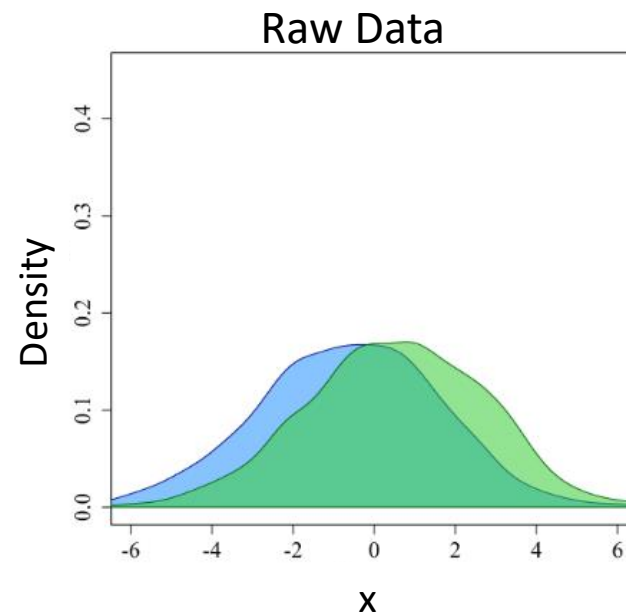Forward SDE. **Training**. Data to Noise.



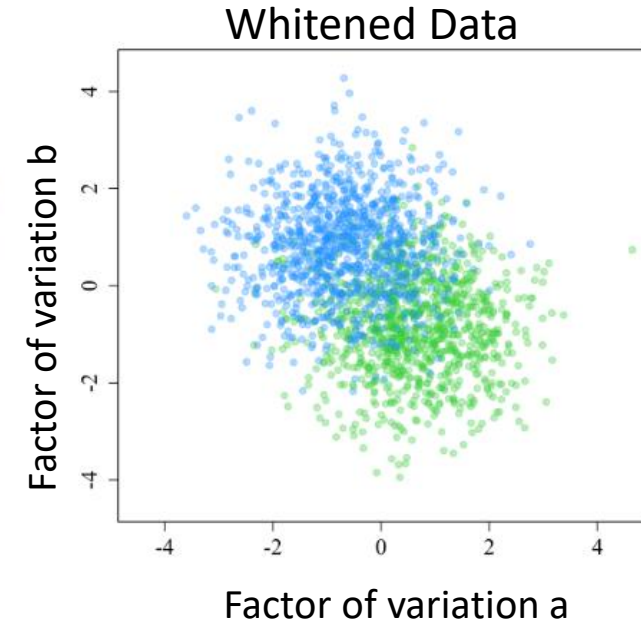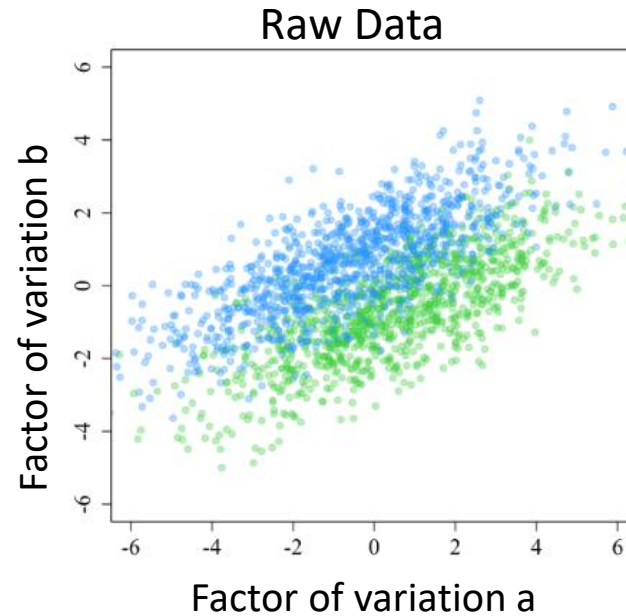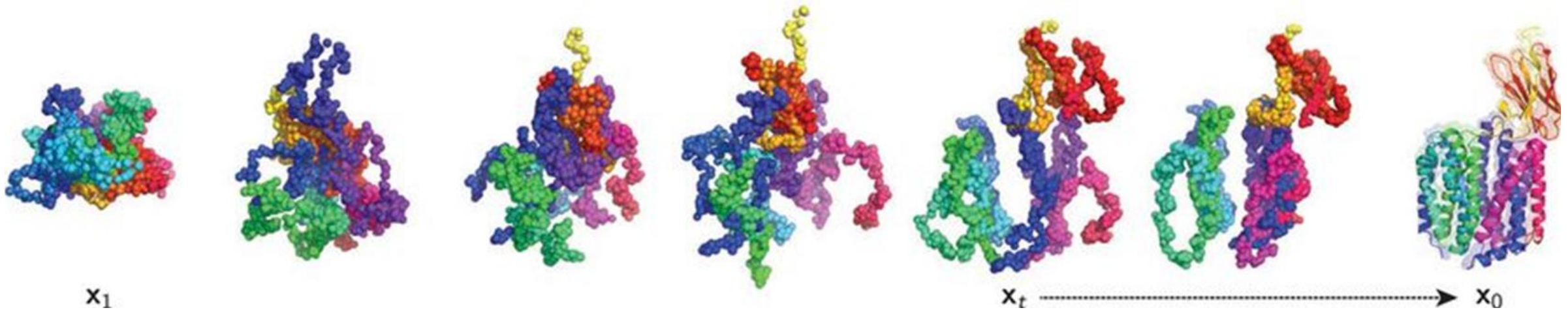Reverse SDE. **Generation**. Noise to Data

# Whitening transformation

Correlated diffusion
as uncorrelated diffusion
in whitened space

$$z = R(x - \mu)$$

# Correlated forward SDE process

$$dx = Rdz = -\frac{\beta_t}{2}Rz\,dt + \sqrt{\beta_t}R\,dw$$



$x_1$                    $x_t \dashrightarrow x_0$

# Constraints as a de-whitening transform



$$F(x) = \sum_{i,j} \boldsymbol{A}_{i,j} \boldsymbol{x}_i \boldsymbol{x}_j$$

$$\mathbb{E}_{p(\boldsymbol{x_t}|\boldsymbol{x_o})}[F(\boldsymbol{x})] = \alpha_t F(\boldsymbol{x}_0) + (1 - \alpha_t)\, \mathbb{E}_{p_{model(x)}}[F(\boldsymbol{x})]$$

# Constraint – Chain Structure



$$r_{i,j} \sim \mathcal{N}(0, \gamma^2 |i-j|)$$

$$\mathbb{E}_{p(x_t|x_0)}\left[D_{ij}^2(x_t)\right] = \alpha_t D_{ij}^2(x_0) + (1-\alpha_t)3\gamma^2 |i-j|$$

# Constraint - Radius of Gyration
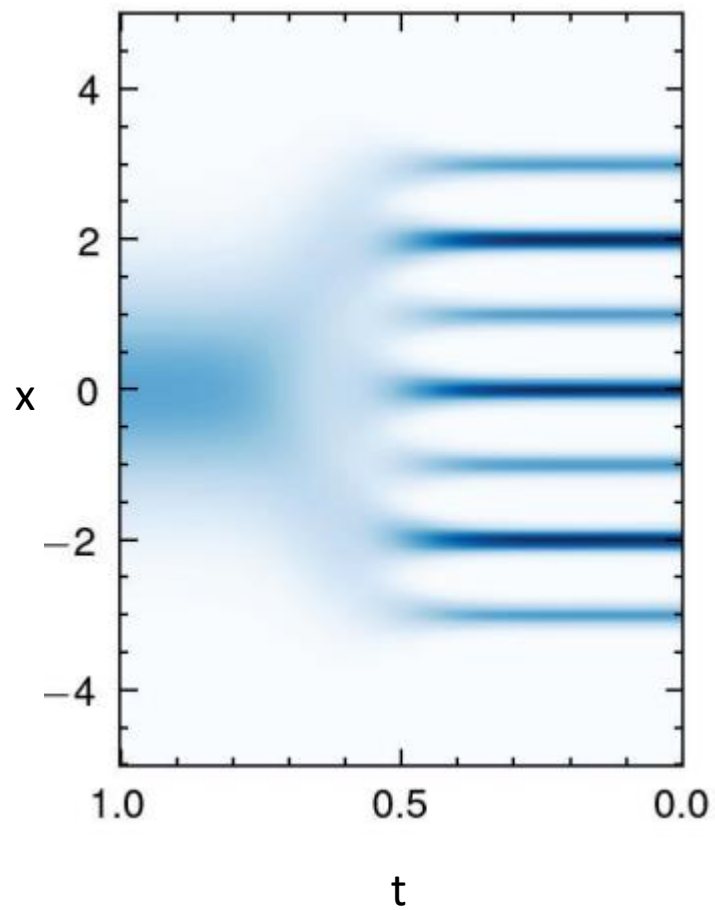


$$R_g^2(x) = \frac{1}{2N^2} \sum_{i,j} D_{i,j}^2(x)$$

# Reverse-time SDE



$$dx = $$

$$\left( -\frac{1}{2}\,x - \boldsymbol{R}\boldsymbol{R}^T \nabla_x \log \boldsymbol{p_t}(\boldsymbol{x}) \right) \beta_t dt$$

$$+\sqrt{\beta_t}\boldsymbol{R}d\widetilde{w}$$

# Correlated Reverse-time SDE



$$dx = \left(-\frac{1}{2}x - RR^T\nabla_x \log p_t(x)\right)\beta_t dt + \sqrt{\beta_t}R\,d\widetilde{w}$$
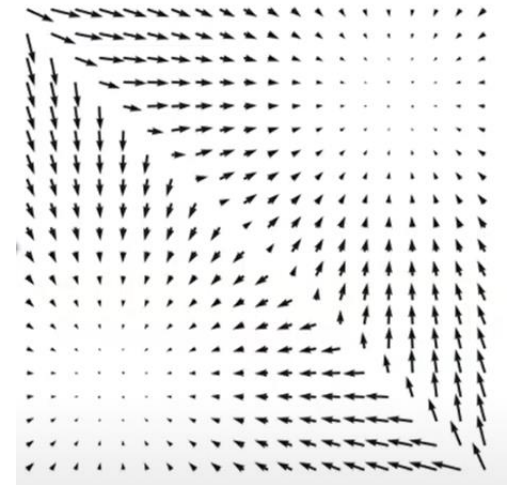
# Score Estimation



$p_{data}(x)$

$\nabla_x \log p_{data}(x)$

Samples

$\sim \nabla_x \log p(x)$

# Optimized denoiser

$$\nabla_{\boldsymbol{x}} \log p_t(\boldsymbol{x}) = \left((1 - \alpha_t)\boldsymbol{R}\boldsymbol{R}^{\boldsymbol{T}}\right)^{-1} (\sqrt{\alpha_t}\widehat{\boldsymbol{x}}_\theta(\boldsymbol{x}, t) - \boldsymbol{x})$$

$$\mathcal{L}_{\mathbf{x}}^{\mathrm{reg}}(\mathbf{x}; \boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t|\mathbf{x}), t \sim \mathrm{Unif}(0,1)} \left[ \frac{\alpha_t \beta_t}{2(1 - \alpha_t)^2} \left\| \left(\mathbf{R}^{-1} + \omega\mathbf{I}\right) (\hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \mathbf{x}) \right\|_2^2 \right]$$

# Optimized denoiser

Noisy structure
$x_t$

Confidence-weighted predicted
inter-residue geometries

Predicted denoised structure
$\hat{x}_\theta(x, t)$



Random Graph
Neural Network

Equivariant
Geometry Solver

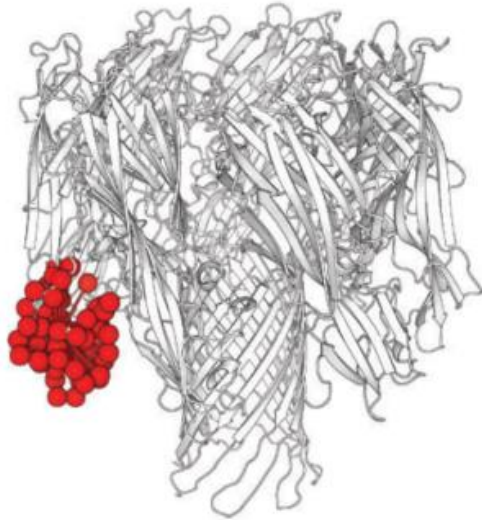# Reduced computational complexity

Random Graph
Neural Network

$\mathcal{O}(NlogN)$
or
$\mathcal{O}(N)$
edges

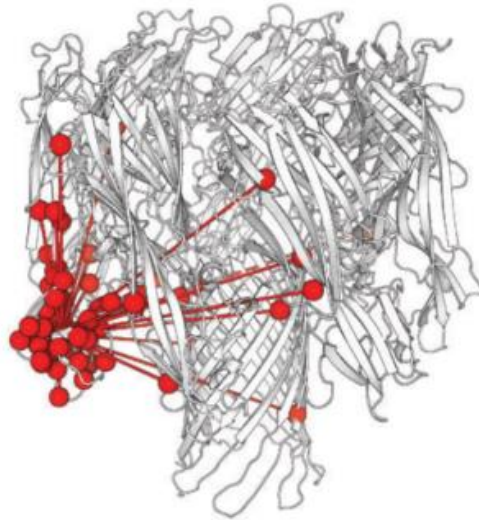# Sub- $\mathcal{O}(N^2)$ scaling - Random edge sampling
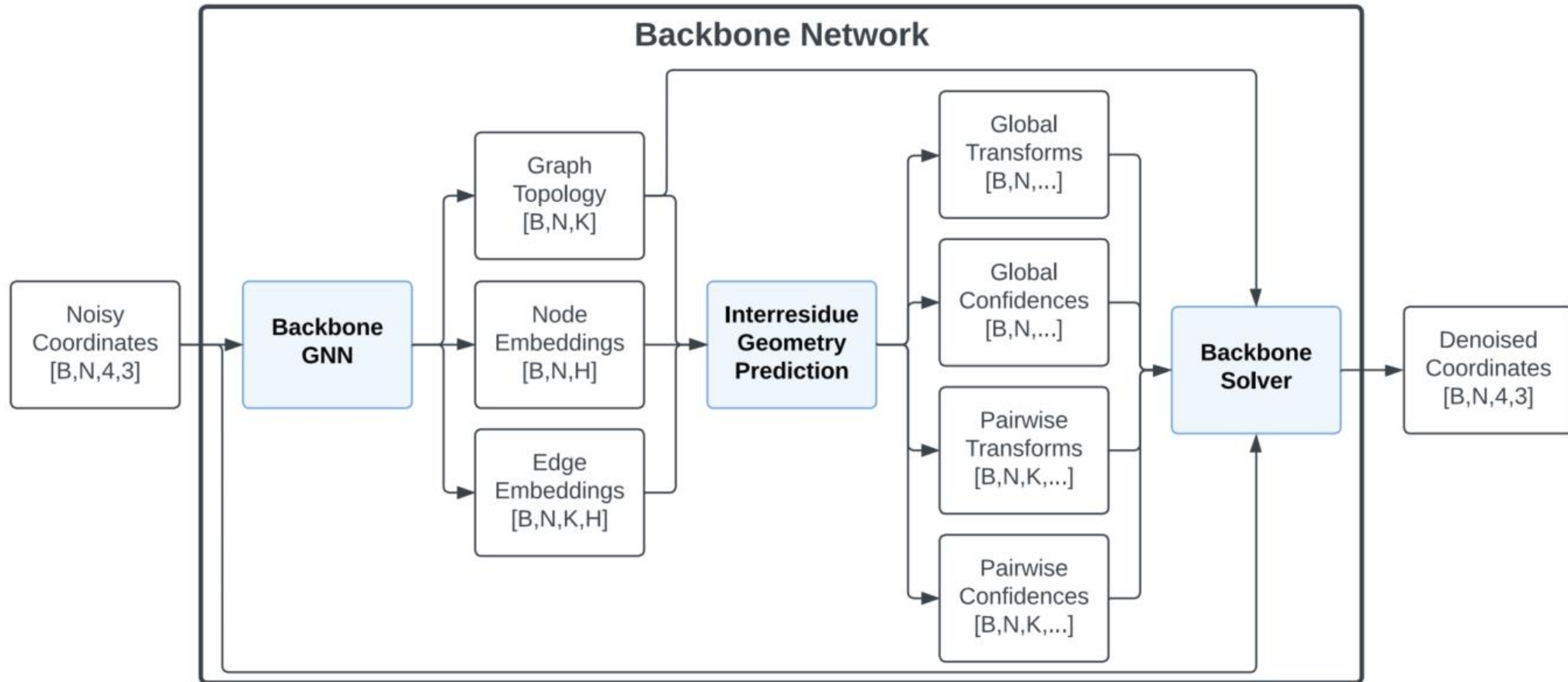


Deterministic graph
k-NN

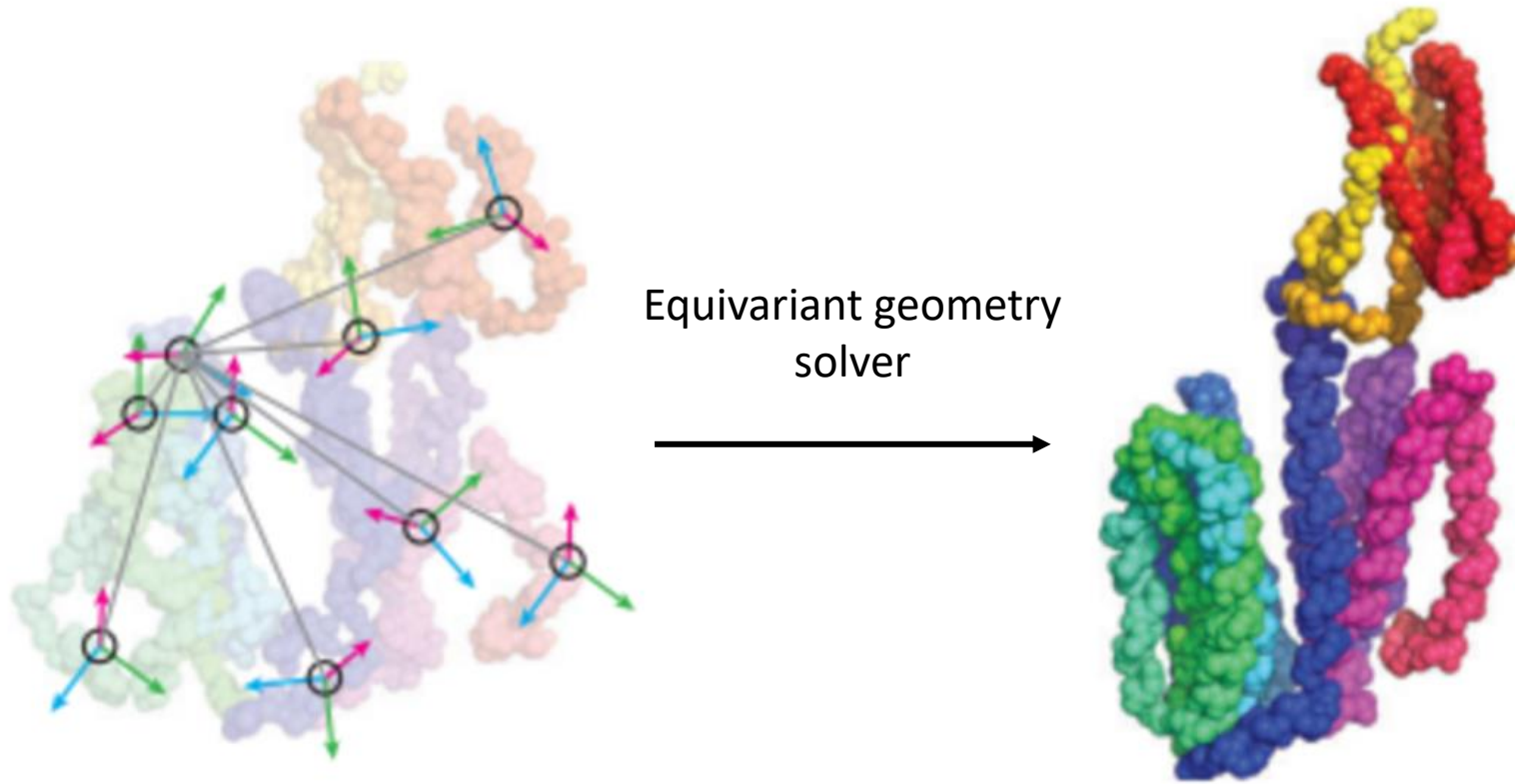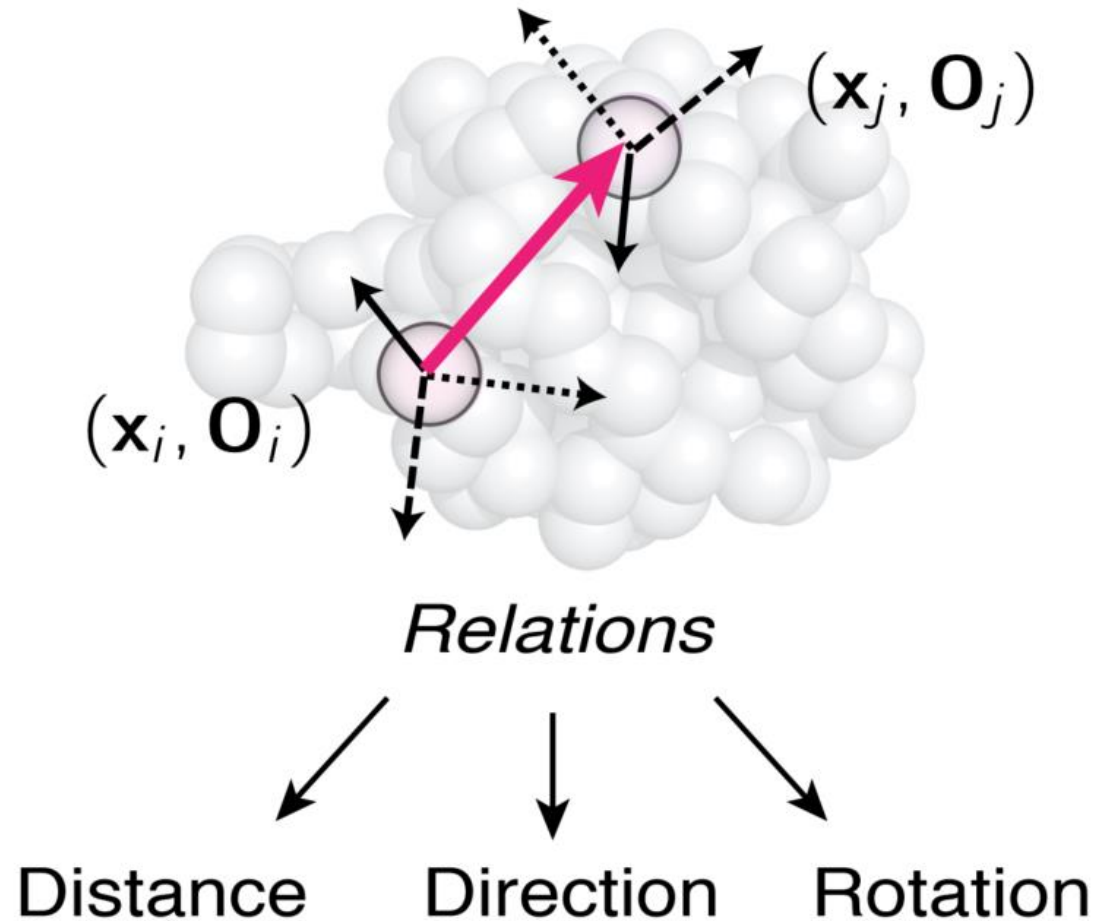Random graphs
Inverse cubic

Mixed graph
20 k-NN + 40 Inverse Cubic

# Backbone graph neural network

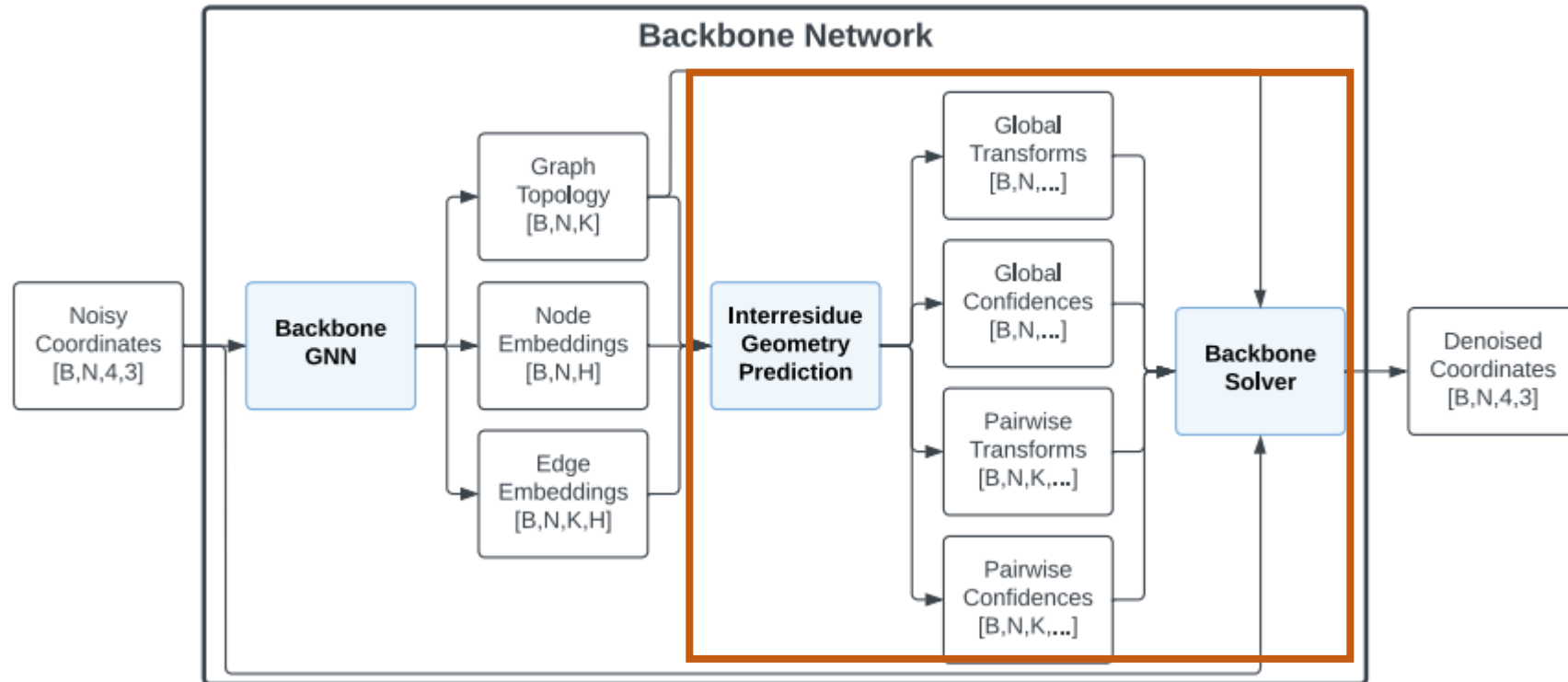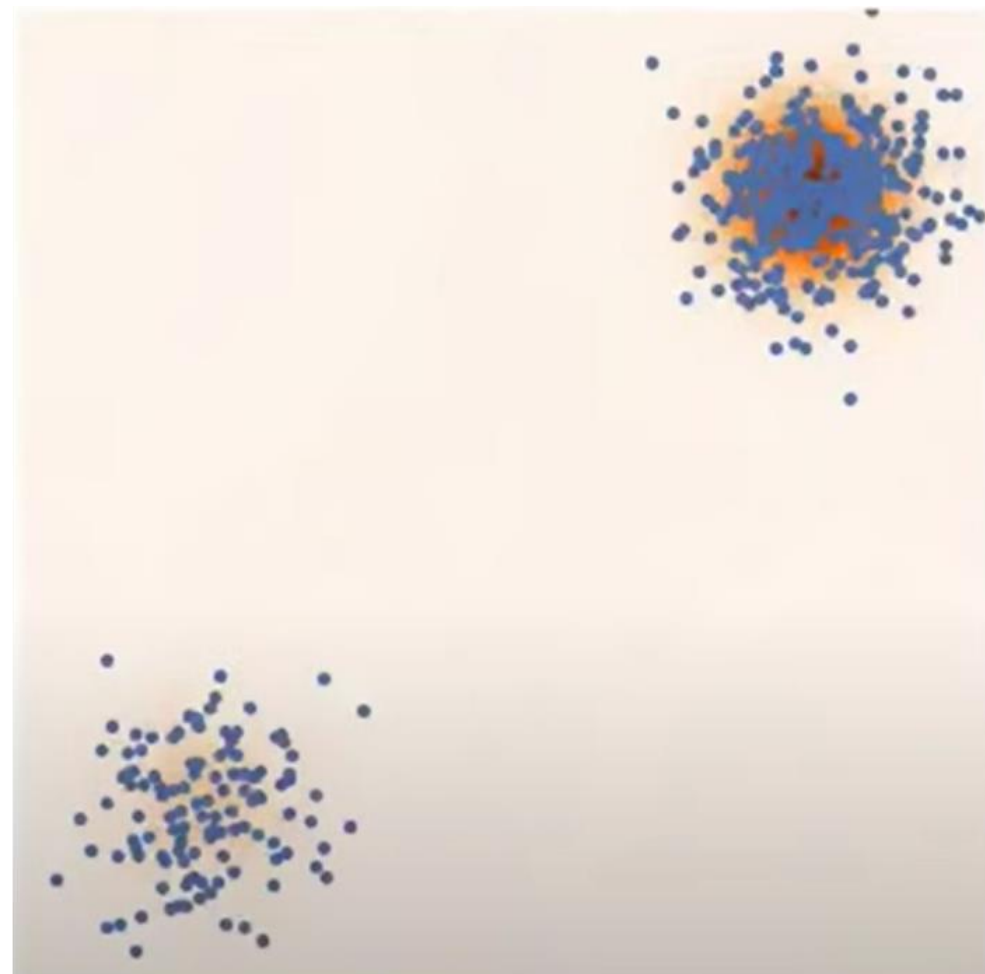# Equivariant geometry solver
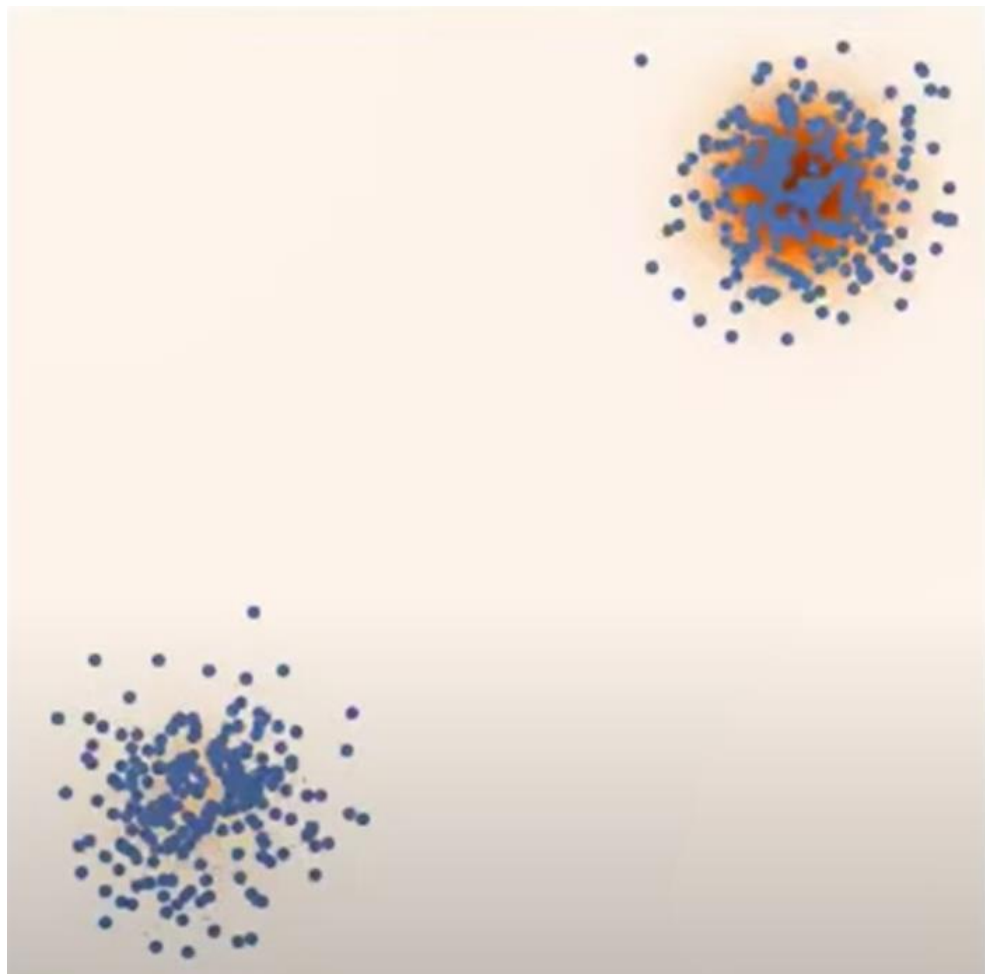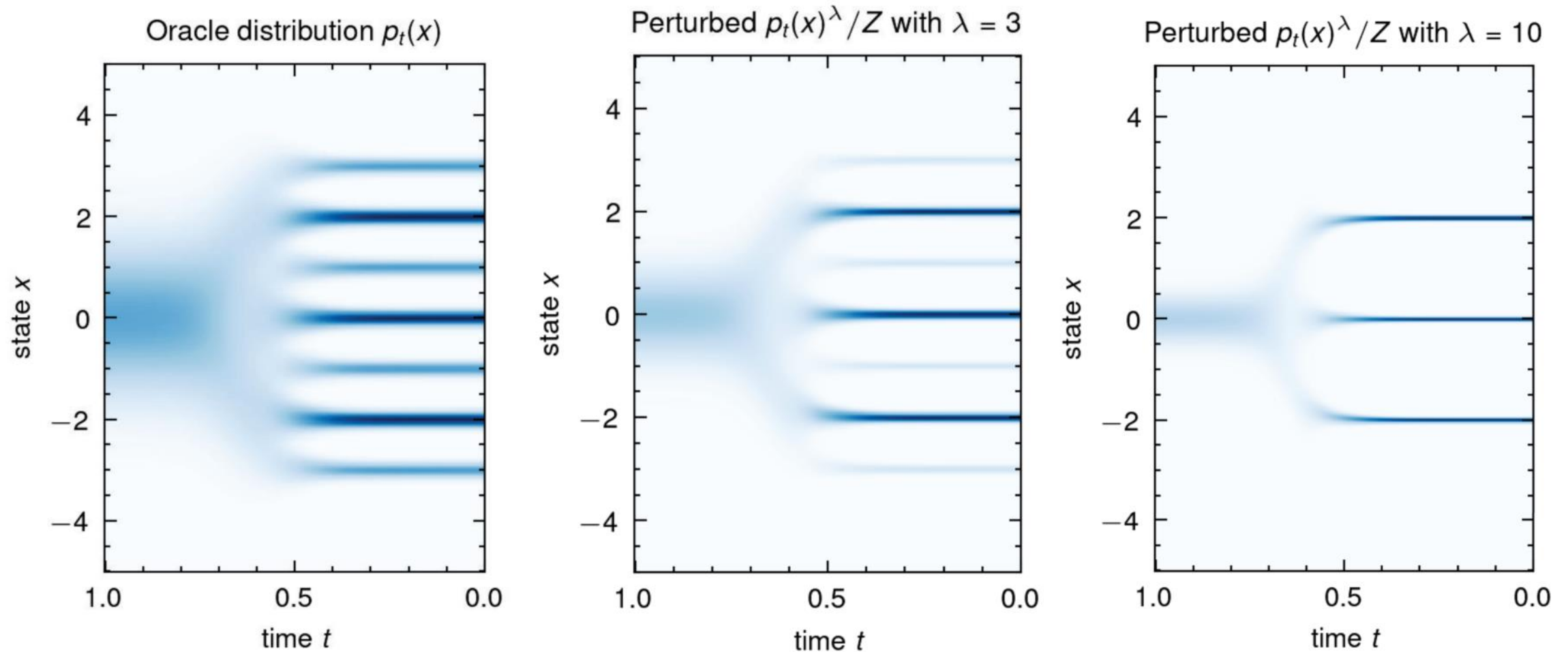


Equivariant geometry solver

# Invariant local frame relations

# Equivariant geometry solver

# Sampling of the backbone - overdispersion

# Low temperature sampling – reweight and concentrate



Oracle distribution $p_t(x)$     Perturbed $p_t(x)^\lambda / Z$ with $\lambda = 3$     Perturbed $p_t(x)^\lambda / Z$ with $\lambda = 10$

# Annealed reverse-time SDE

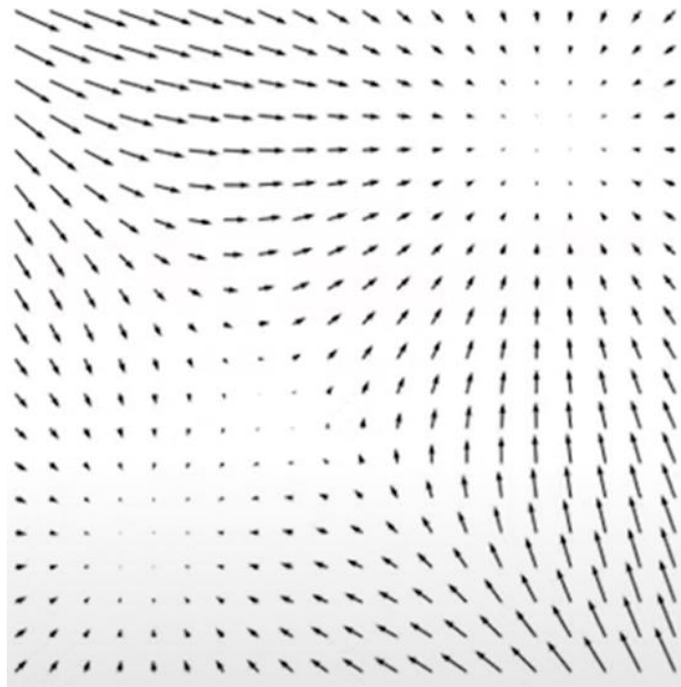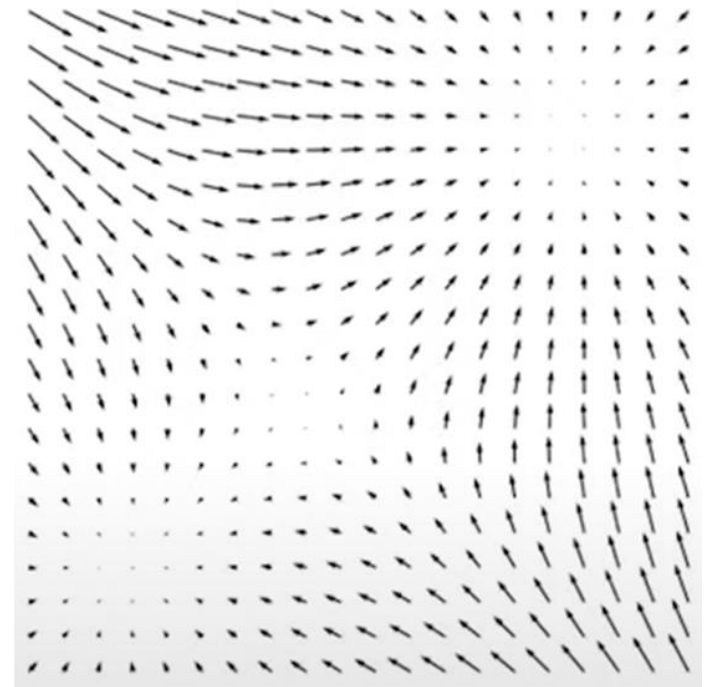$$\lambda_t \approx \frac{\lambda_0}{\alpha_t + (1 - \alpha_t)\lambda_0}$$
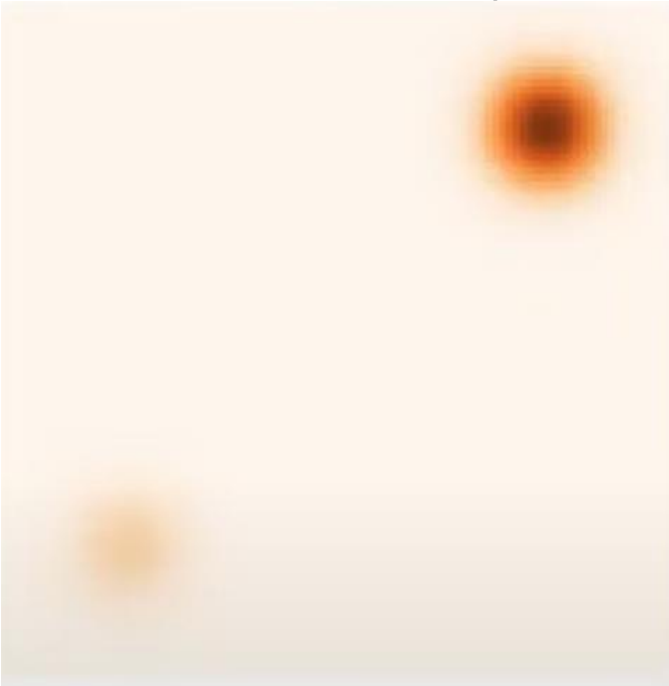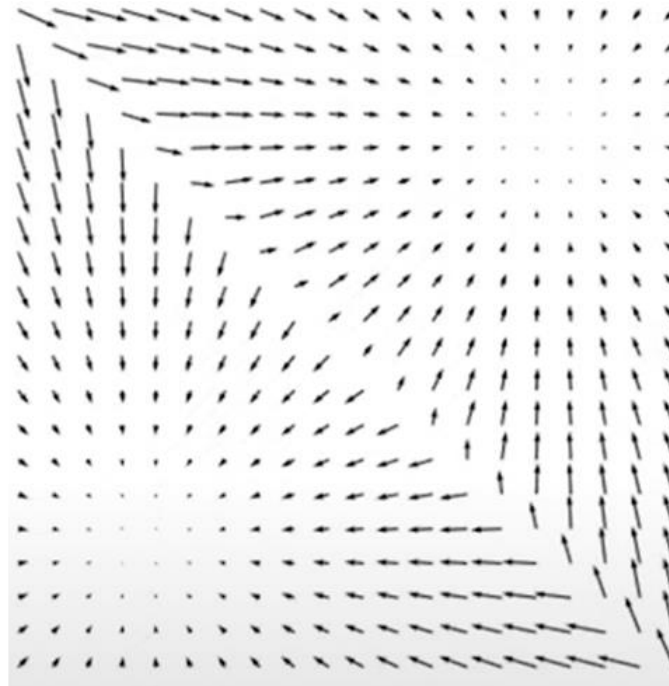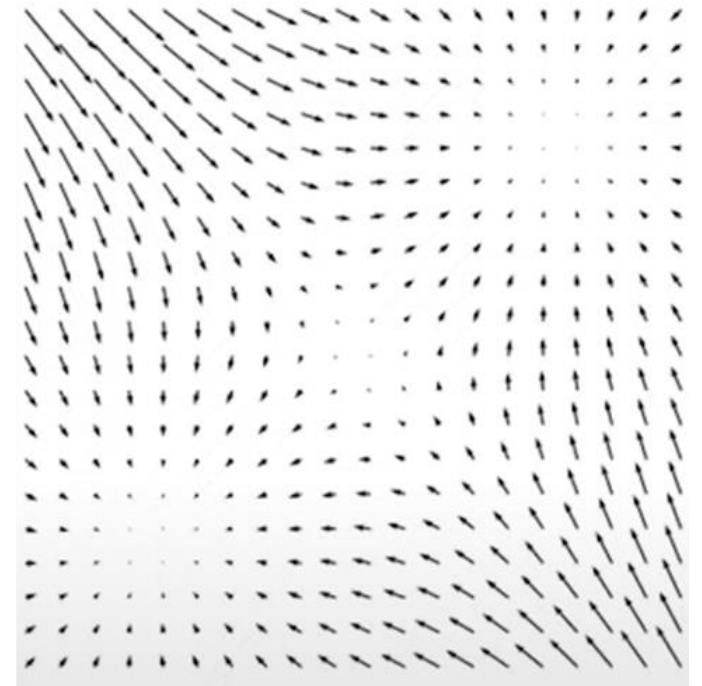
# Score evolution



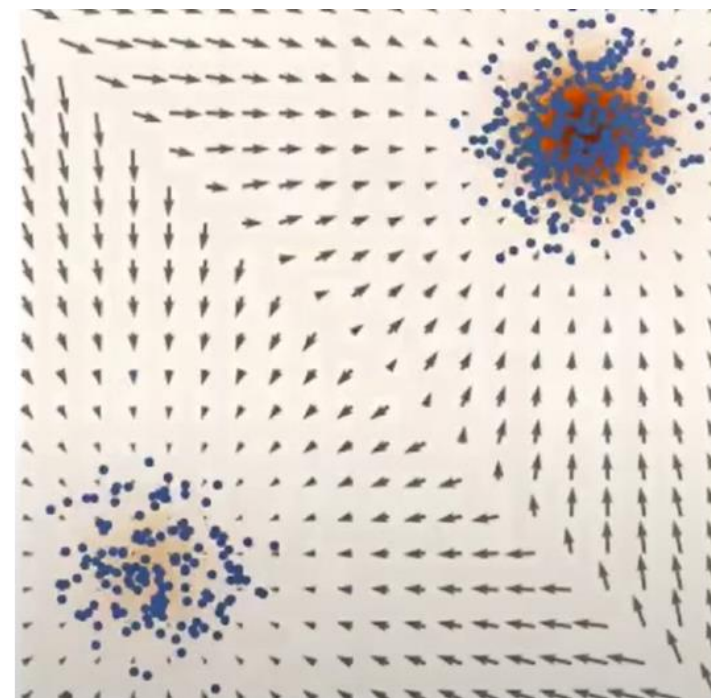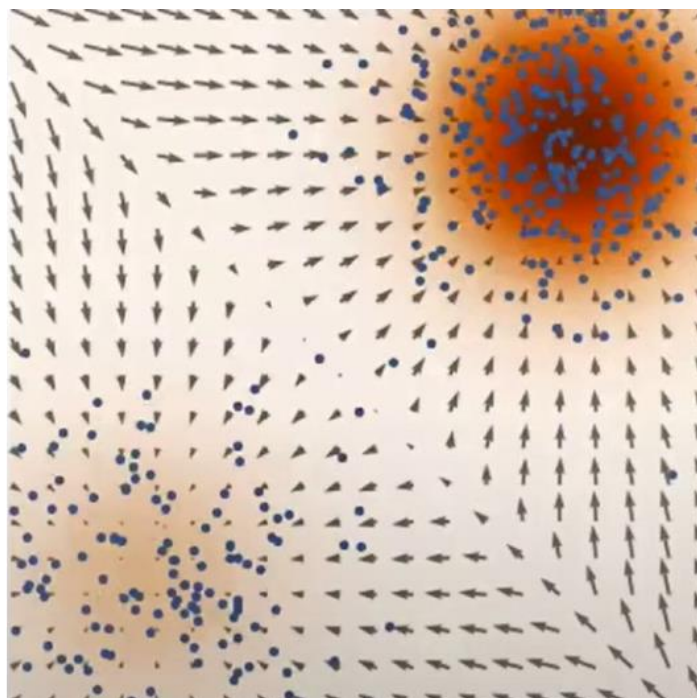Data Density       Data Score       Estimated Score
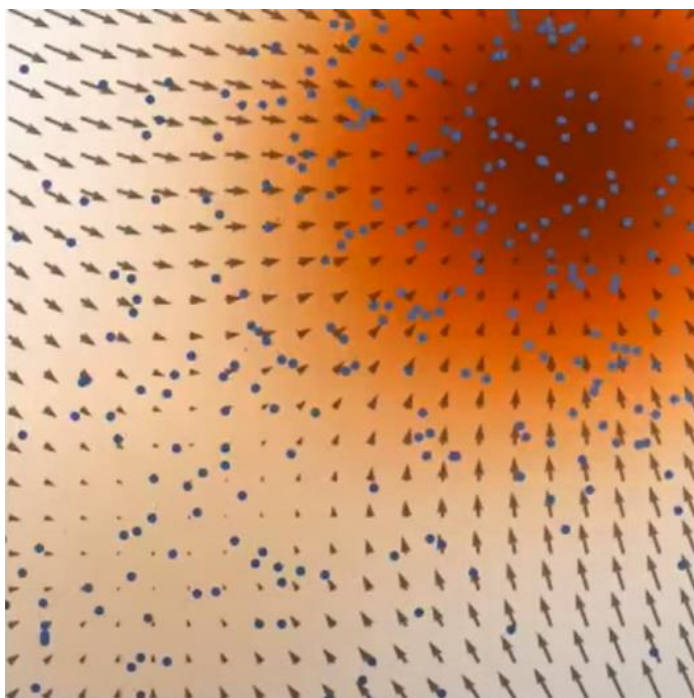
# Score evolution

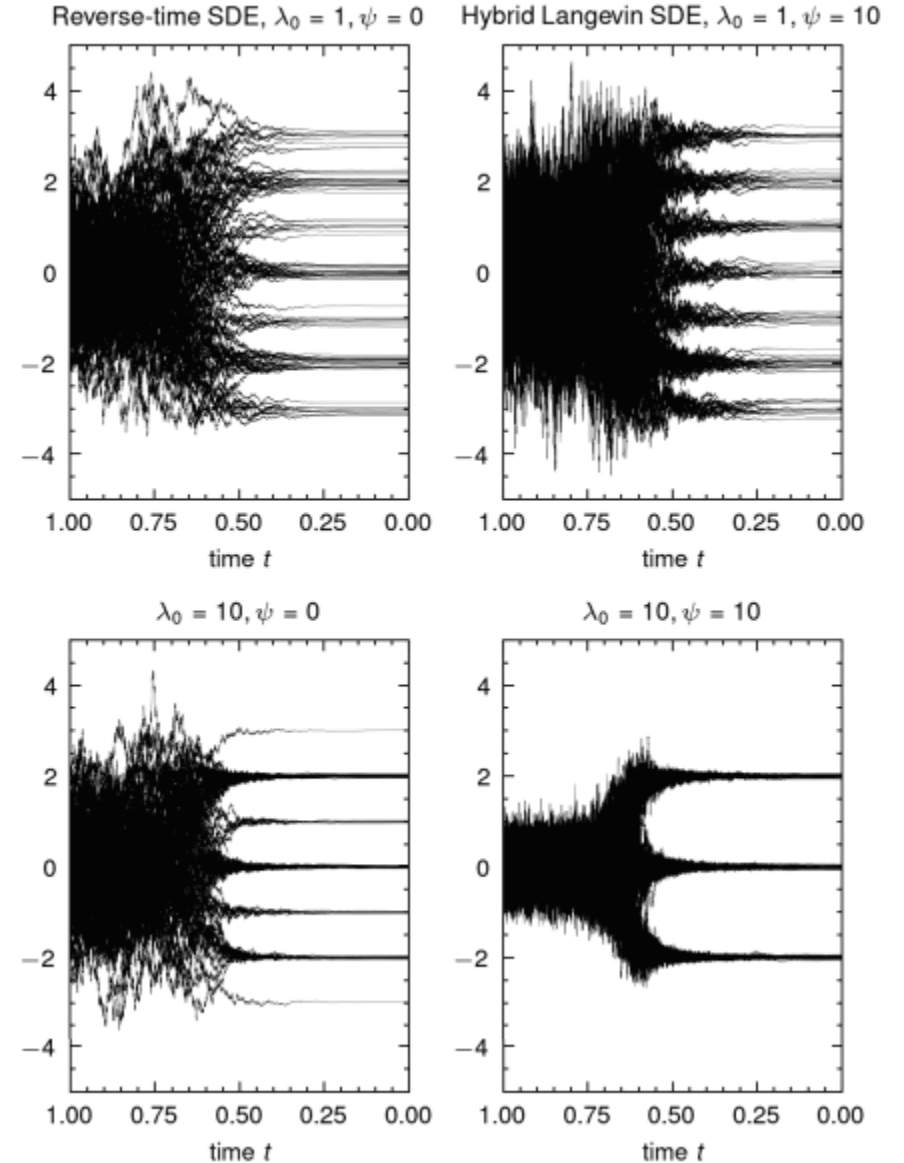Data Density

Data Score

Estimated Score
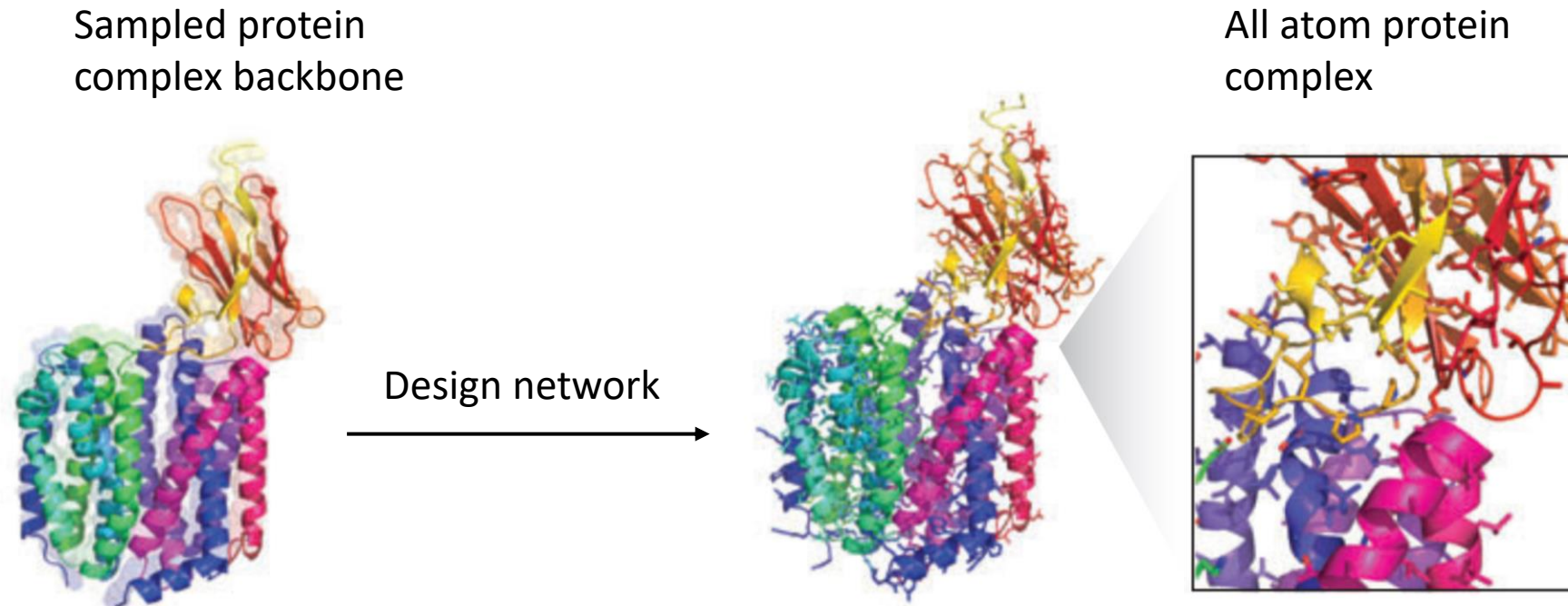
# Annealed Langevin dynamics
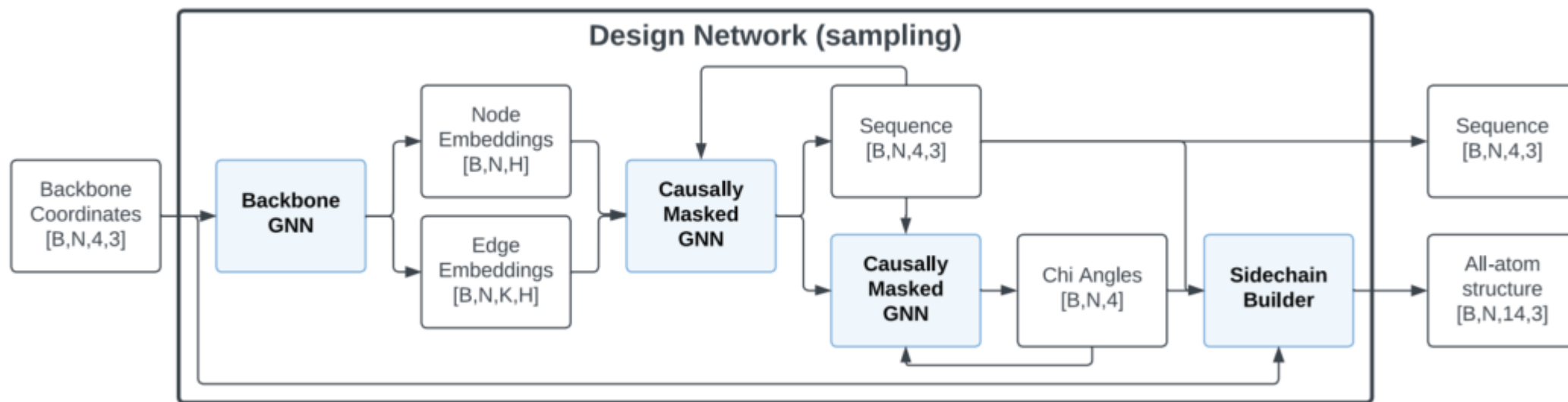
# Hybrid Langevin SDE

$$dx = \left( -\frac{1}{2}x - \left( \lambda_t + \frac{\lambda_0 \psi}{2} \right) \frac{\sqrt{\alpha_t}\widehat{x_{\theta}}(x,t) - x}{1 - \alpha_t} \right) \beta_t dt + \sqrt{\beta_t(1 + \psi)}\, R\, d\overline{w}$$
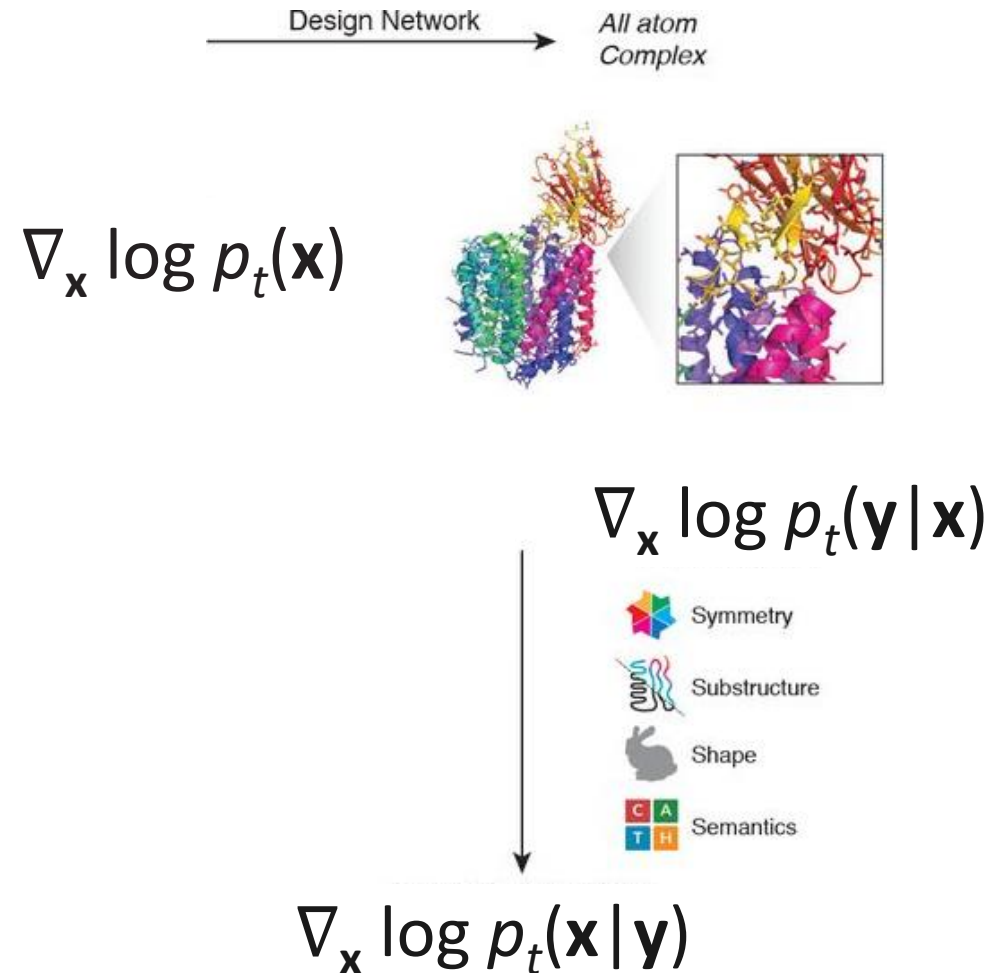
# From backbone to sequence and heavy atom position



Sampled protein complex backbone

Design network

All atom protein complex

# Design Network

# Conditional modeling



$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$$

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{y}|\mathbf{x})$$

Symmetry

Substructure

Shape

Semantics

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}|\mathbf{y})$$

# Conditional modeling

Bayes' rule

$$p(x|y) = \frac{p(x)p(y|x)}{p(y)}$$

Bayes' rule for score functions

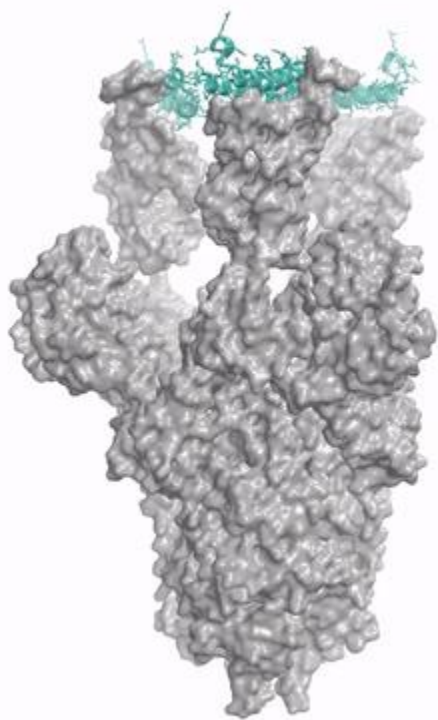$$\boxed{\nabla x \log p(x)} + \boxed{\nabla x \log p(y|x)} - \cancel{\nabla x \log p(y)}$$
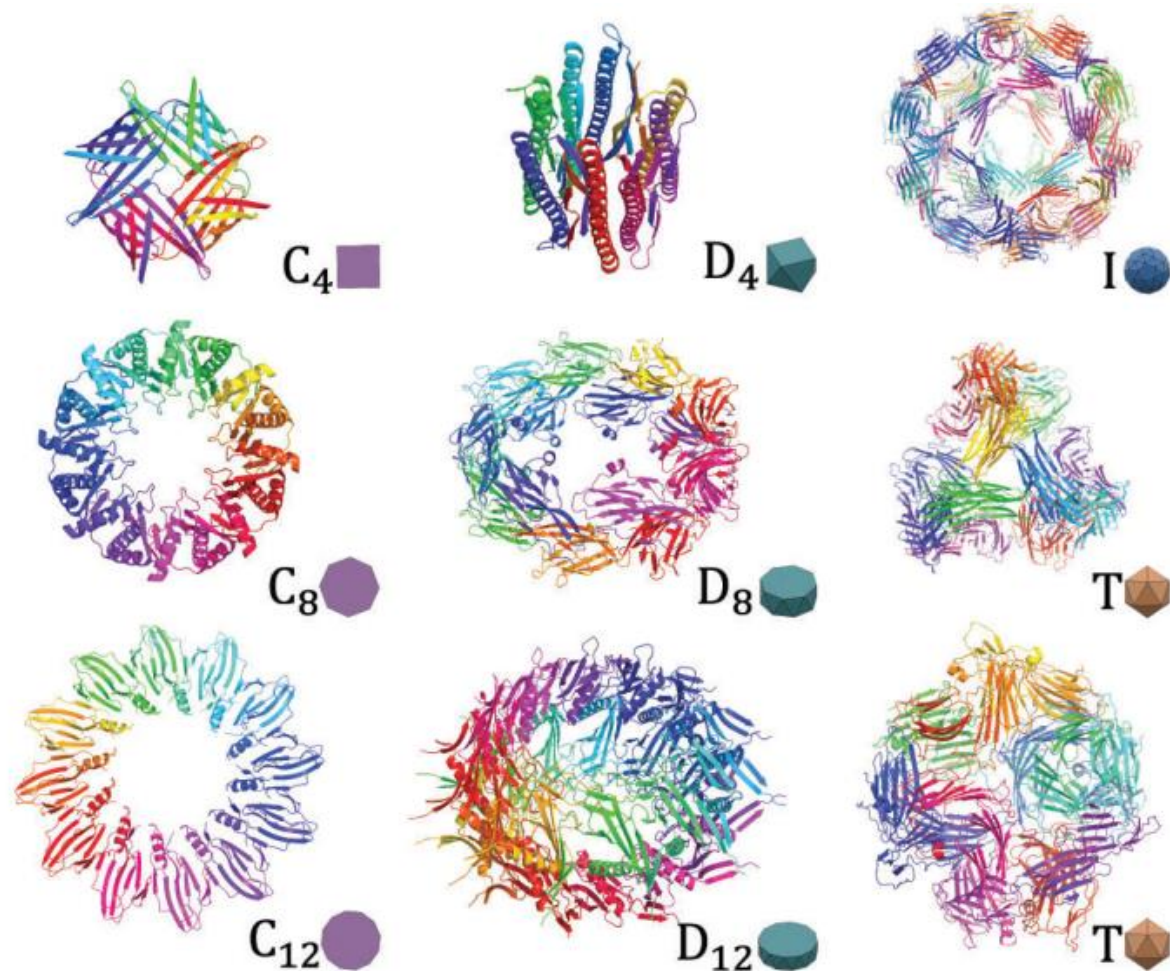
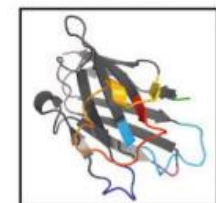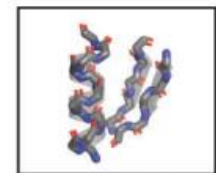Unconditional
score

e.g. classifier

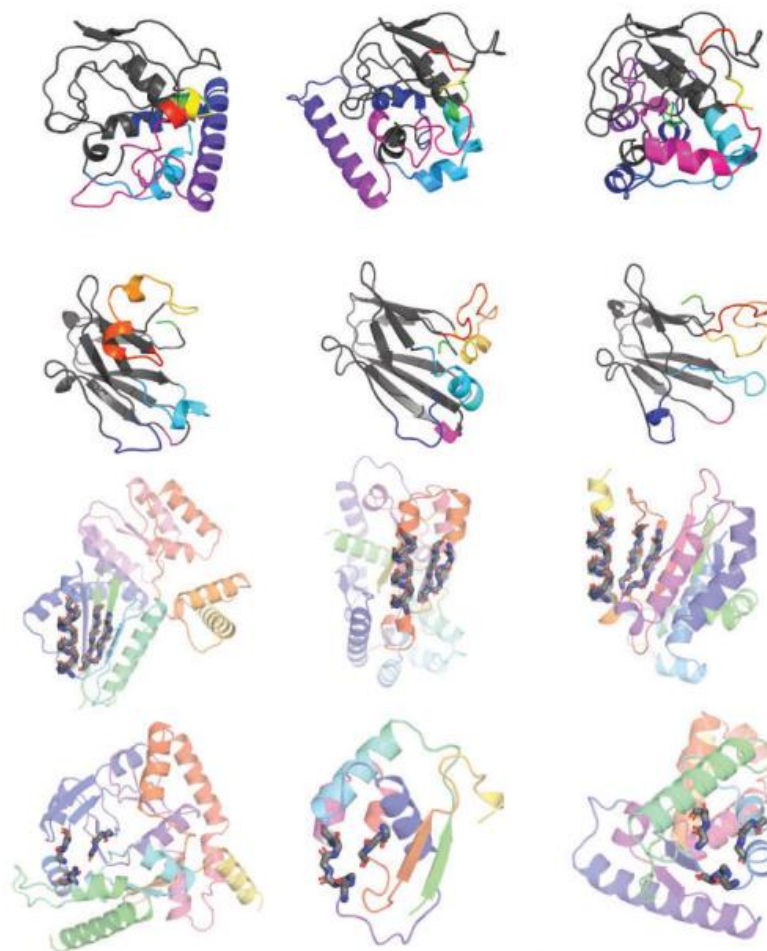# Conditional modeling

# Symmetry and substructure guided diffusion
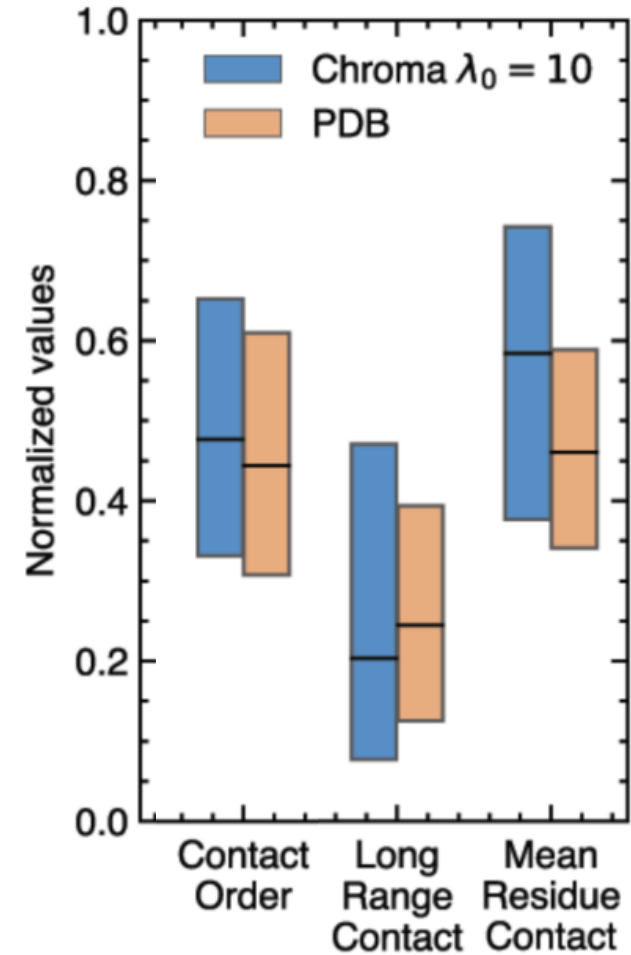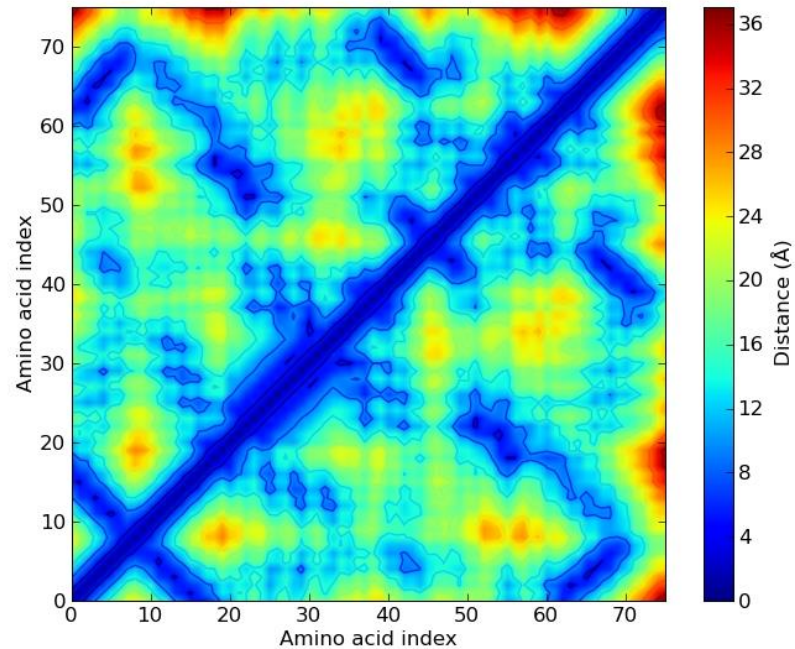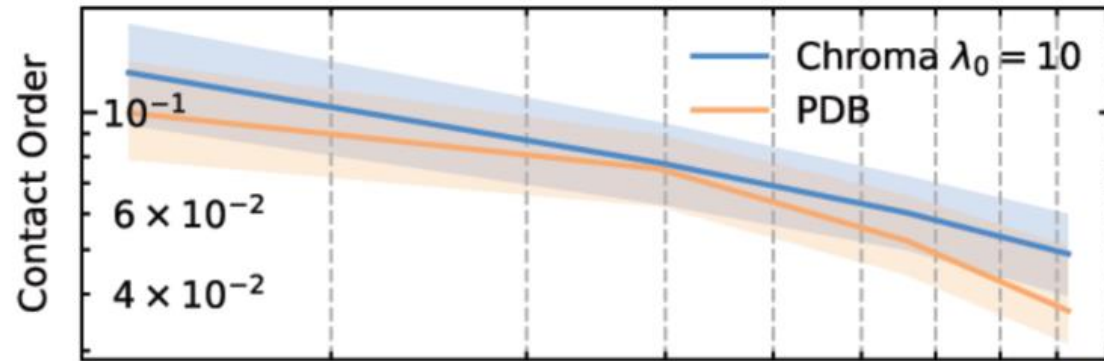
# Evaluation

- 50,000 single chains, 10,000 complexes - qualitative
- 10,000 single chain proteins - quantitative
- $\lambda_0 = 10$
- $\psi = 2$
- 200 steps
- Single chain lengths N: p(N) = 1/N
- Complex # chain and N = # chain and N of random complex from PDB
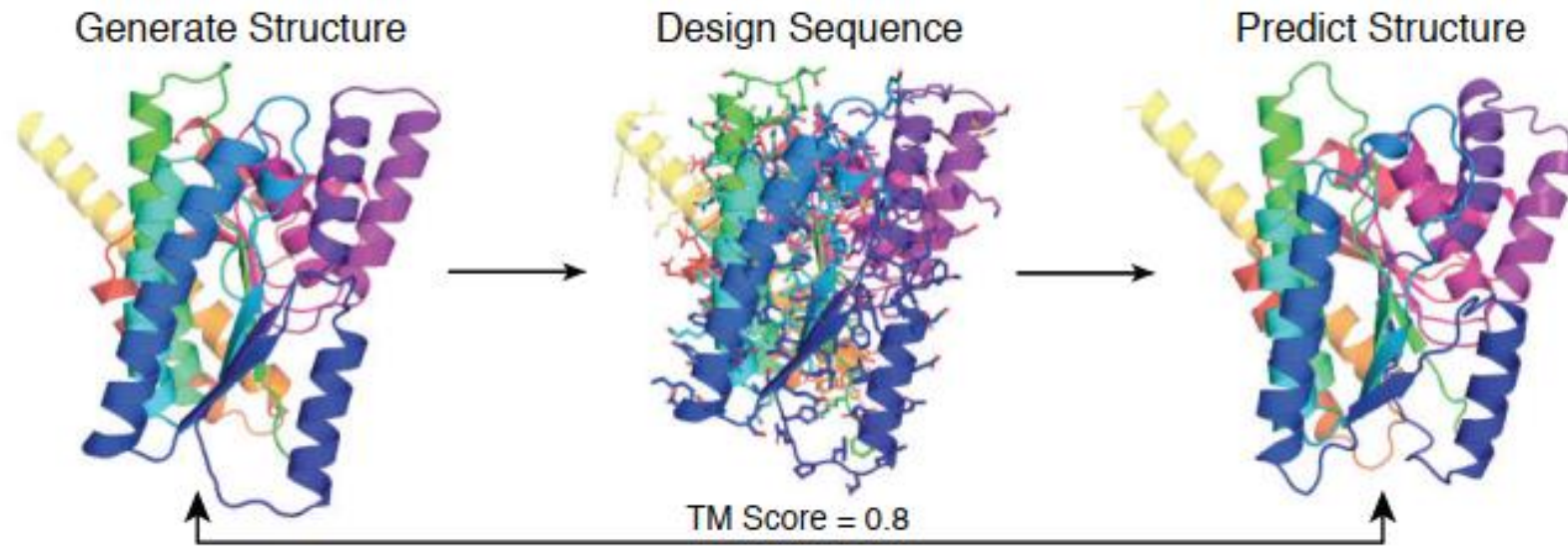
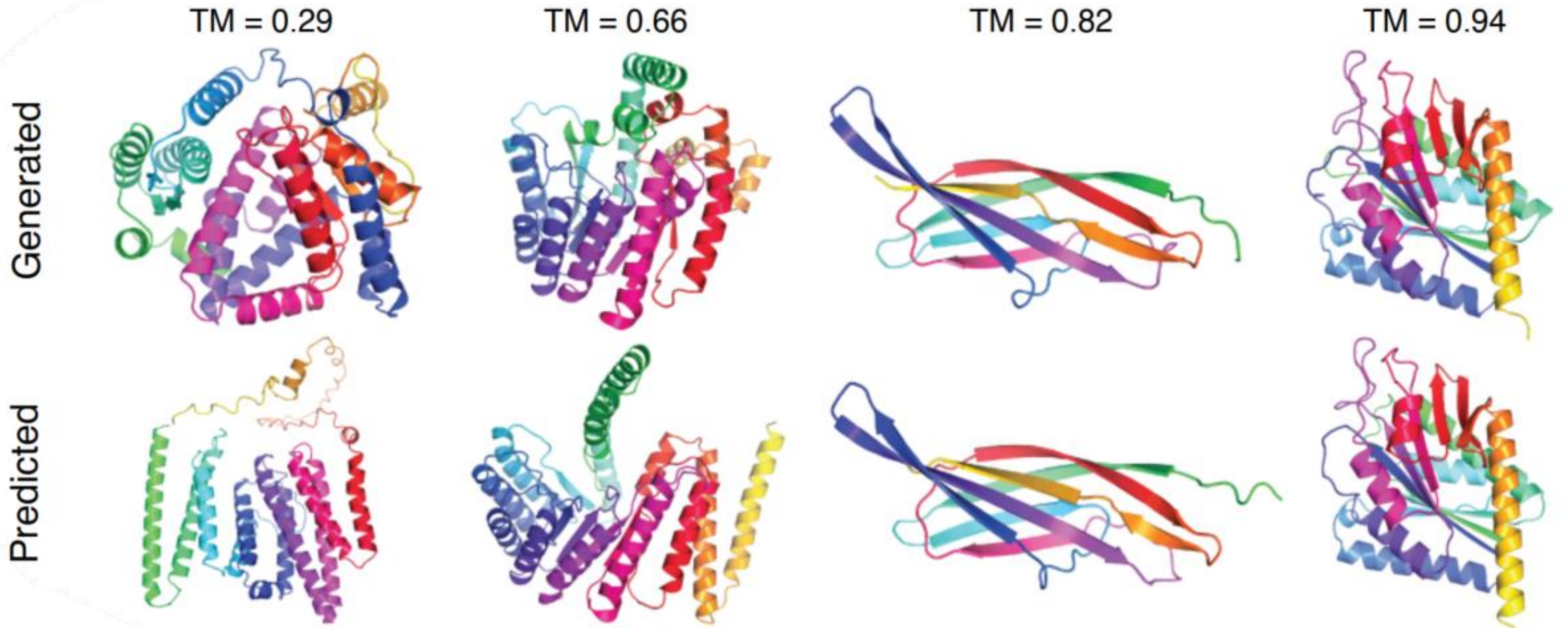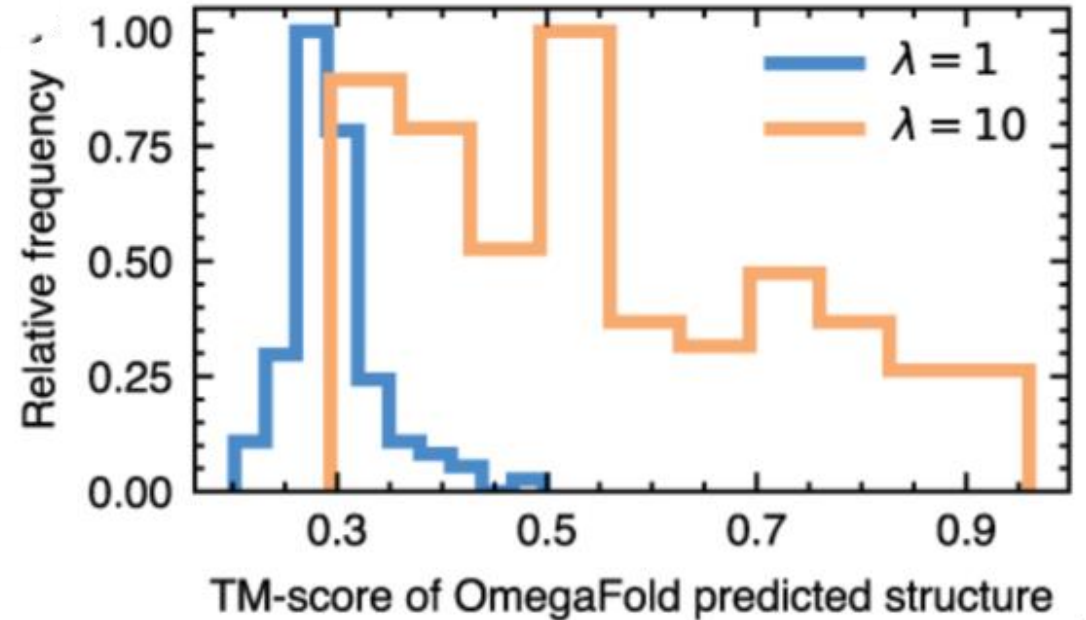# Evaluation – Secondary structurs

# Evaluation – Residue interactions

# Evaluating Chroma by structure prediction with OmegaFold



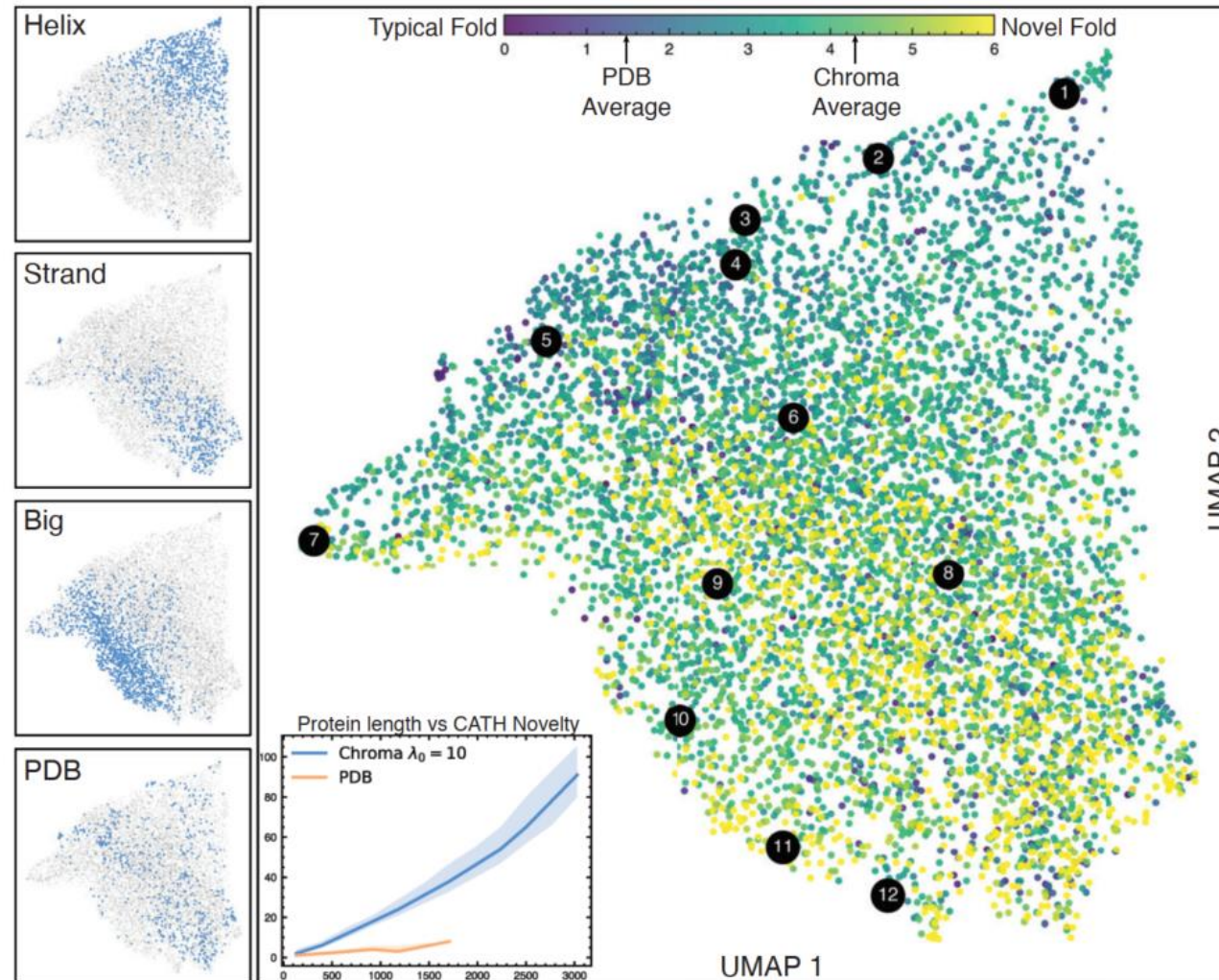Generate Structure → Design Sequence → Predict Structure
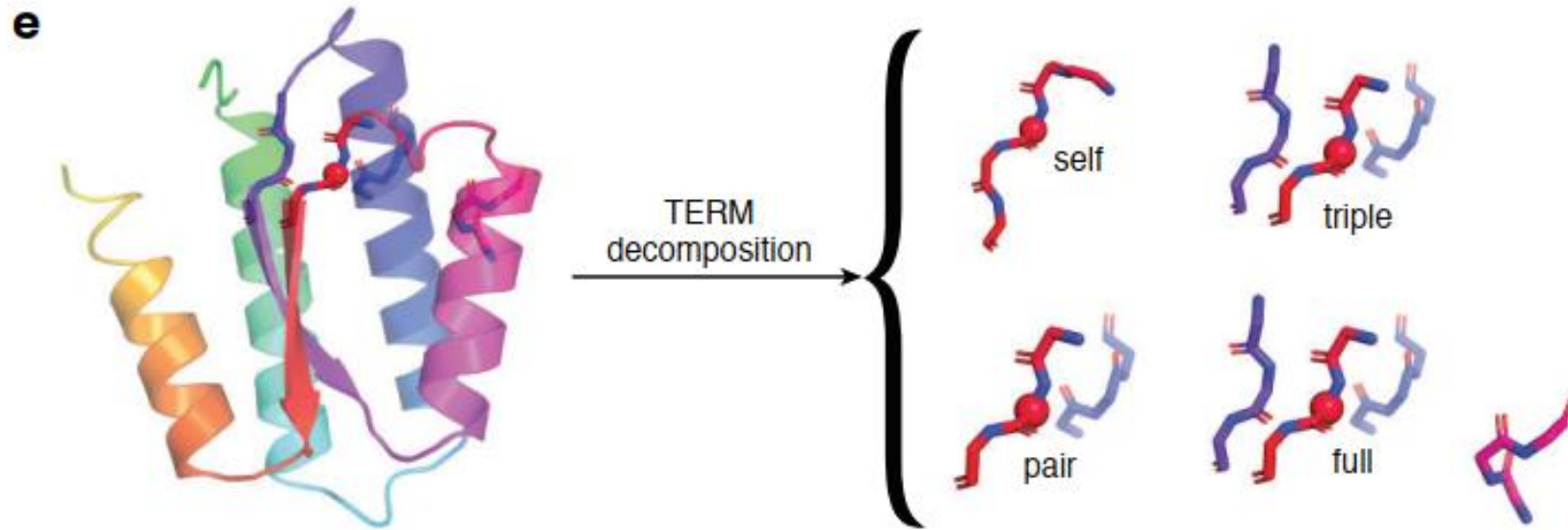
TM Score = 0.8

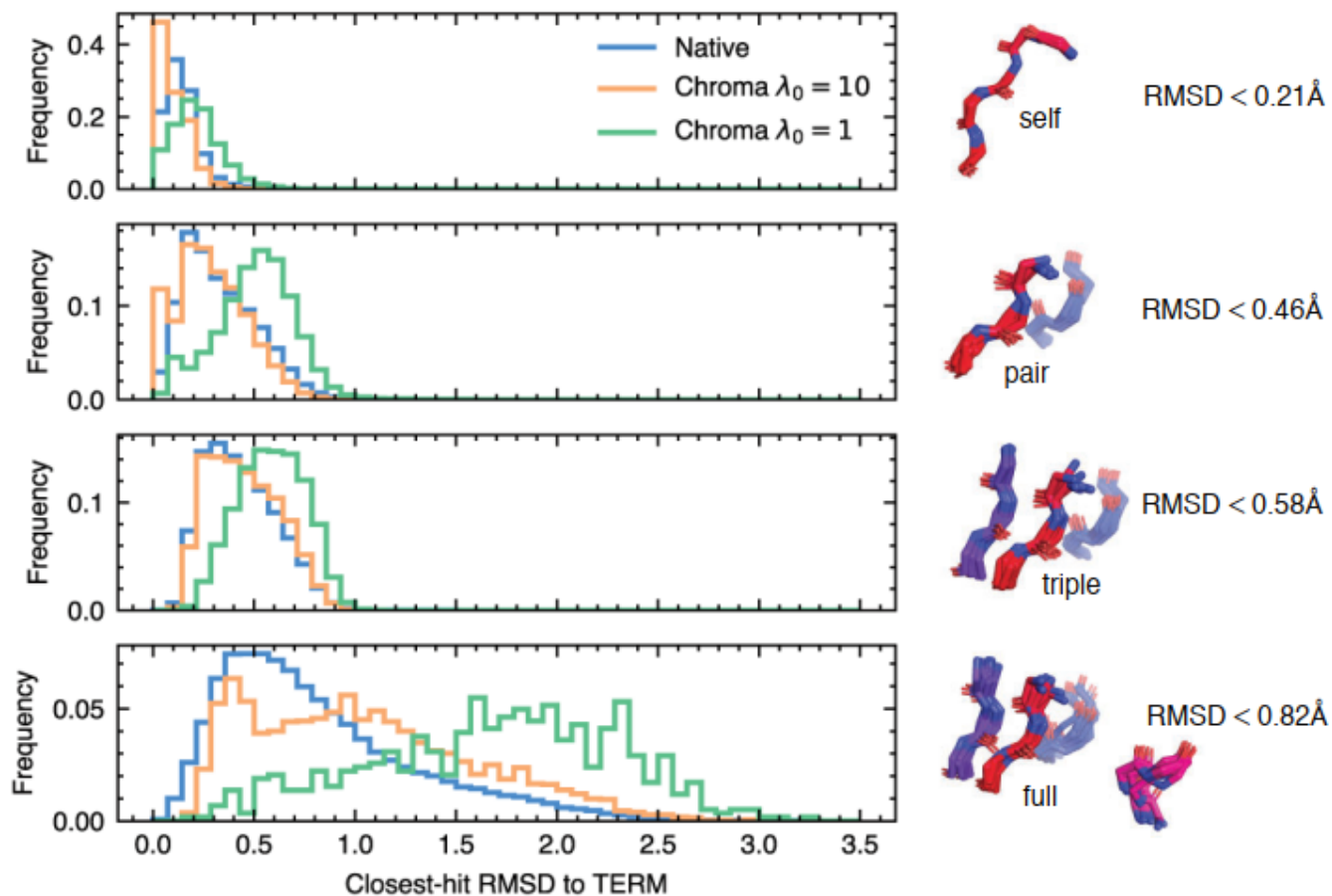# Evaluation - TM-scores

# Evaluation - TM-scores

# Evaluation - Novelty and structural homology

# Evaluation - TERMs

# Evaluation - Closest-match RMSD for TERMs

# Limitations

+ combination of promising maturities in GDMs

+ elegant way they implement empirical knowledge


- missing experimental characterization

- no quantitative evaluations for many designs and design choices

- no benchmarking

- sequential generation of backbone, sequence and rotamers

- choice of model to evaluate folding

Illuminating chroma?!