

Generative Models of protein sequences

Marion Chauveau, *2nd* year Phd student

Supervisors: Ivan Junier & Olivier Rivoire / Collaborator: Yaakov Kleeorin

Generative Models

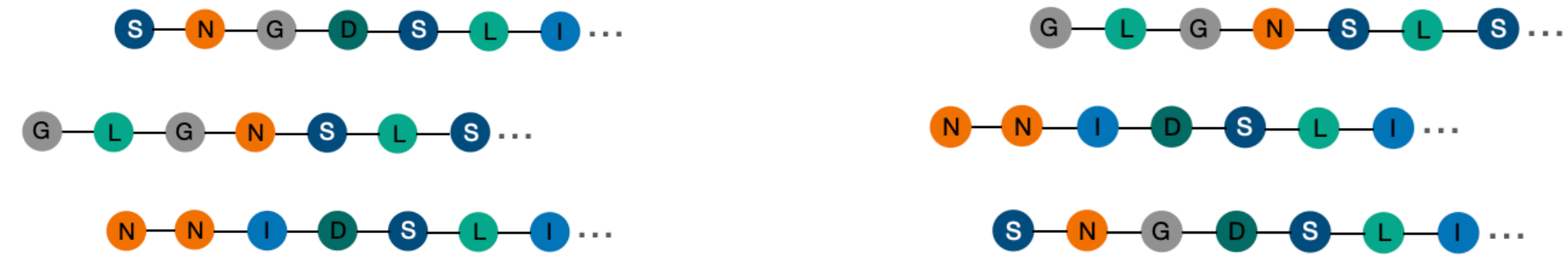


Training data $\sim P_{data}$

Artificial data $\sim P_{model}$

1. Learn P_{model} similar to P_{data}
2. Generate samples from P_{model}

Generative Models



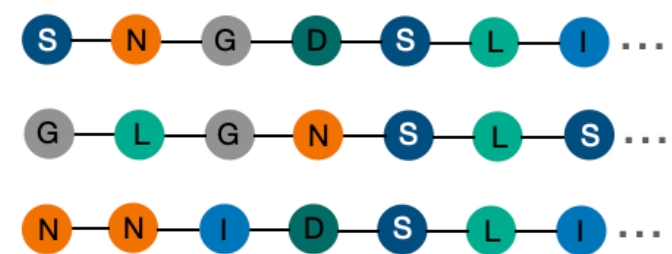
Training data $\sim P_{data}$

Artificial data $\sim P_{model}$

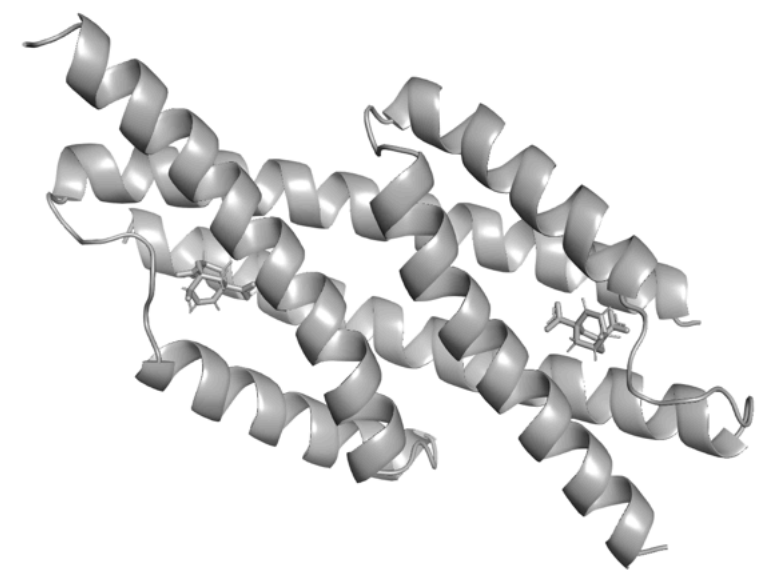
1. Learn P_{model} similar to P_{data}
2. Generate samples from P_{model}

Homologous sequences

Amino-acid chains



3d structure



Function

Catalyse a specific
reaction

Generative Models



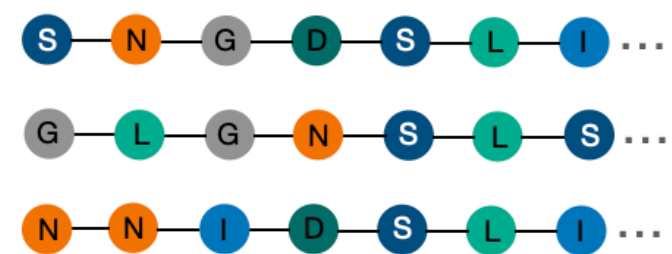
Training data $\sim P_{data}$

Artificial data $\sim P_{model}$

1. Learn P_{model} similar to P_{data}
2. Generate samples from P_{model}

Homologous sequences

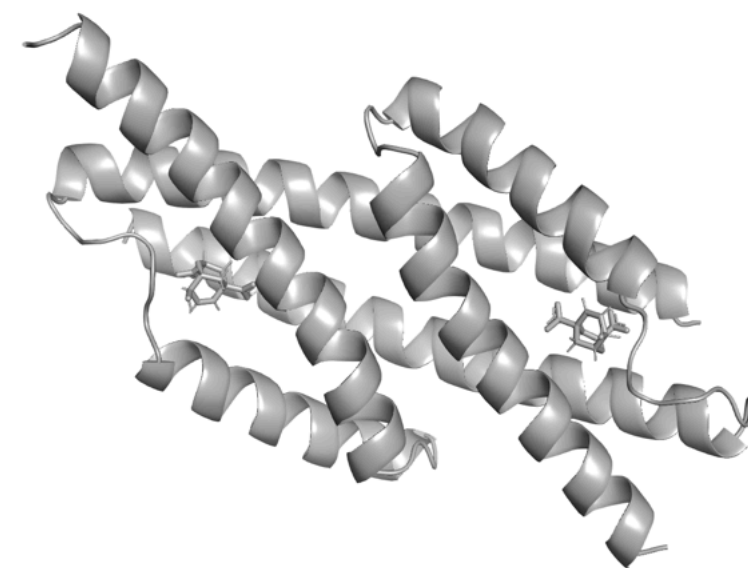
Amino-acid chains



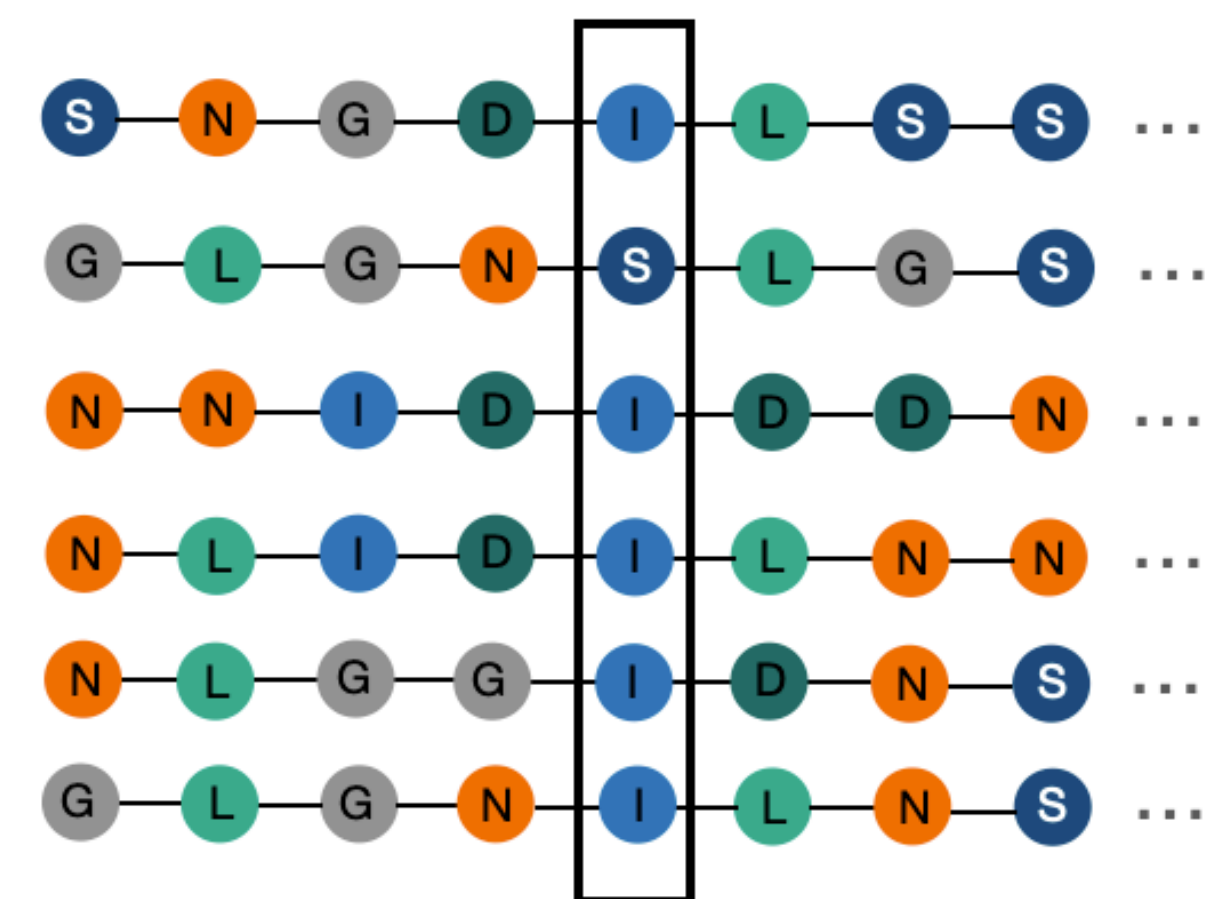
Function

Catalyse a specific reaction

3d structure



Conservation



Generative Models



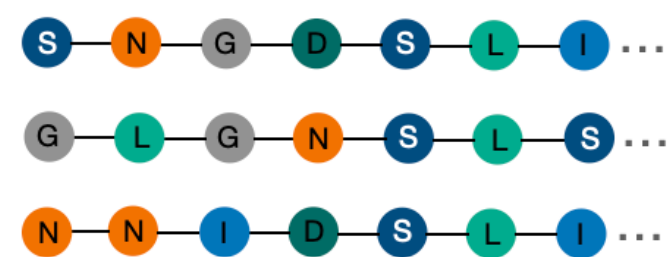
Training data $\sim P_{data}$

Artificial data $\sim P_{model}$

1. Learn P_{model} similar to P_{data}
2. Generate samples from P_{model}

Homologous sequences

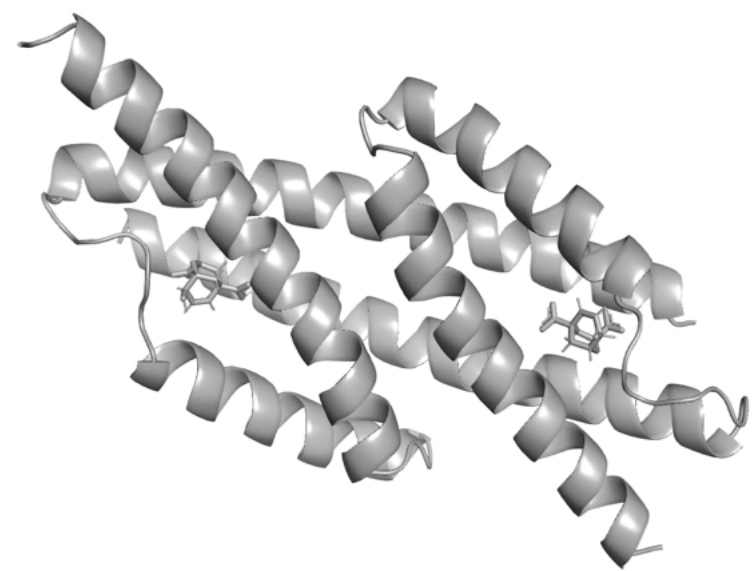
Amino-acid chains



Function

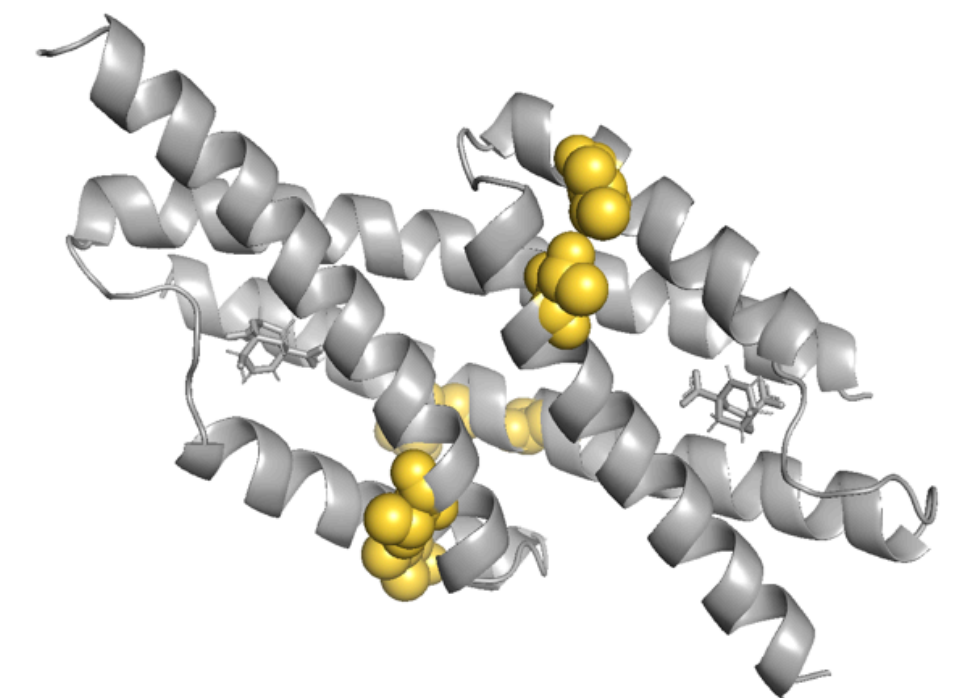
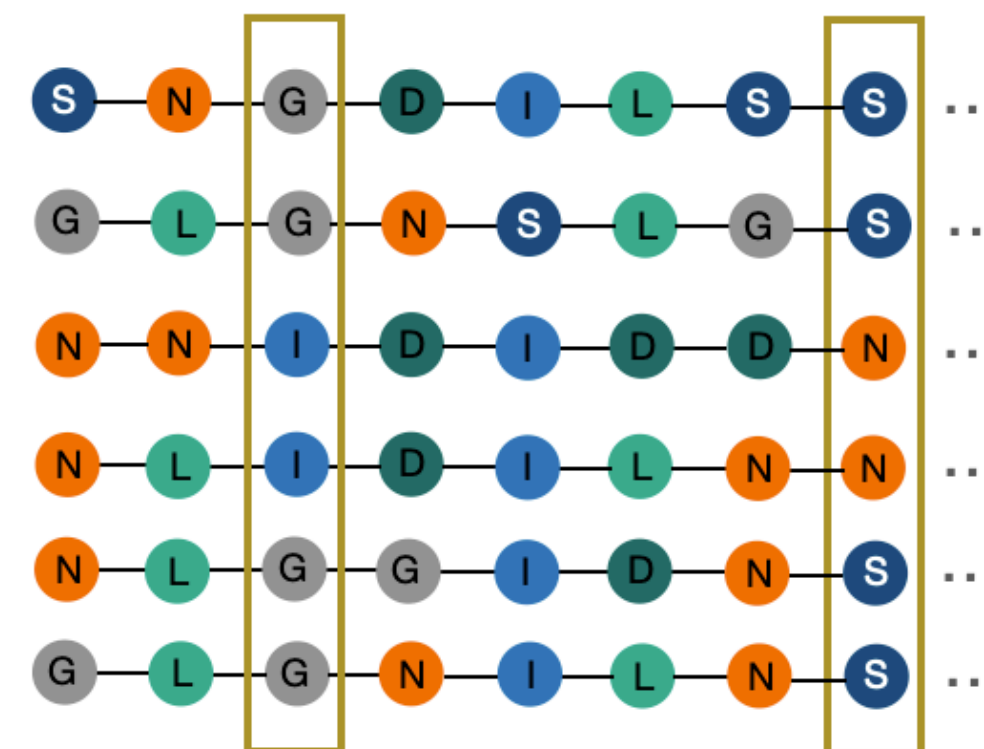
Catalyse a specific reaction

3d structure



Correlations

Contacts



Generative Models



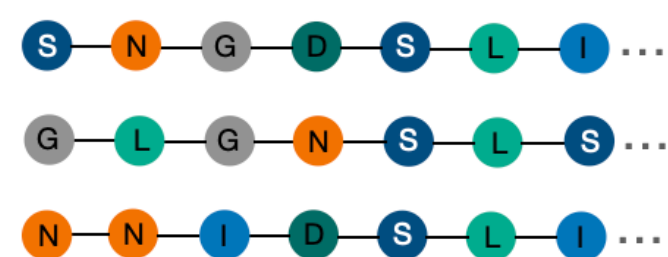
Training data $\sim P_{data}$

Artificial data $\sim P_{model}$

1. Learn P_{model} similar to P_{data}
2. Generate samples from P_{model}

Homologous sequences

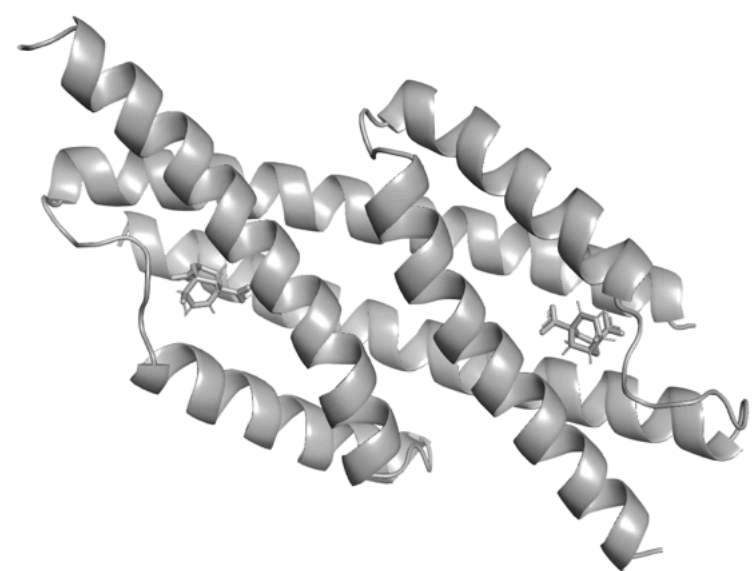
Amino-acid chains



Function

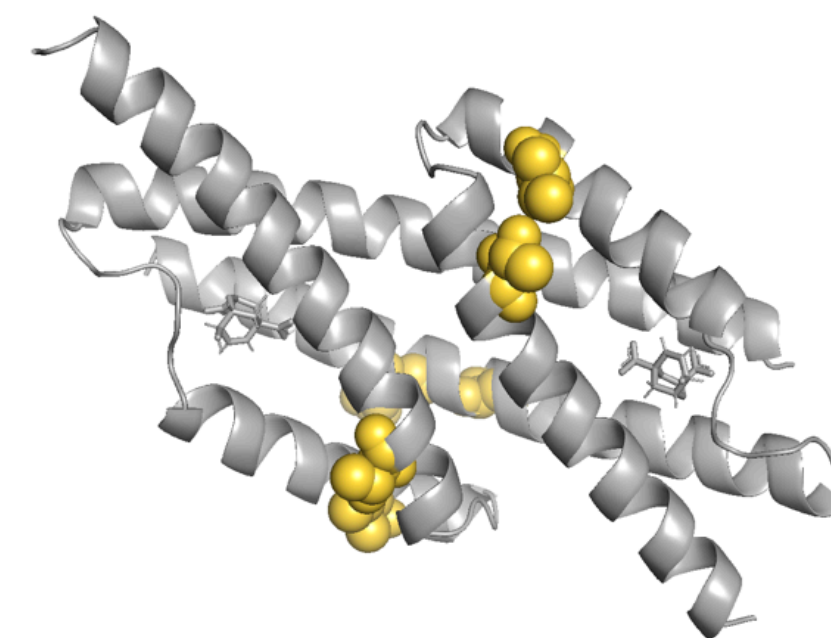
Catalyse a specific reaction

3d structure



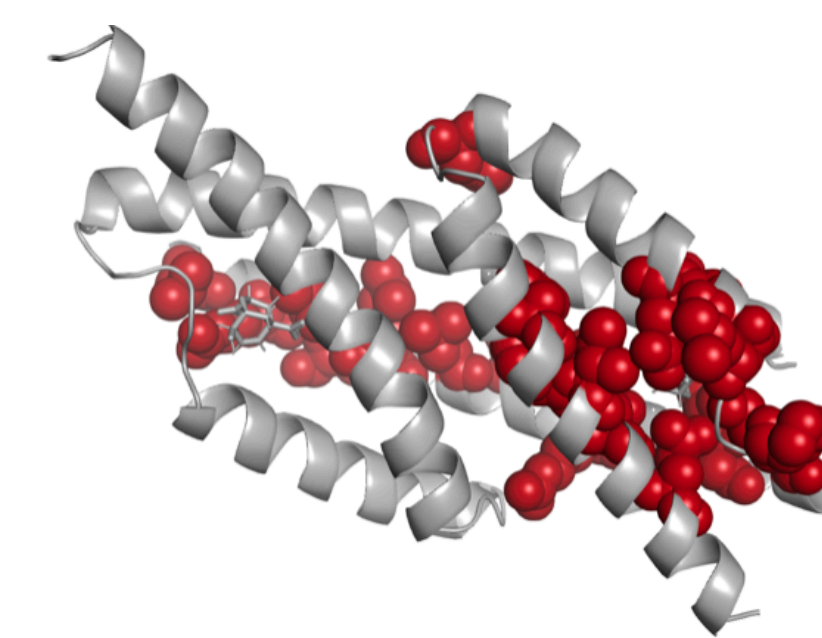
Correlations

Contacts



Local

Functional Positions



Collective

Generative Models



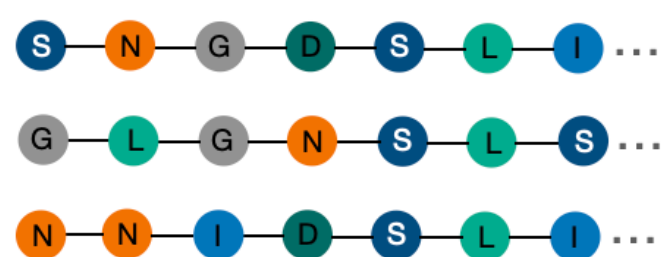
Training data $\sim P_{data}$

Artificial data $\sim P_{model}$

1. Learn P_{model} similar to P_{data}
2. Generate samples from P_{model}

Homologous sequences

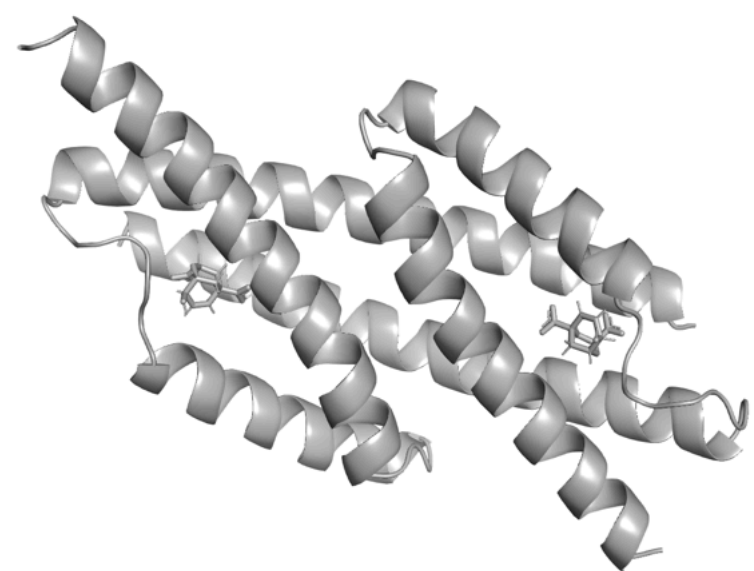
Amino-acid chains



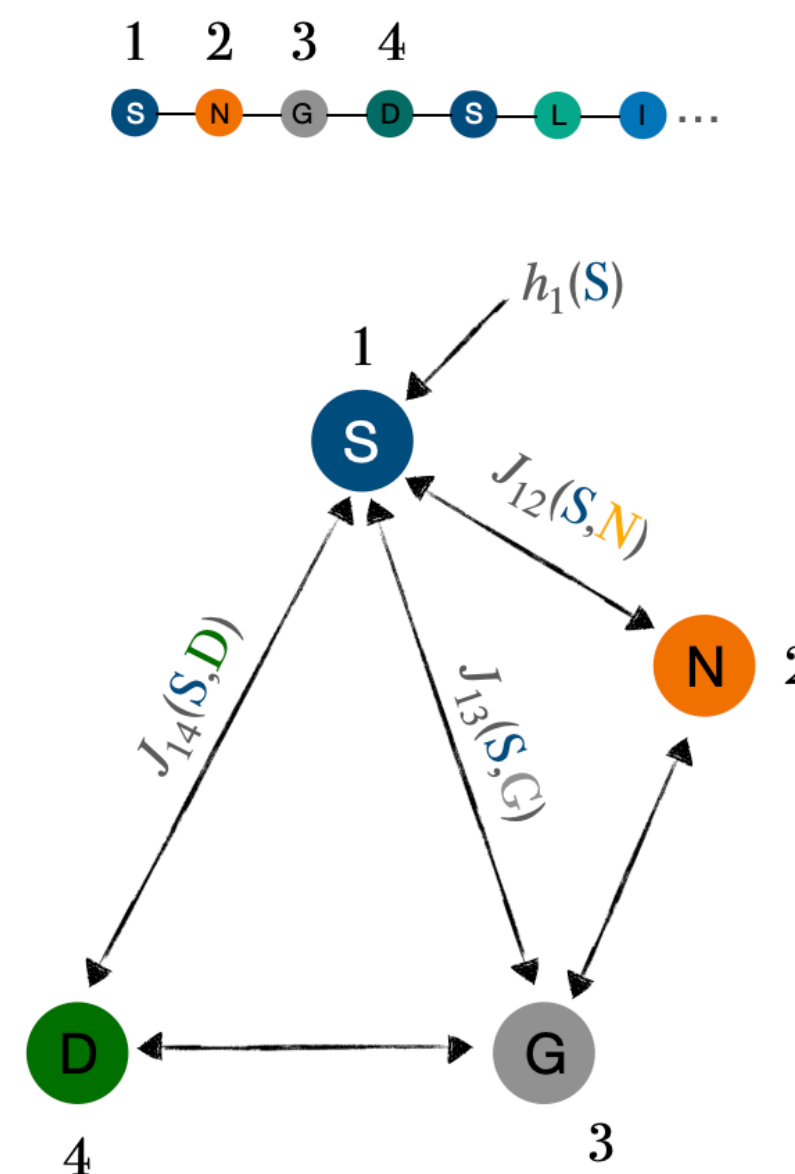
Function

Catalyse a specific reaction

3d structure



Potts Model

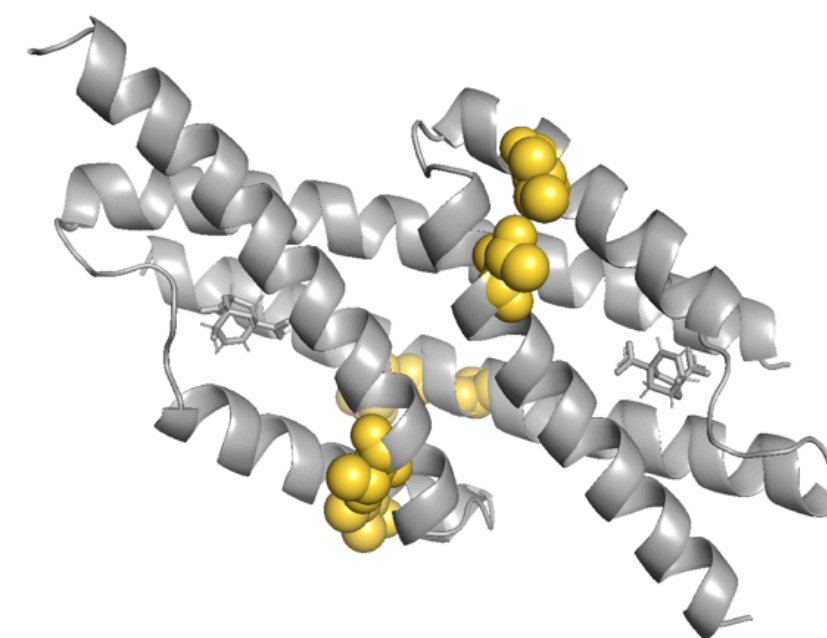


$$p(\{\sigma_i\}_{i=1,\dots,L}) = \frac{1}{Z(\mathbf{h}, \mathbf{J})} \prod_{i=1}^L e^{h_i(\sigma_i)} \prod_{i < j} e^{J_{ij}(\sigma_i, \sigma_j)}$$

- ☞ Parameters inferred with Gradient descent algorithm
- ☞ Boltzmann Machine algorithm (BM)

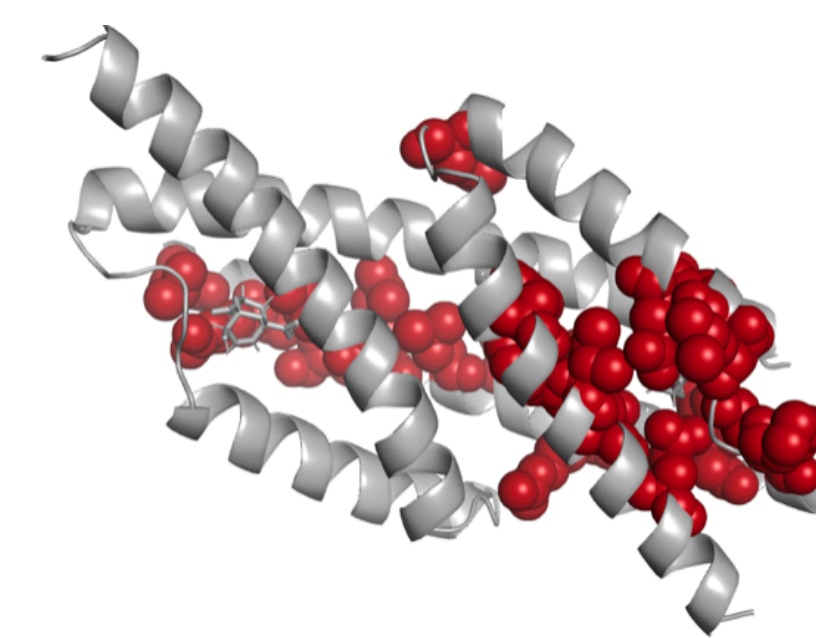
Correlations

Contacts



Local

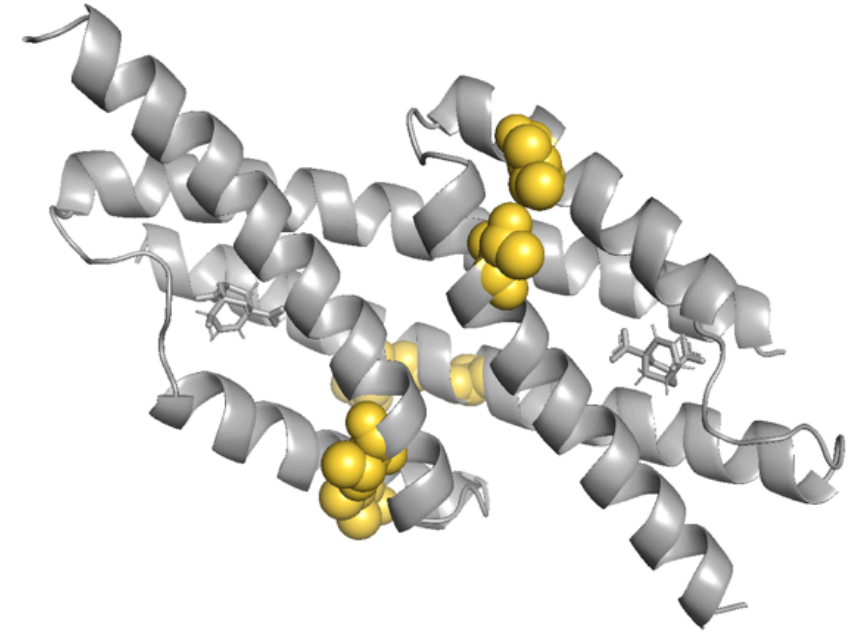
Functional Positions



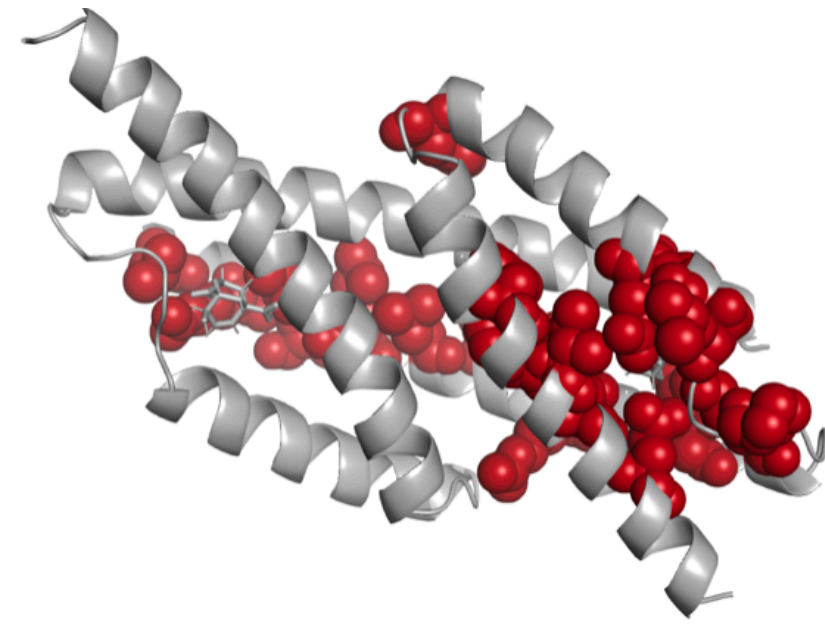
Collective

Different features

Contacts

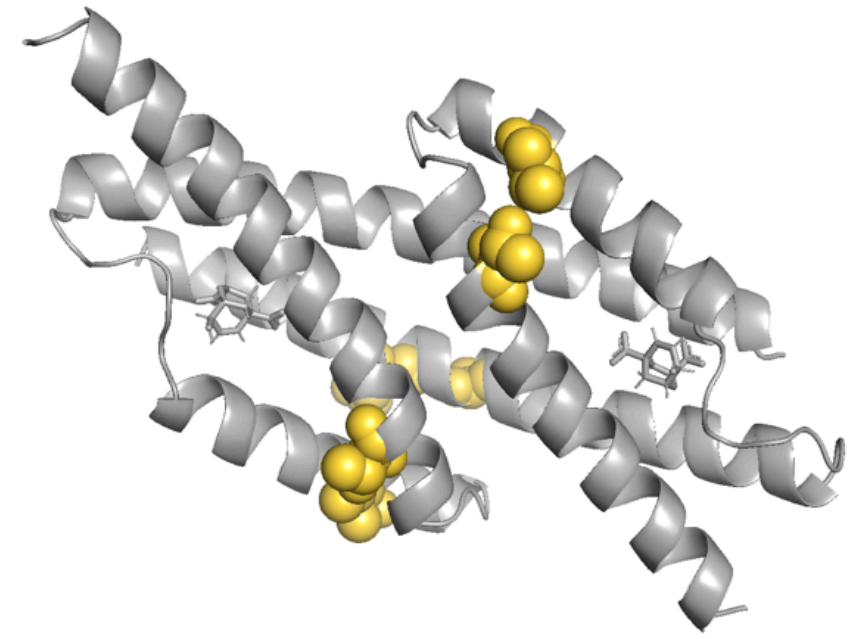


Functional Positions

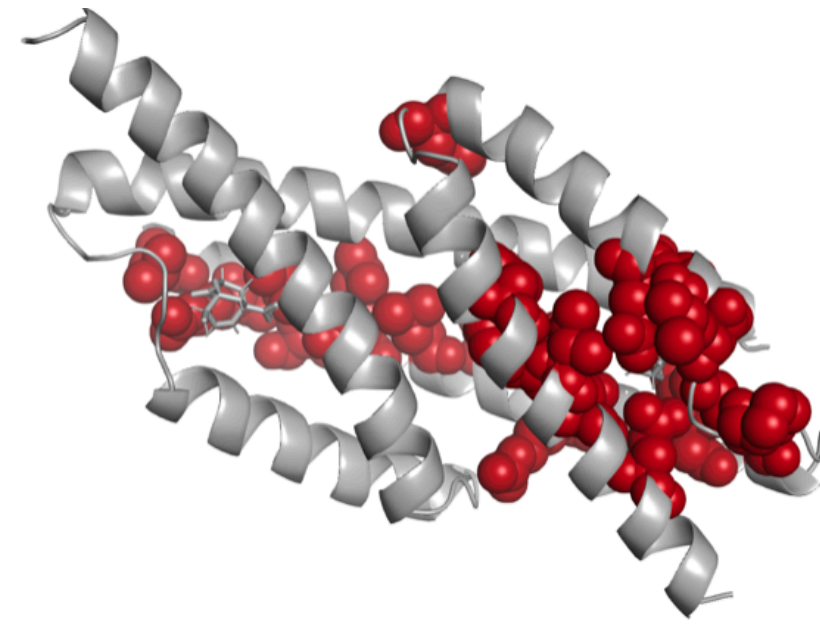


Different features

Contacts



Functional Positions



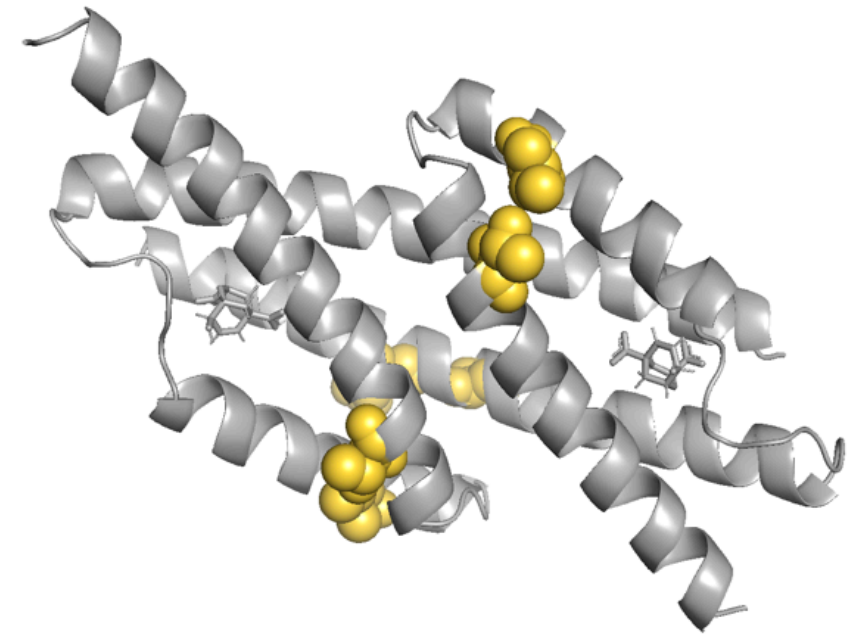
Undersampling

Number of natural sequences \ll $\begin{cases} \text{Number of possible sequences} \\ \text{Number of parameters} \end{cases}$

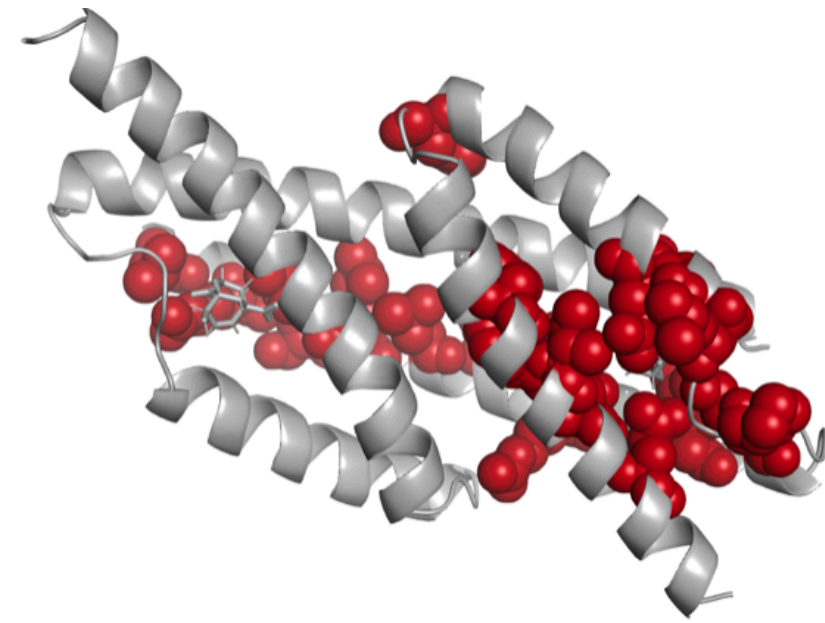
L2-Regularization: $\text{Loss} + \lambda \times \text{Penalty}$ L2-norm of the parameters

Different features

Contacts



Functional Positions



Undersampling

Number of natural sequences \ll $\begin{cases} \text{Number of possible sequences} \\ \text{Number of parameters} \end{cases}$

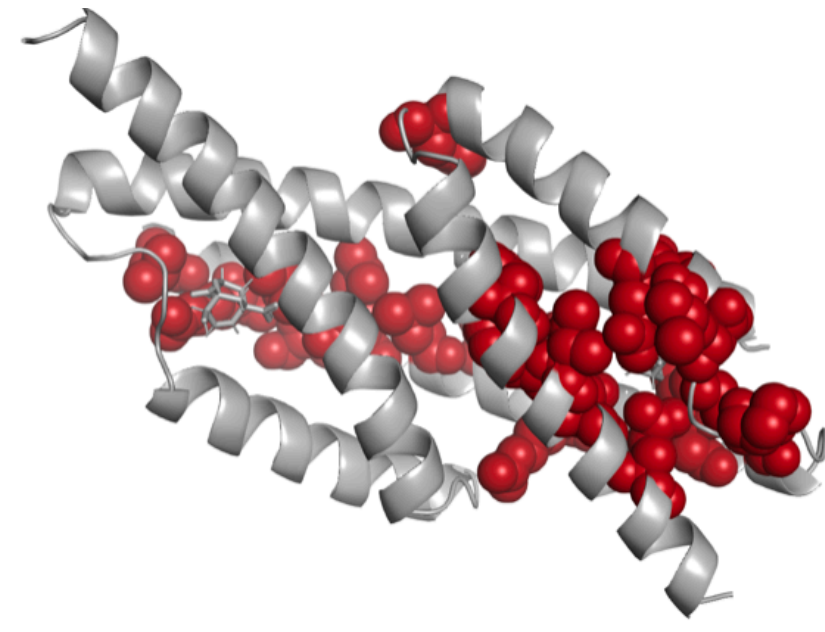
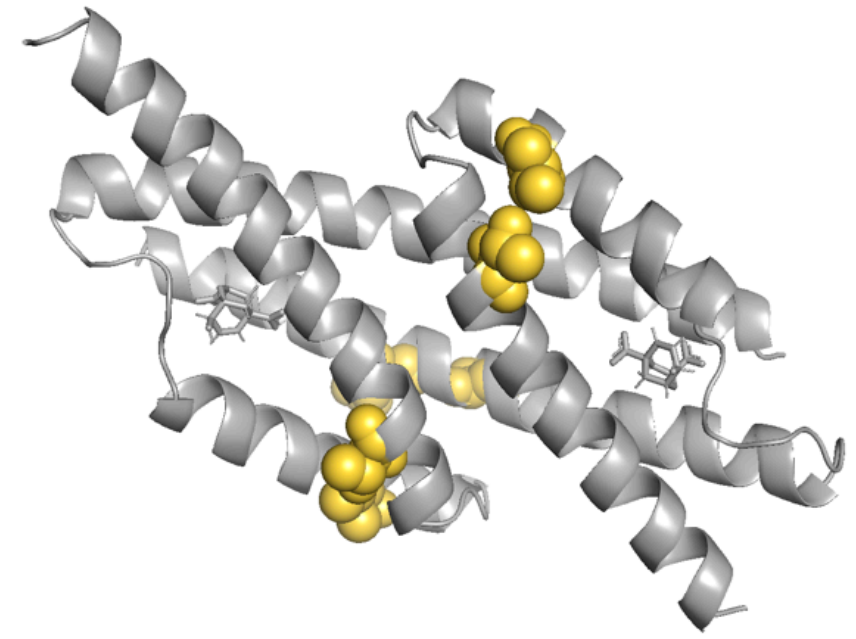
L2-Regularization: $\text{Loss} + \lambda \times \text{Penalty}$ L2-norm of the parameters

→ Introduce a bias

Different features

Contacts

Functional Positions



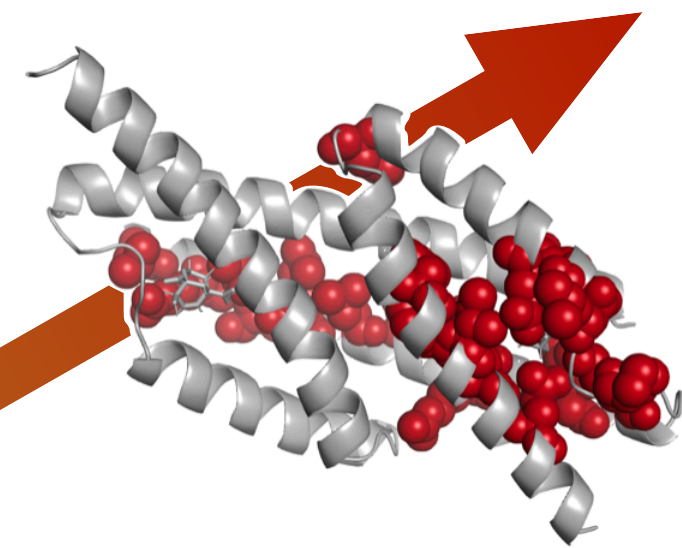
Undersampling

Number of natural sequences \ll $\begin{cases} \text{Number of possible sequences} \\ \text{Number of parameters} \end{cases}$

L2-Regularization: $\text{Loss} + \lambda \times \text{Penalty}$ L2-norm of the parameters

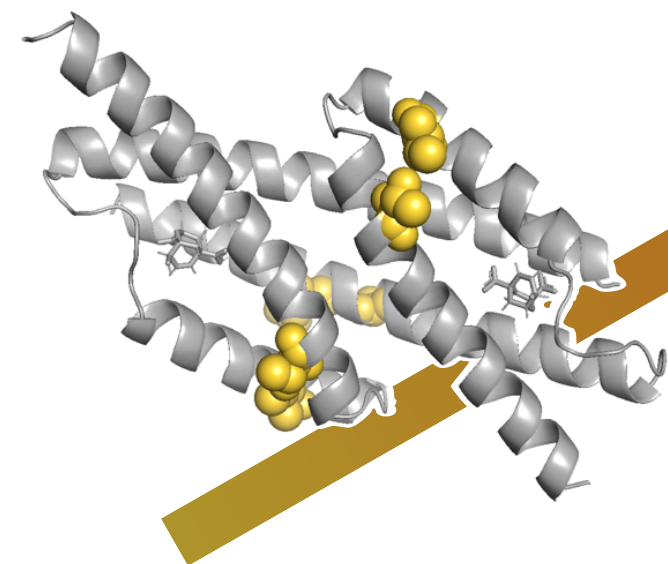
→ Introduce a bias

High regularization



Functional Positions

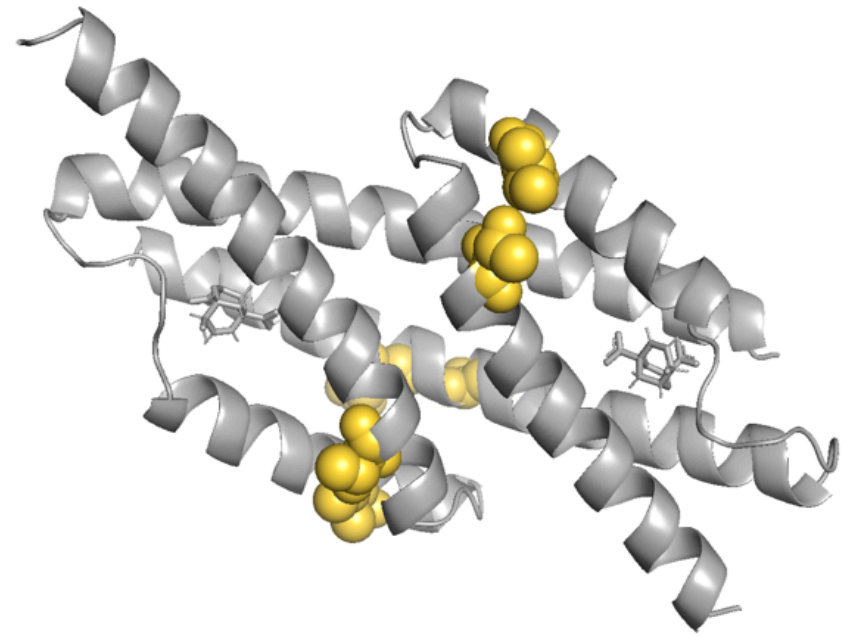
Contacts



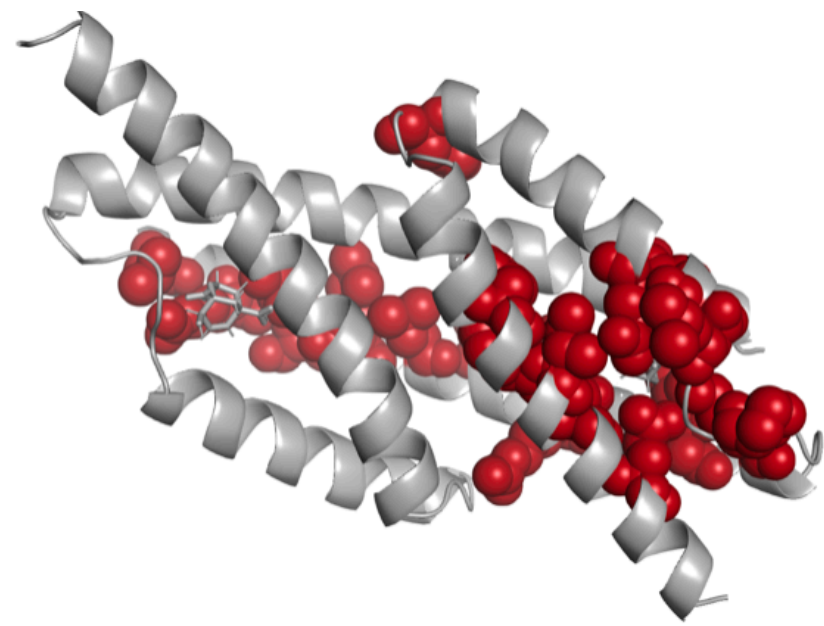
Low regularization

Different features

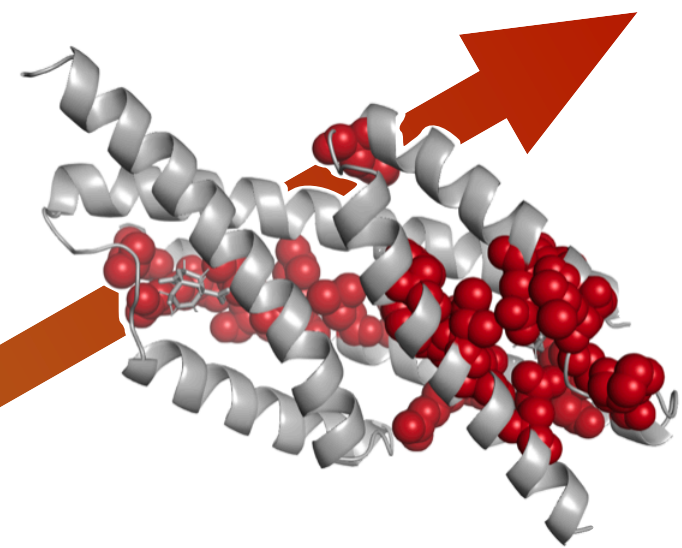
Contacts



Functional Positions

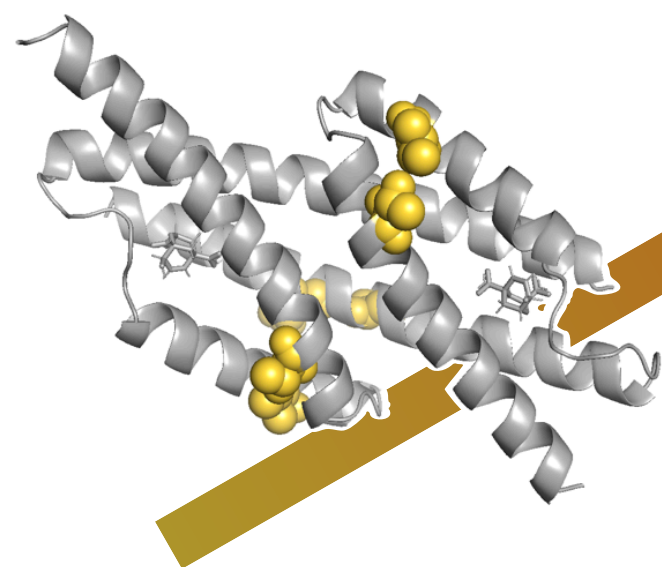


High regularization



Functional Positions

Contacts



Low regularization

Undersampling

Number of natural sequences \ll $\begin{cases} \text{Number of possible sequences} \\ \text{Number of parameters} \end{cases}$

L2-Regularization: $\text{Loss} + \lambda \times \text{Penalty}$ L2-norm of the parameters

→ Introduce a bias

Infer Generative Models that:

- Combine the inference of both **local** & **Collective** features
- Reproduce the diversity of natural protein families
- Capture other statistics
(*1st, 2nd, 3rd order statistics, PCA, energy distributions...*)

→ **SBM** (stochastic Boltzmann Machine) *