

# Modèles génératifs pour les séquences de protéines

## Generative models for protein sequences

Marion Chauveau

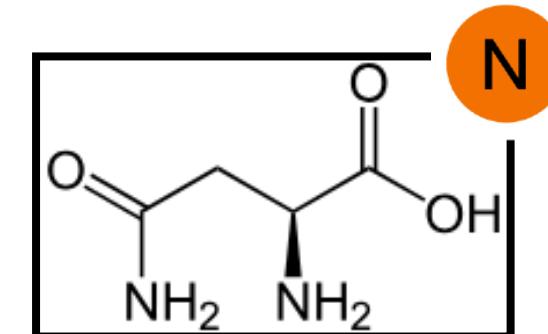
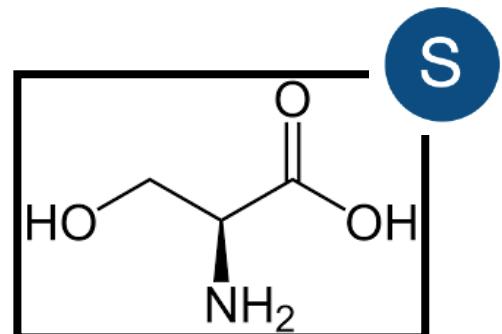
Thesis advisors: Ivan Junier and Olivier Rivoire

PhD Thesis Defense | October 1<sup>st</sup>, 2025 | ESPCI Paris

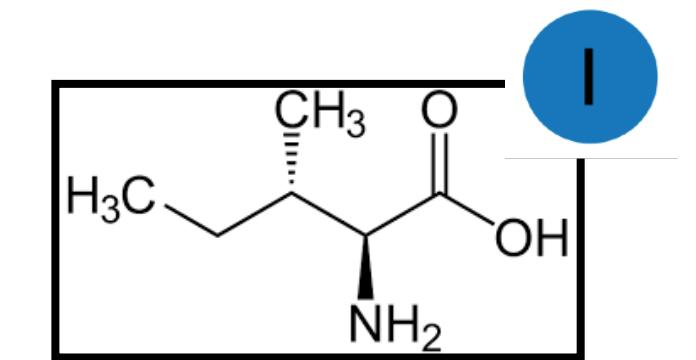
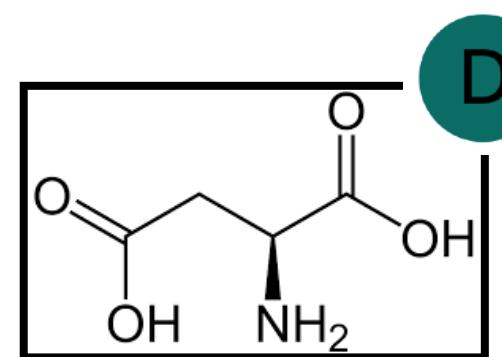
# Proteins

## The basics

- 20 building blocks: amino acids



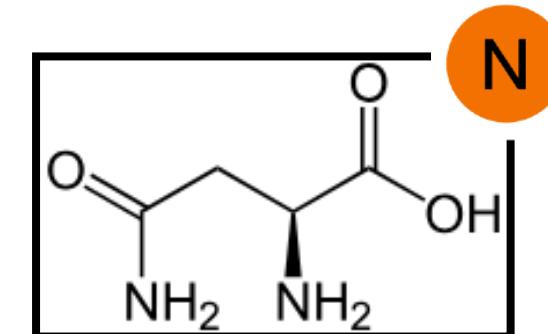
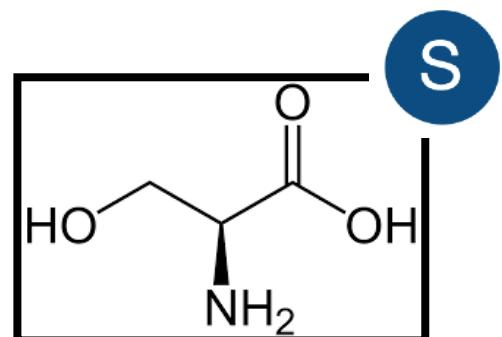
• • •



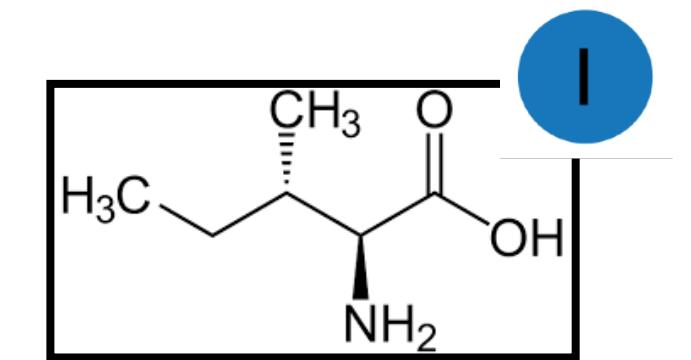
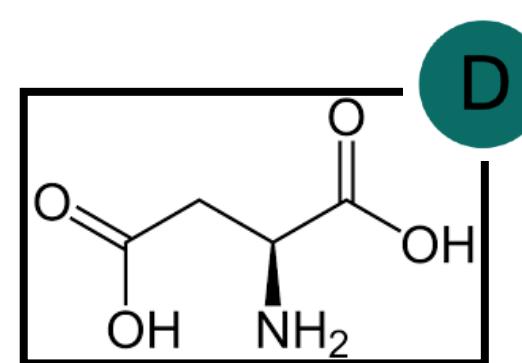
# Proteins

## The basics

- ▶ 20 building blocks: amino acids
- ▶ Assembled into chains
- ▶ Sequence encoded in the DNA



• • •



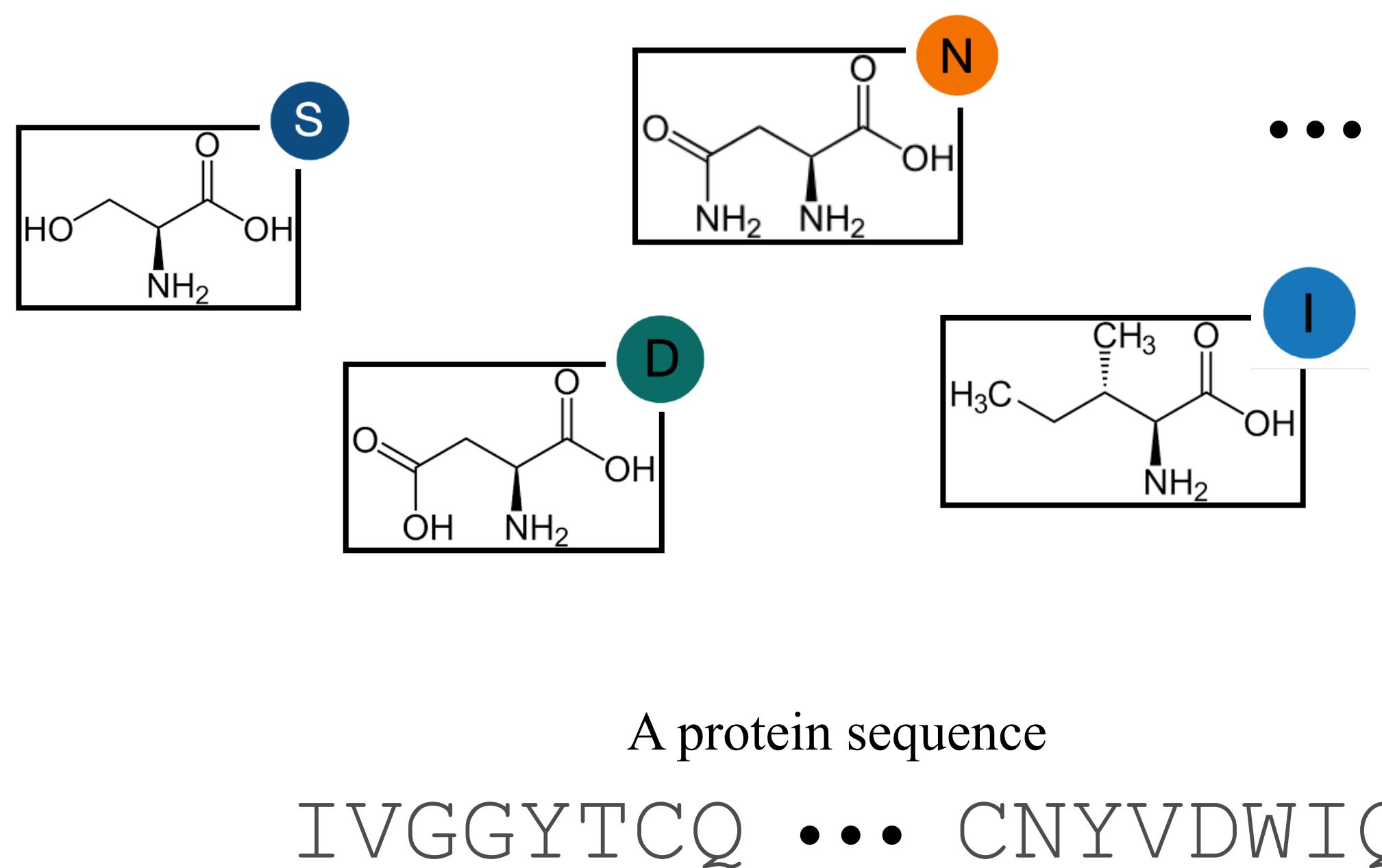
A protein sequence

IVGGYTCQ • • • CNYVDWIQ

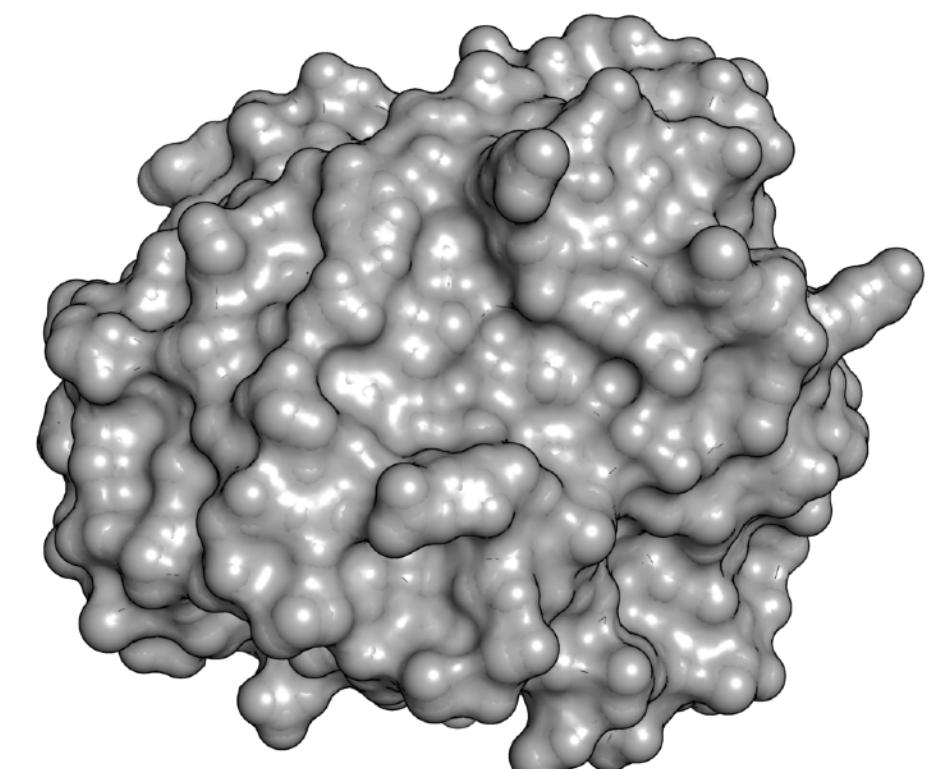
# Proteins

## The basics

- ▶ 20 building blocks: amino acids
- ▶ Assembled into chains
- ▶ Sequence encoded in the DNA
- ▶ Non-trivially ordered matter
- ▶ Life's essential machines (catalysis, transport, signaling, immune defense, structure...)



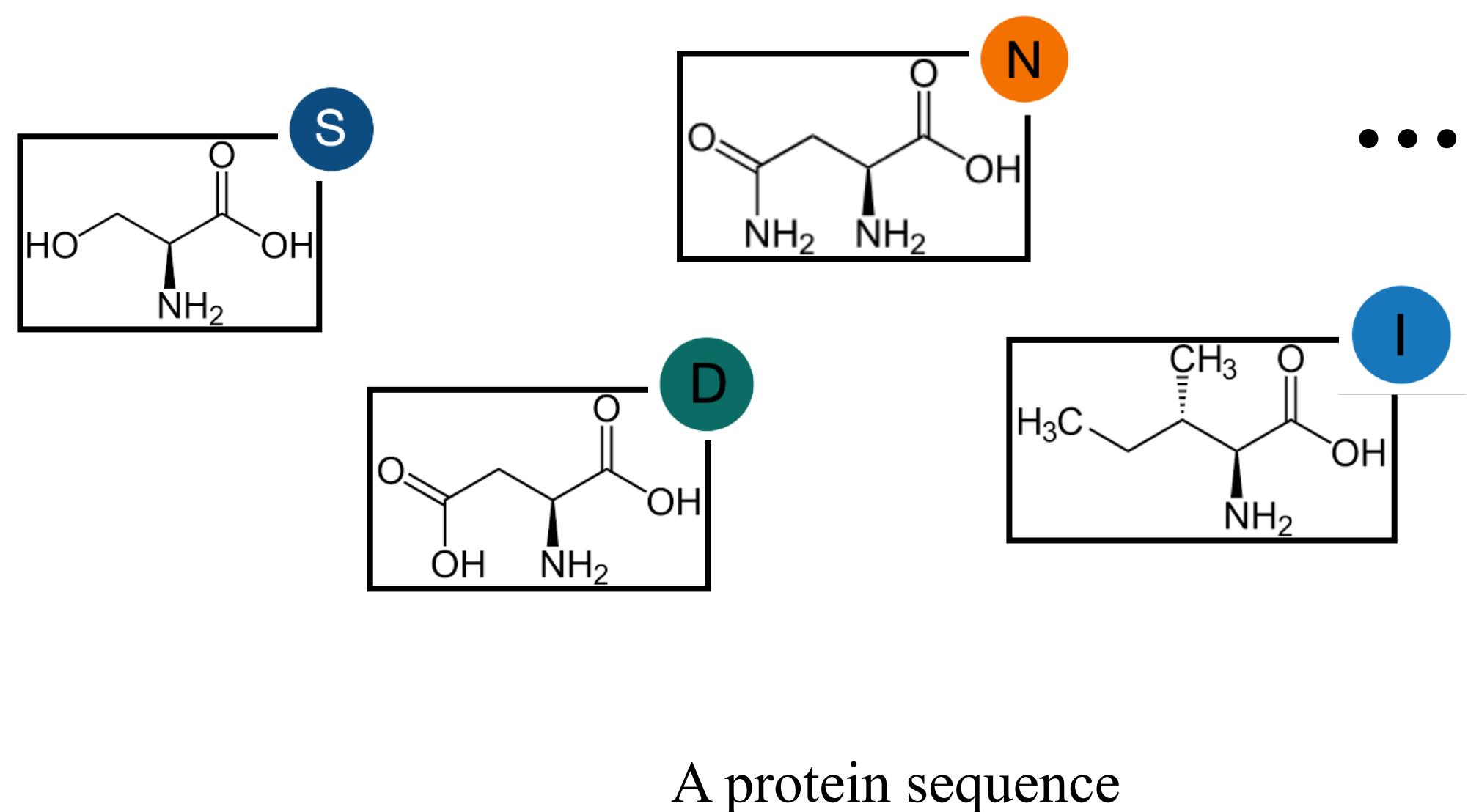
A protein structure representation



# Proteins

## The basics

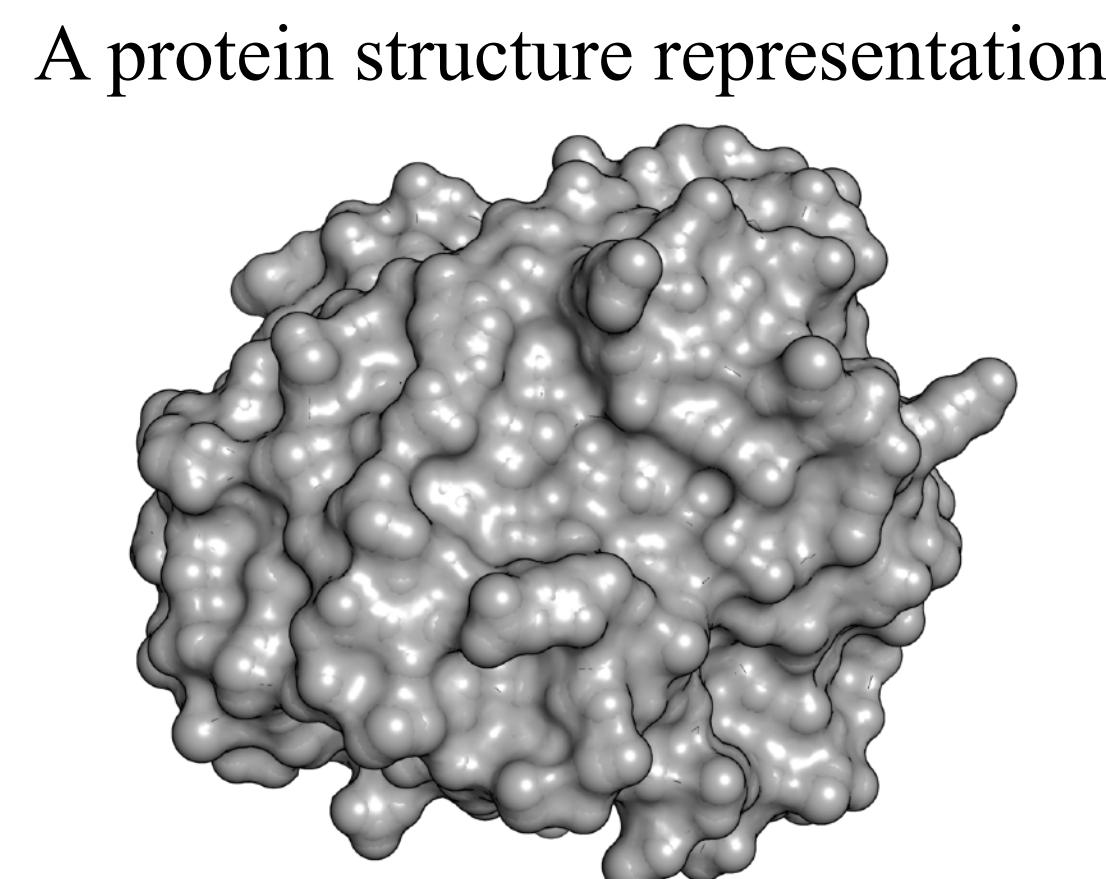
- ▶ 20 building blocks: amino acids
- ▶ Assembled into chains
- ▶ Sequence encoded in the DNA
- ▶ Non-trivially ordered matter
- ▶ Life's essential machines (catalysis, transport, signaling, immune defense, structure...)



## Performance & adaptability

- ▶ High catalytic efficiency (enzymes)
- ▶ Specific interaction with target molecules
- ▶ Capacity to evolve new functions

IVGGYTCQ ... CNYVDWIQ

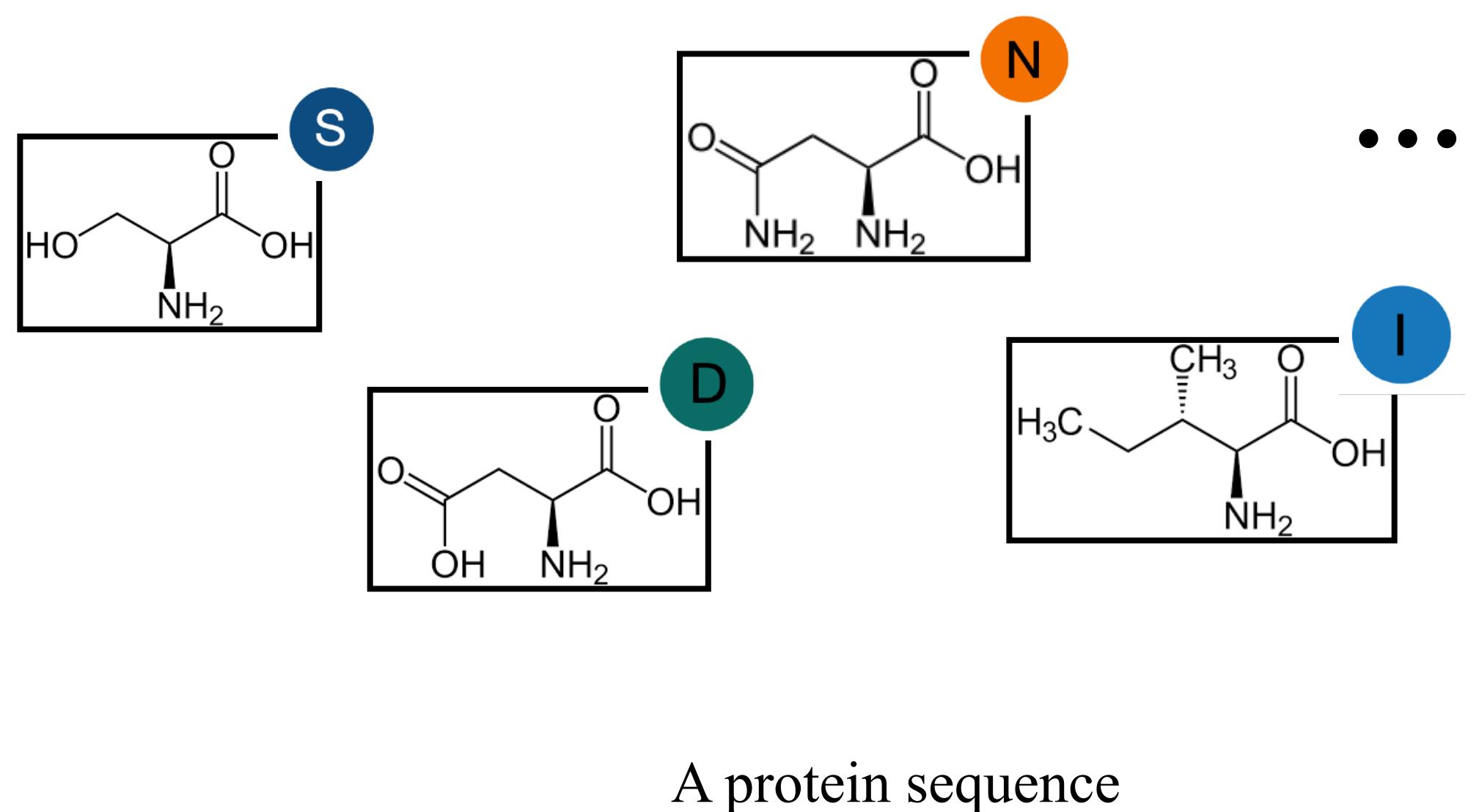


A protein structure representation

# Proteins

## The basics

- ▶ 20 building blocks: amino acids
- ▶ Assembled into chains
- ▶ Sequence encoded in the DNA
- ▶ Non-trivially ordered matter
- ▶ Life's essential machines (catalysis, transport, signaling, immune defense, structure...)

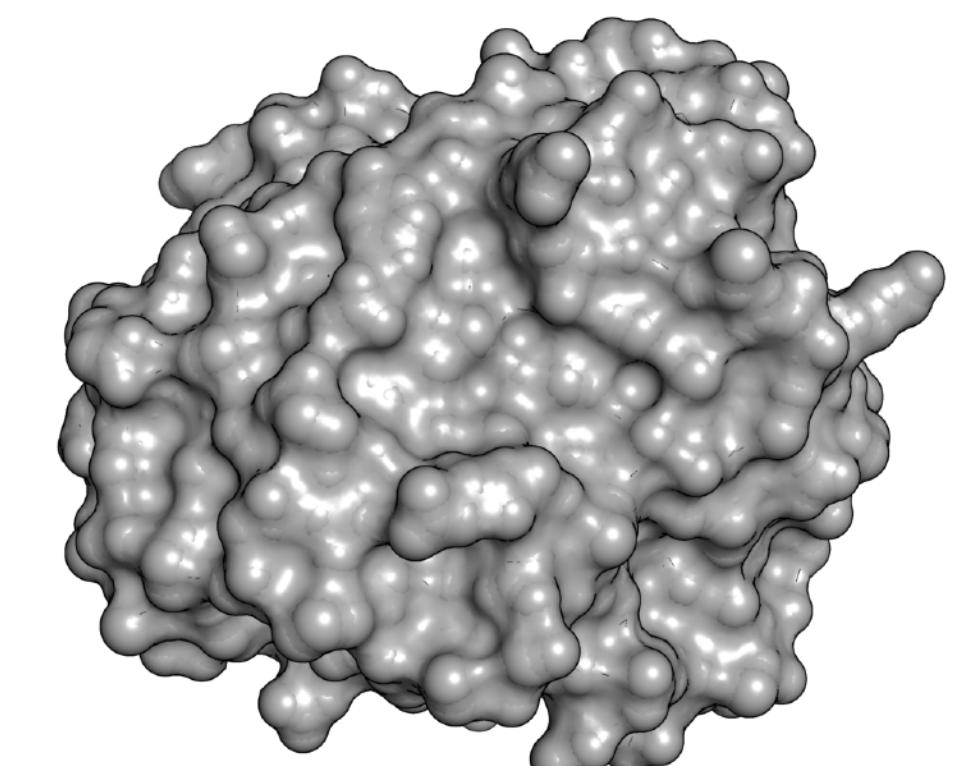


## Performance & adaptability

- ▶ High catalytic efficiency (enzymes)
- ▶ Specific interaction with target molecules
- ▶ Capacity to evolve new functions

IVGGYTCQ ... CNYVDWIQ

A protein structure representation



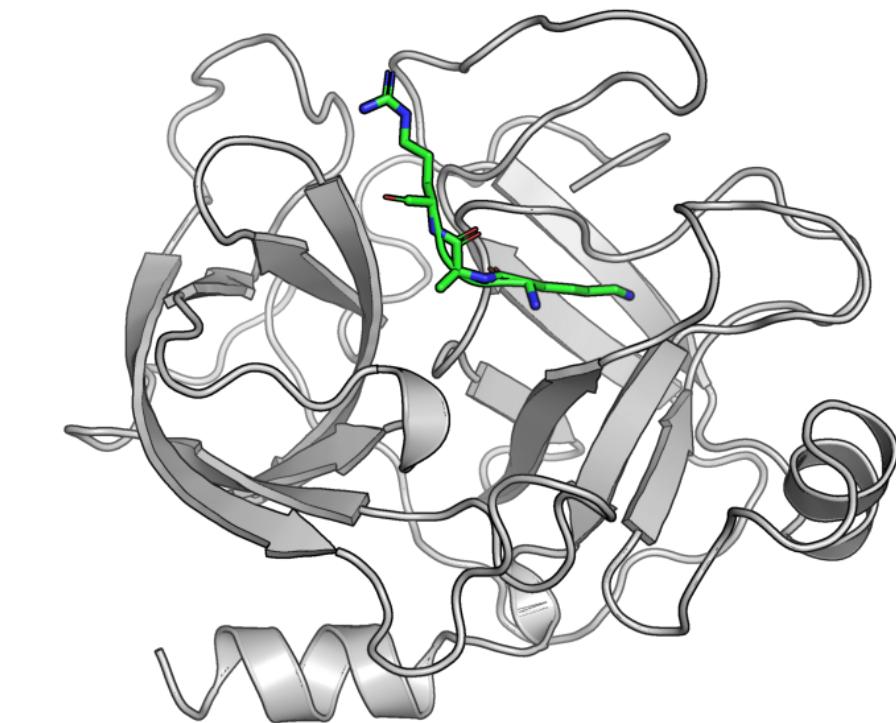
## Why are proteins worth studying?

- ▶ Understand how proteins work
- ▶ Applications in medicine, industry...

# Study proteins by learning from evolution

## Different approaches

- ▶ Observe their structure (X-ray crystallography...)
- ▶ Simulate them using structure (Molecular Dynamics...)
- ▶ Modify them (mutational scans...)
- ▶ Learn from evolution (statistical models...) Etc...

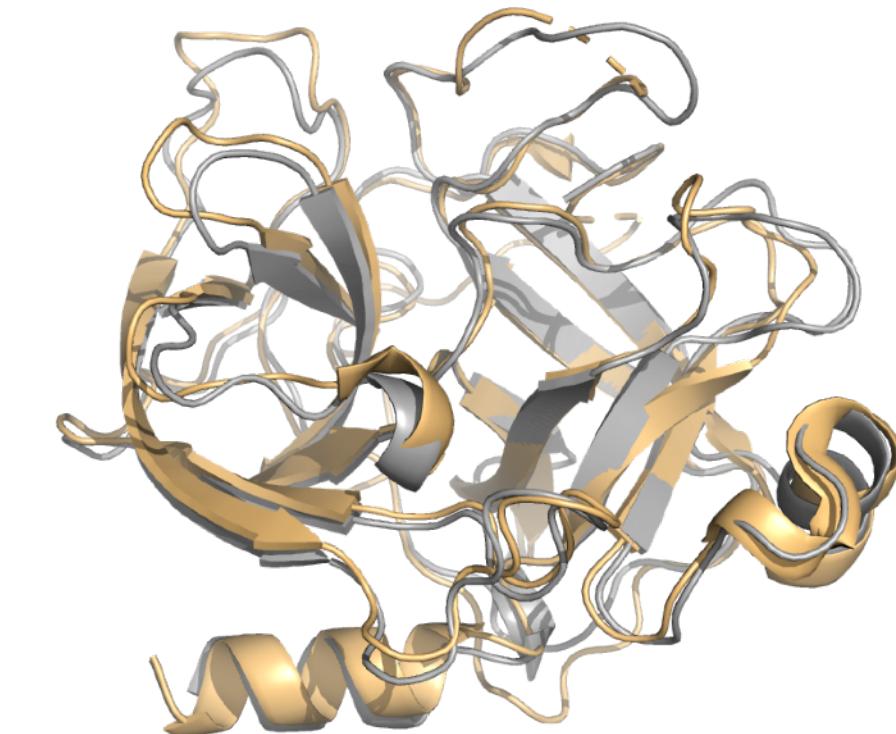


Rat Trypsin structure

# Study proteins by learning from evolution

## Different approaches

- ▶ Observe their structure (X-ray crystallography...)
- ▶ Simulate them using structure (Molecular Dynamics...)
- ▶ Modify them (mutational scans...)
- ▶ Learn from evolution (statistical models...) Etc...



## Why statistical models?

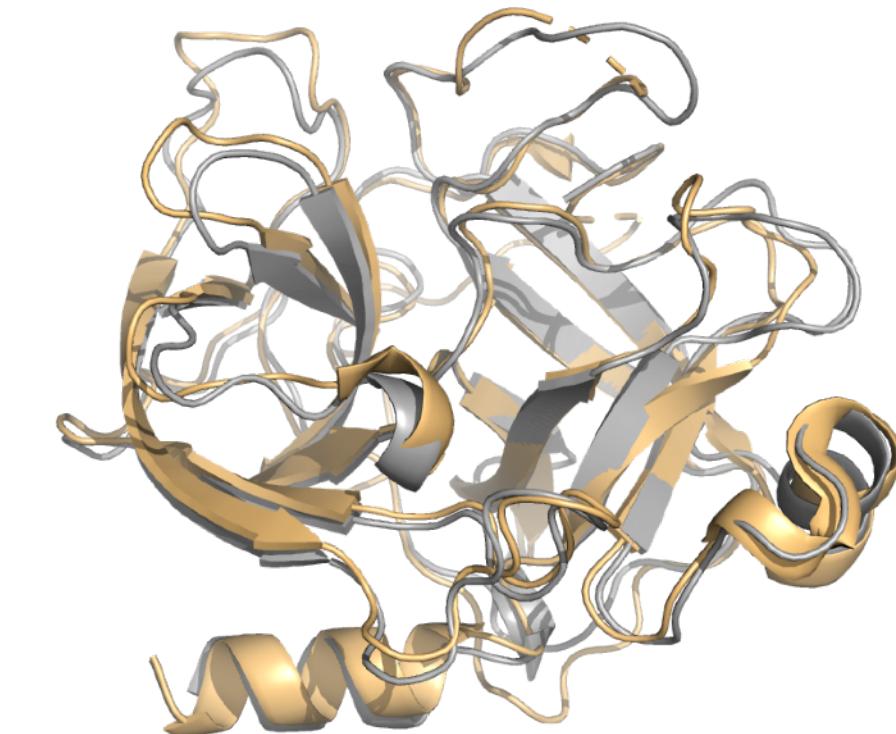
- ▶ Protein sequences are shaped by evolutionary pressure
- ▶ Many sequences can have similar structure and function: family

 IVGGYTCQ ... CNYVDWIQ  
 IVGGR--R ... AQFVNWID

# Study proteins by learning from evolution

## Different approaches

- ▶ Observe their structure (X-ray crystallography...)
- ▶ Simulate them using structure (Molecular Dynamics...)
- ▶ Modify them (mutational scans...)
- ▶ Learn from evolution (statistical models...) Etc...



## Why statistical models?

- ▶ Protein sequences are shaped by evolutionary pressure
- ▶ Many sequences can have similar structure and function: family

IVGGYTCQ ... CNYVDWIQ  
 IVGGR--R ... AQFVNWID

Multiple Sequence Alignment

## How?

- ▶ Collect sequences with shared structure and function
- ▶ Build a Multiple Sequence Alignment (MSA)
- ▶ Look for statistical signatures

Sequences	Positions		260
	1	5	
1	IVGGYTCQ	...	CNYVDWIQ
2	IVGGR--R	...	AQFVNWID
3	IIGGH-AK	...	STFLSWIK
	:		:
	ITNGAYDG	...	TSQIINWIR
	VNGNFDG	...	GLYSGWIQ

# Study proteins by learning from evolution

## Different approaches

- ▶ Observe their structure (X-ray crystallography...)
- ▶ Simulate them using structure (Molecular Dynamics...)
- ▶ Modify them (mutational scans...)
- ▶ **Learn from evolution (statistical models...)** Etc...

## Why statistical models?

- ▶ Protein sequences are shaped by evolutionary pressure
- ▶ Many sequences can have similar structure and function: family

## How?

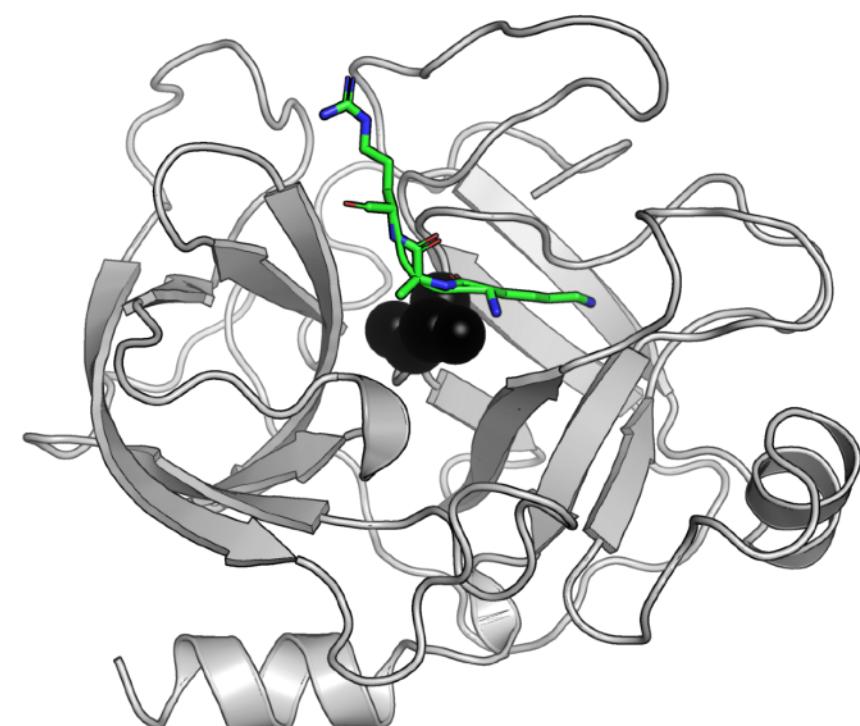
- ▶ Collect sequences with shared structure and function
- ▶ Build a Multiple Sequence Alignment (MSA)
- ▶ Look for statistical signatures

## To find

- ▶ Conservations (important amino-acids)
- ▶ Correlations (important interactions)

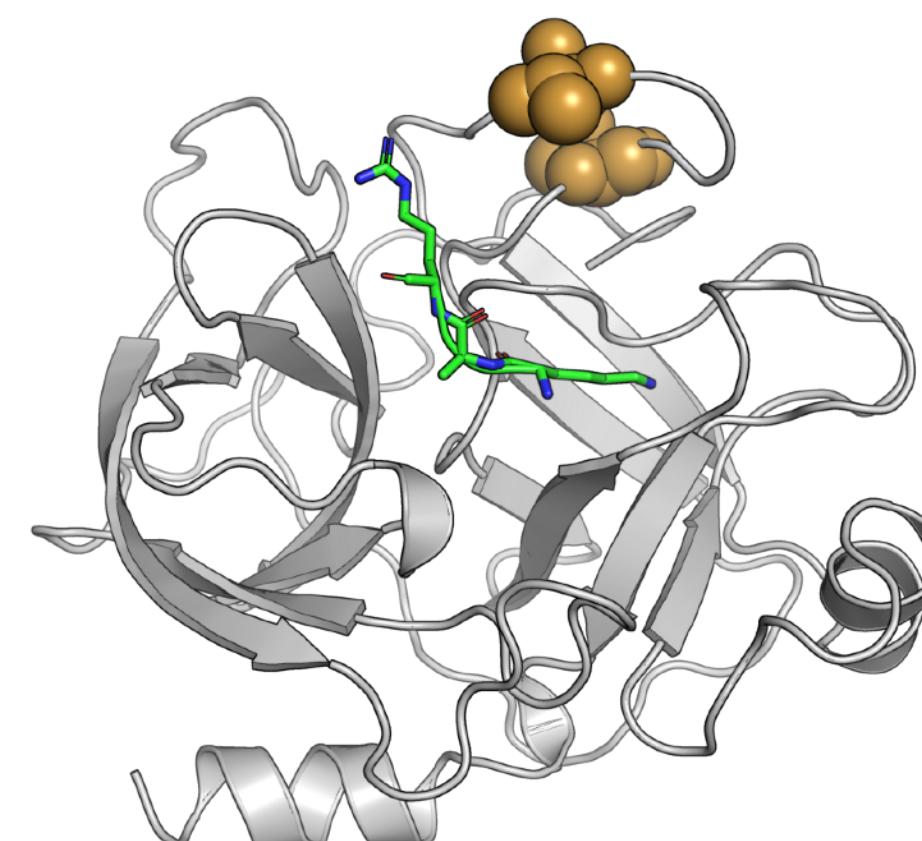
## Conservation

	Positions
Sequences	
1	• • • CQGD SGGP • • •
2	• • • CNGD SGAS • • •
3	• • • CTGD SGGP • • •
4	• • • GYGD SGGP • • •
5	• • • CQGD SGGP • • •



## Correlations

	Positions
Sequences	
1	• • • S N G D I Y S S • • •
2	• • • G L G N S F S G • • •
3	• • • N N I D I D N D • • •
4	• • • N L I D I L N N • • •
5	• • • N L G G I D S R • • •



Kendrew *et al.*, Nature, 1958

C Edgar *et al.*, Current opinion in structural biology, 2006

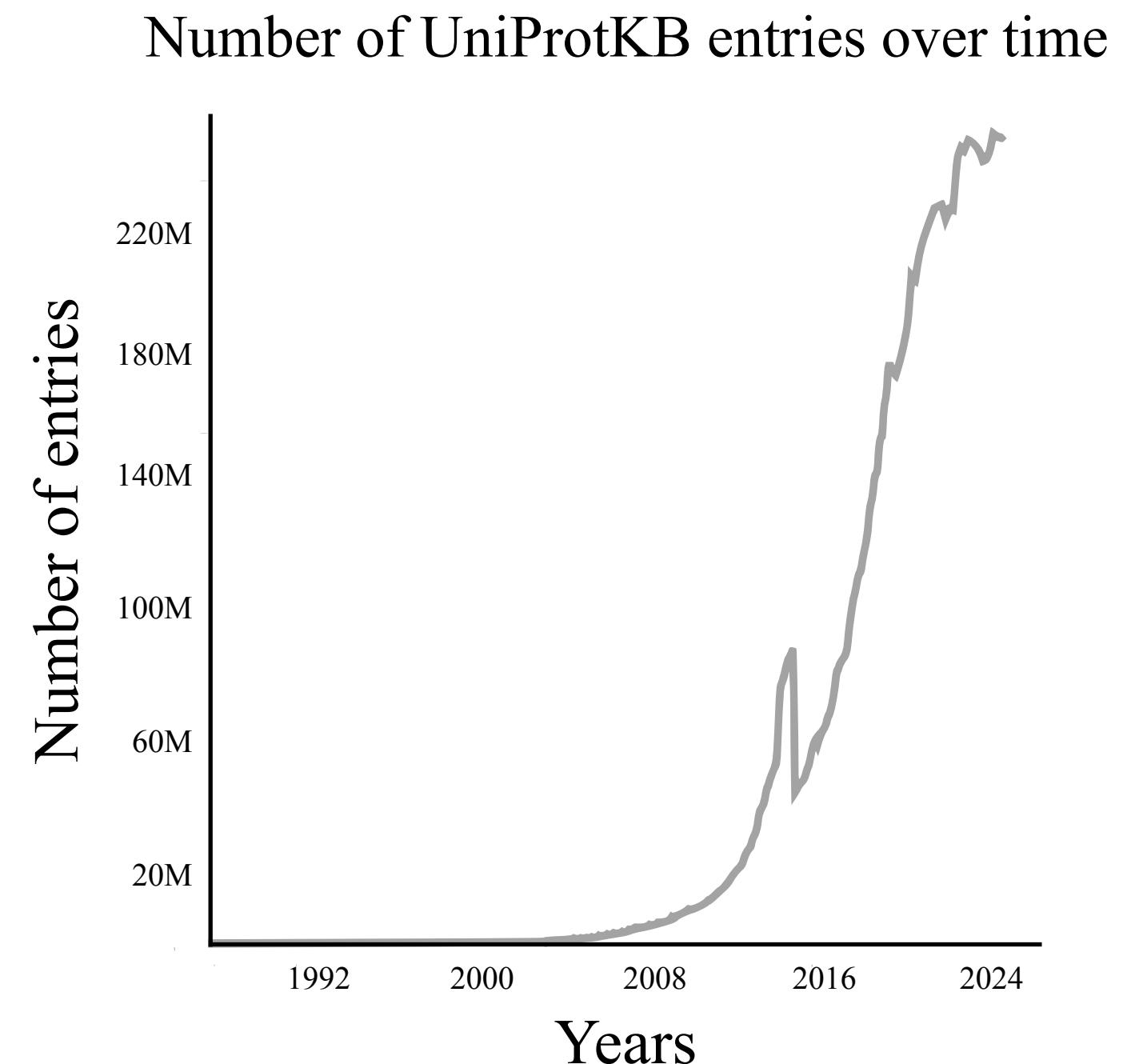
M Fowler *et al.*, Nature methods, 2014

Cocco *et al.*, Rep. Prog. Phys., 2018

# Statistical learning on a protein family

## Enough data for advanced methods

- ▶ Massive expansion of available sequences in the last two decades
- ▶ Approaches grounded in statistical physics
- ▶ Deep learning methods

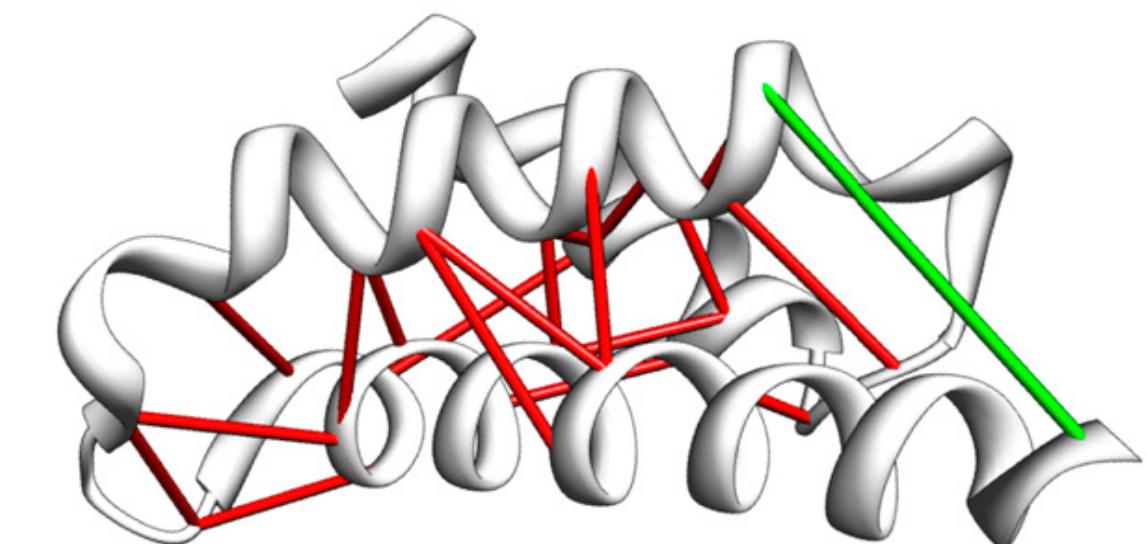


# Statistical learning on a protein family

## Enough data for advanced methods

- ▶ Massive expansion of available sequences in the last two decades
- ▶ Approaches grounded in statistical physics
- ▶ Deep learning methods

Contact prediction

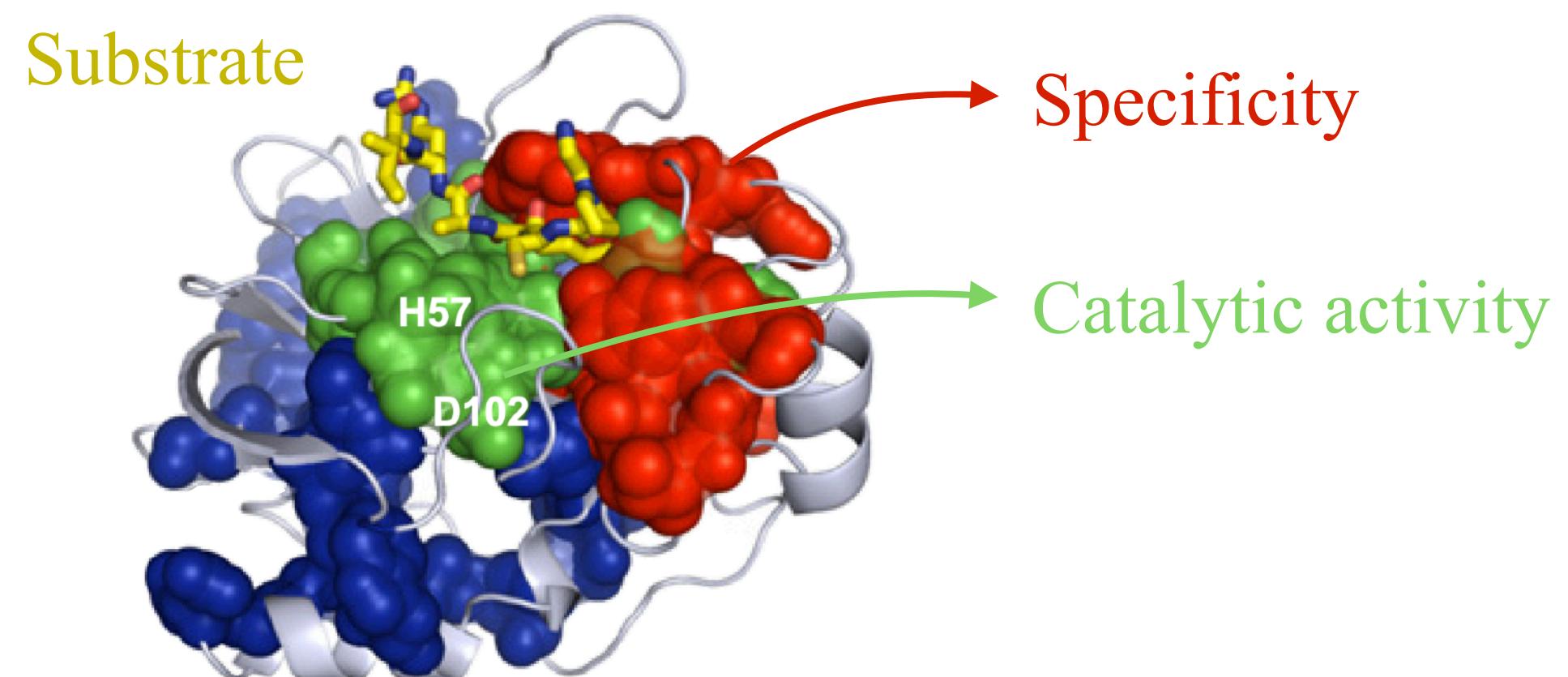


Morcos F *et al.*, PNAS, 2011

## Different methods capture different interaction scales

- ▶ From contacting pairs (ex: Direct Coupling Analysis)
- ▶ To coevolving groups (ex: Statistical Coupling Analysis)

S1A sectors mapped onto rat trypsin structure



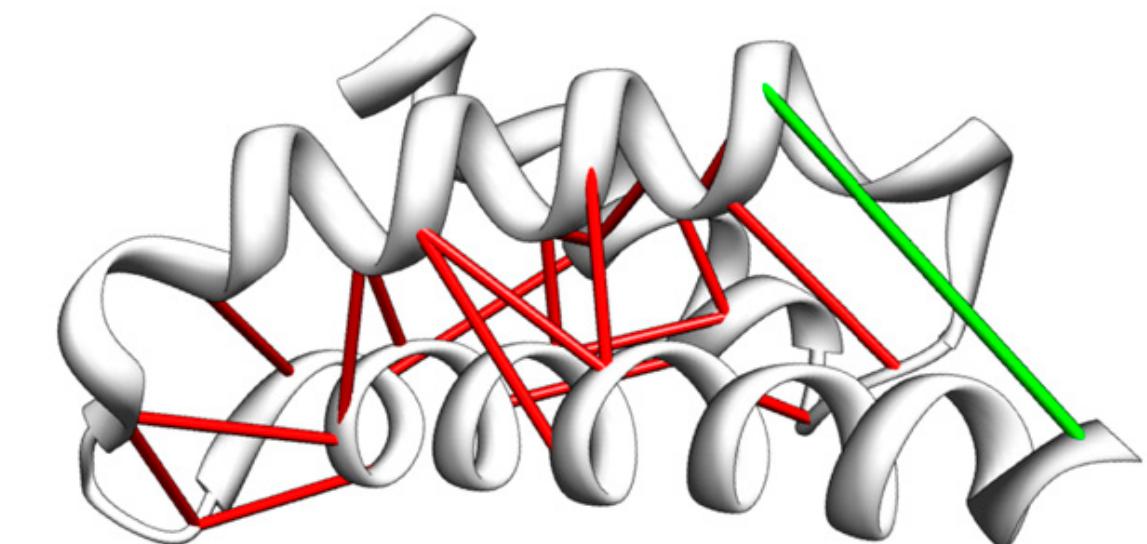
Halabi *et al.*, Cell, 2009

# Statistical learning on a protein family

## Enough data for advanced methods

- ▶ Massive expansion of available sequences in the last two decades
- ▶ Approaches grounded in statistical physics
- ▶ Deep learning methods

Contact prediction

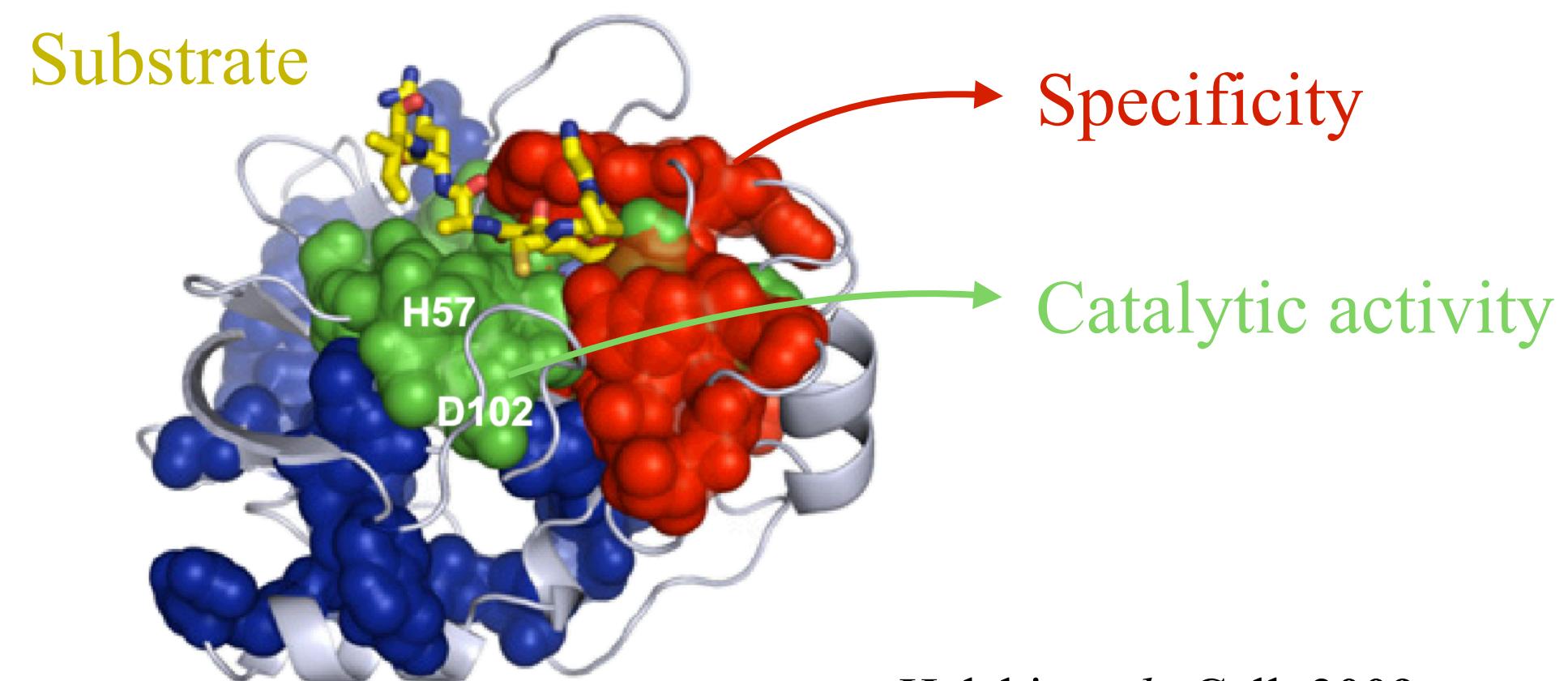


Morcos F *et al.*, PNAS, 2011

## Different methods capture different interaction scales

- ▶ From contacting pairs (ex: Direct Coupling Analysis)
- ▶ To coevolving groups (ex: Statistical Coupling Analysis)

S1A sectors mapped onto rat trypsin structure



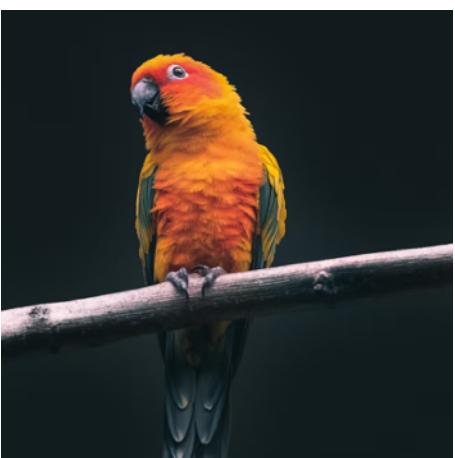
Halabi *et al.*, Cell, 2009

## Statistical Coupling Analysis applied to S1A family

- ▶ 3 coevolving groups (sectors)
- ▶ Structurally connected
- ▶ Functionally independent (mutagenesis experiments)

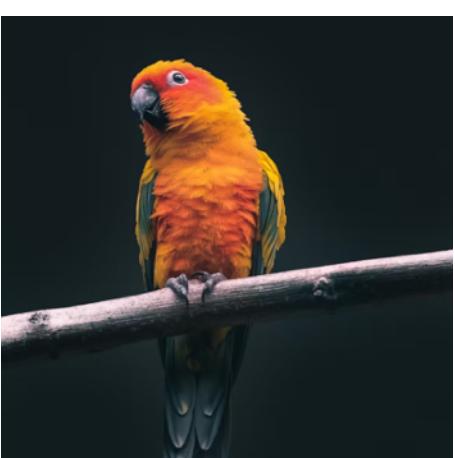
# Generative models for protein sequences

Training data

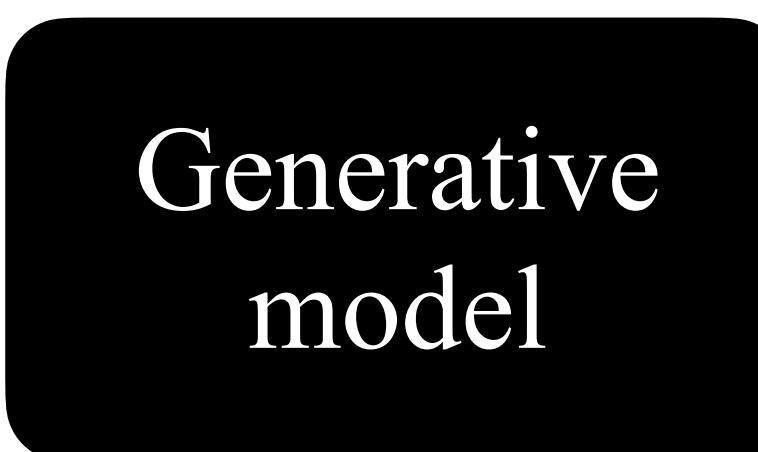


# Generative models for protein sequences

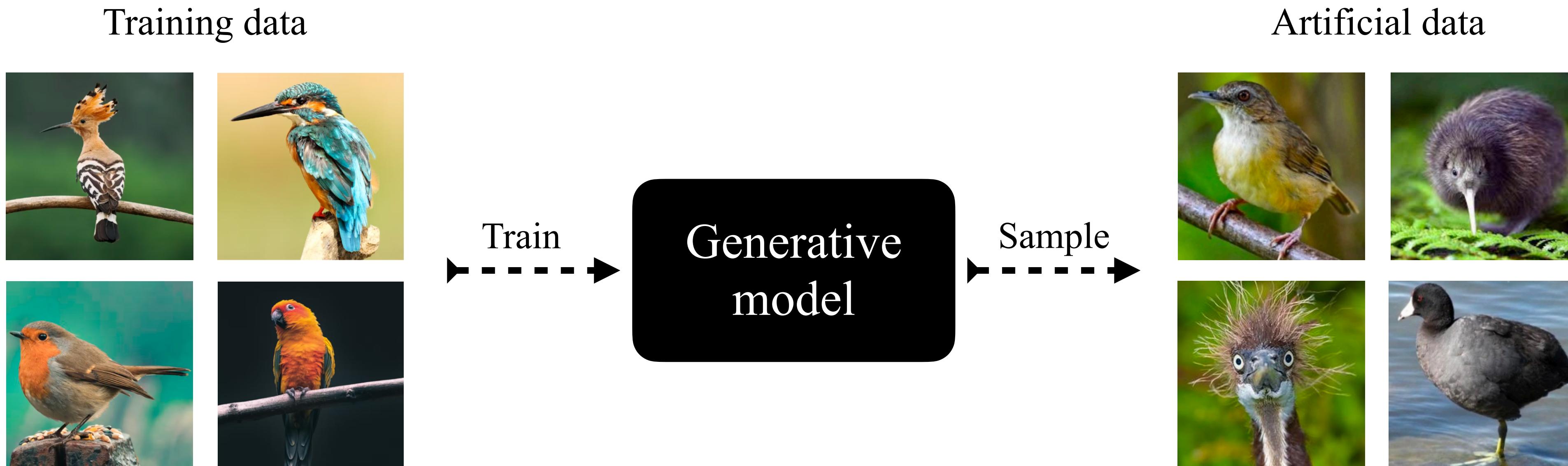
Training data



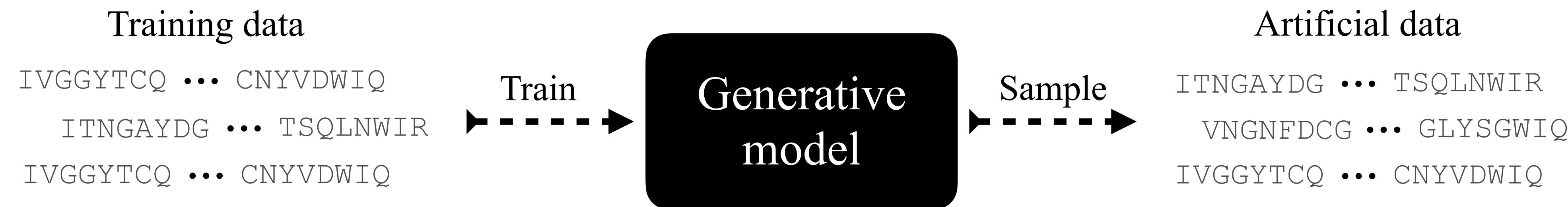
Train  
→ - - - →



# Generative models for protein sequences



# Generative models for protein sequences



# Generative models for protein sequences

## Principle

- ▶ Probability distribution over protein sequences

Training data

IVGGYTCQ ... CNYVDWIQ

ITNGAYDG ... TSQLNWIR

IVGGYTCQ ... CNYVDWIQ

Train  
→ → →

Generative  
model

↓  
Sample

Artificial data  $\sim P_{\text{model}}$

ITNGAYDG ... TSQLNWIR

VNGNFDG ... GLYSGWIQ

IVGGYTCQ ... CNYVDWIQ

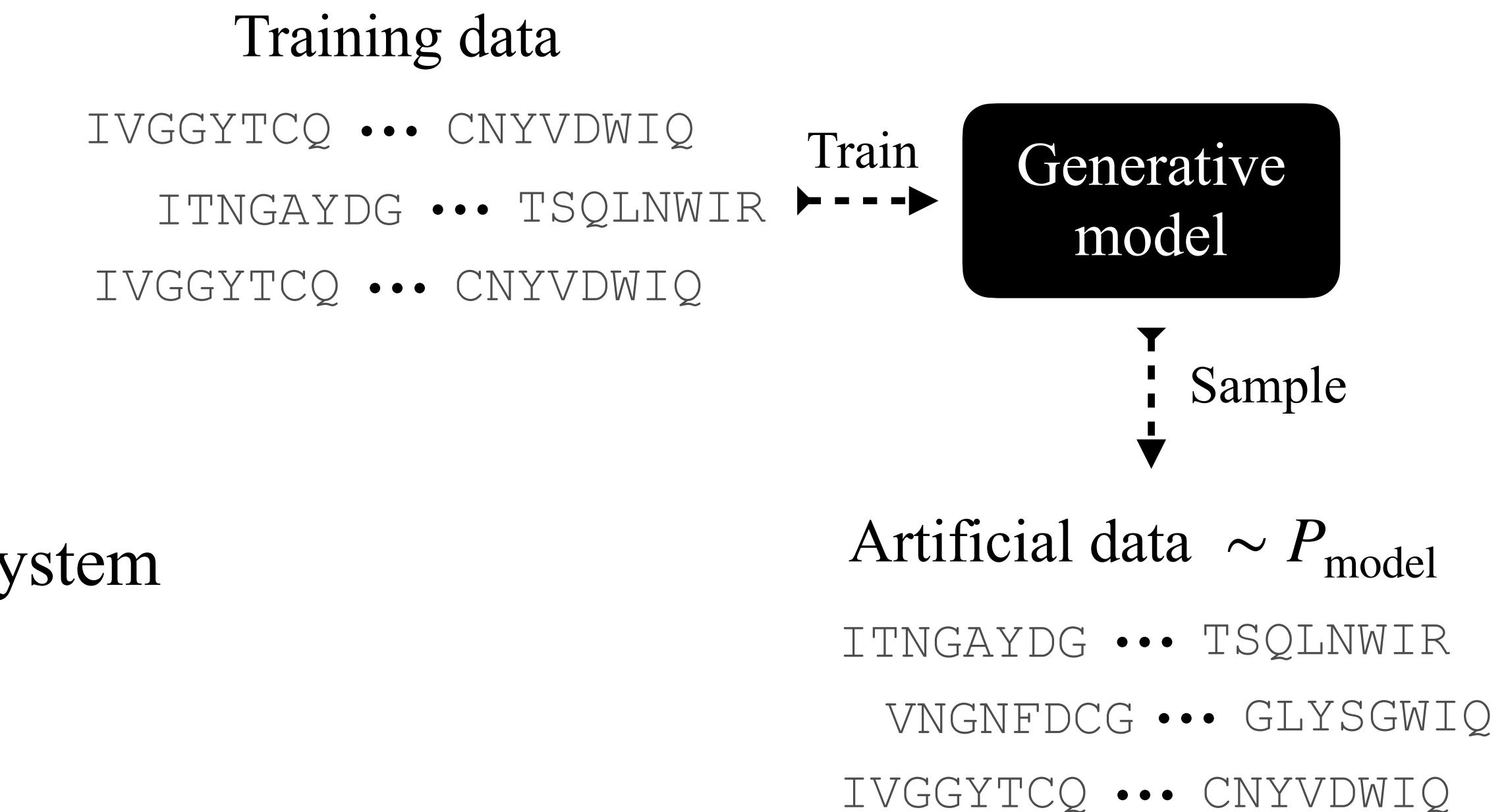
# Generative models for protein sequences

## Principle

- ▶ Probability distribution over protein sequences

## For what?

- ▶ Design-oriented goal: make new proteins
- ▶ Framework for understanding: parametrization of the system



# Generative models for protein sequences

## Principle

- ▶ Probability distribution over protein sequences

## For what?

- ▶ Design-oriented goal: make new proteins
- ▶ Framework for understanding: parametrization of the system

## Modeling a protein family

- ▶ High dimensional sequence space  $\sim 10^{66} - 10^{650}$
- ▶ Model evaluation: wet-lab characterization, in silico evaluations...

Training data

IVGGYTCQ ... CNYVDWIQ

ITNGAYDG ... TSQLNWIR

IVGGYTCQ ... CNYVDWIQ

Train  
→-----→

Generative  
model

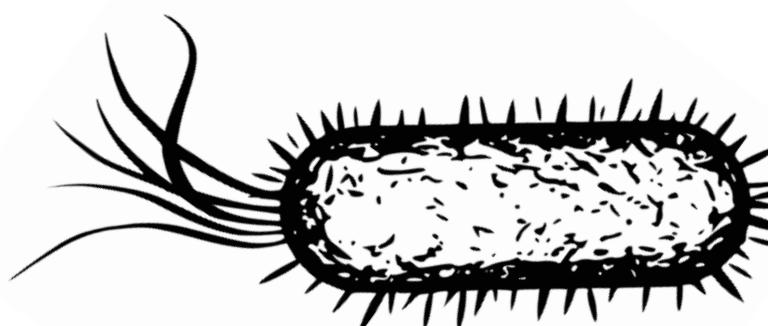
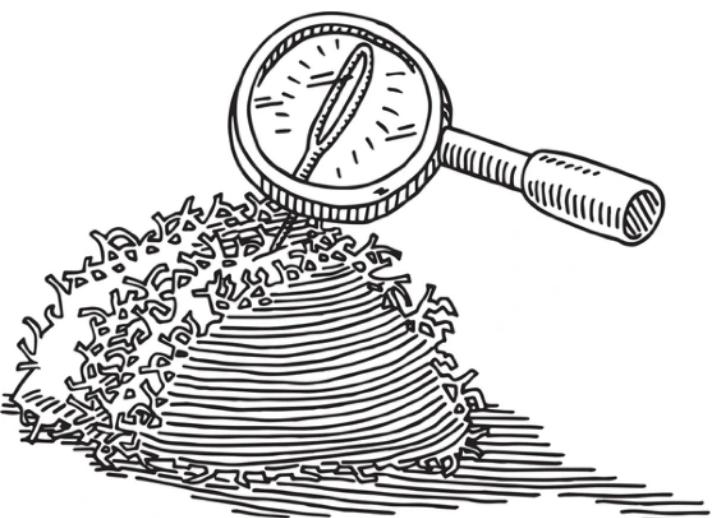
↓  
Sample

Artificial data  $\sim P_{\text{model}}$

ITNGAYDG ... TSQLNWIR

VNGNFDG ... GLYSGWIQ

IVGGYTCQ ... CNYVDWIQ



Hawkins-Hooker *et al.*, PLoS CB, 2021

Repecka *et al.*, Nature Machine Intelligence, 2021

Sgarbossa *et al.*, Elife, 2023

Watson L *et al.*, Nature, 2023

# Generative models for protein sequences

## Principle

- ▶ Probability distribution over protein sequences

## For what?

- ▶ Design-oriented goal: make new proteins
- ▶ Framework for understanding: parametrization of the system

## Modeling a protein family

- ▶ High dimensional sequence space  $\sim 10^{66} - 10^{650}$
- ▶ Model evaluation: wet-lab characterization, in silico evaluations...

## Model

- ▶ Variational Autoencoders, Diffusion models, Transformers, Restricted Boltzmann Machine, **Boltzmann Machine**...

## Training data

IVGGYTCQ ... CNYVDWIQ

ITNGAYDG ... TSQLNWIR

IVGGYTCQ ... CNYVDWIQ

Train  
→ → →

Generative  
model

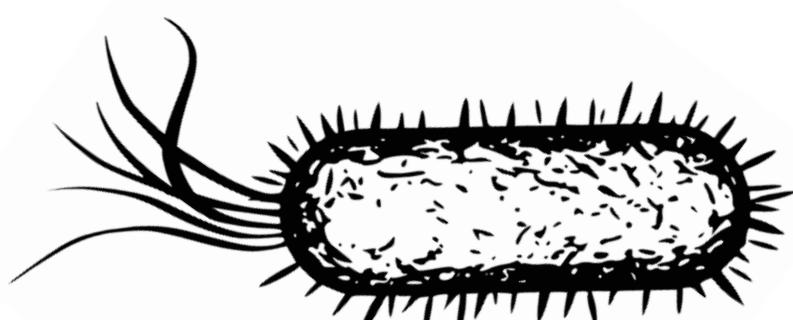
↓  
Sample

Artificial data  $\sim P_{\text{model}}$

ITNGAYDG ... TSQLNWIR

VNGNFDG ... GLYSGWIQ

IVGGYTCQ ... CNYVDWIQ



Hawkins-Hooker *et al.*, PLoS CB, 2021

Repecka *et al.*, Nature Machine Intelligence, 2021

Sgarbossa *et al.*, Elife, 2023

Watson L *et al.*, Nature, 2023

# Generative models for protein sequences

## Principle

- ▶ Probability distribution over protein sequences

## For what?

- ▶ Design-oriented goal: make new proteins
- ▶ Framework for understanding: parametrization of the system

## Modeling a protein family

- ▶ High dimensional sequence space  $\sim 10^{66} - 10^{650}$
- ▶ Model evaluation: wet-lab characterization, in silico evaluations...

## Model

- ▶ Variational Autoencoders, Diffusion models, Transformers, Restricted Boltzmann Machine, **Boltzmann Machine**...

## Boltzmann Machine

- ▶ Interpretability
- ▶ Mapping with other models
- ▶ Generative capacity experimentally tested (Russ *et al.* 2020)

## Training data

IVGGYTCQ ... CNYVDWIQ

ITNGAYDG ... TSQLNWIR

IVGGYTCQ ... CNYVDWIQ

Train  
→ → →

Generative  
model

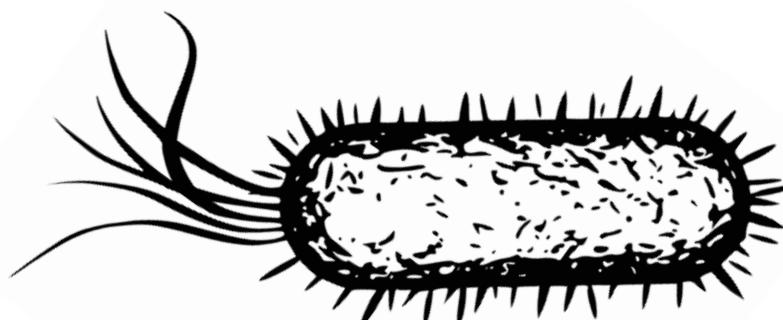
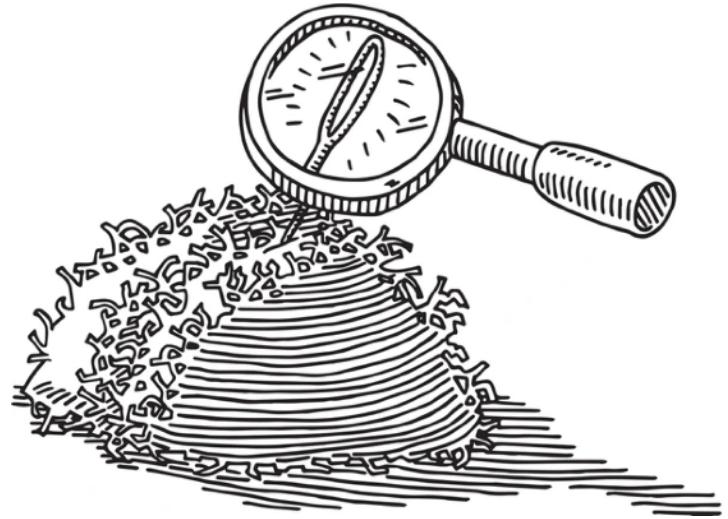
↓  
Sample

Artificial data  $\sim P_{\text{model}}$

ITNGAYDG ... TSQLNWIR

VNGNFDG ... GLYSGWIQ

IVGGYTCQ ... CNYVDWIQ



Hawkins-Hooker *et al.*, PLoS CB, 2021

Repecka *et al.*, Nature Machine Intelligence, 2021

Sgarbossa *et al.*, Elife, 2023

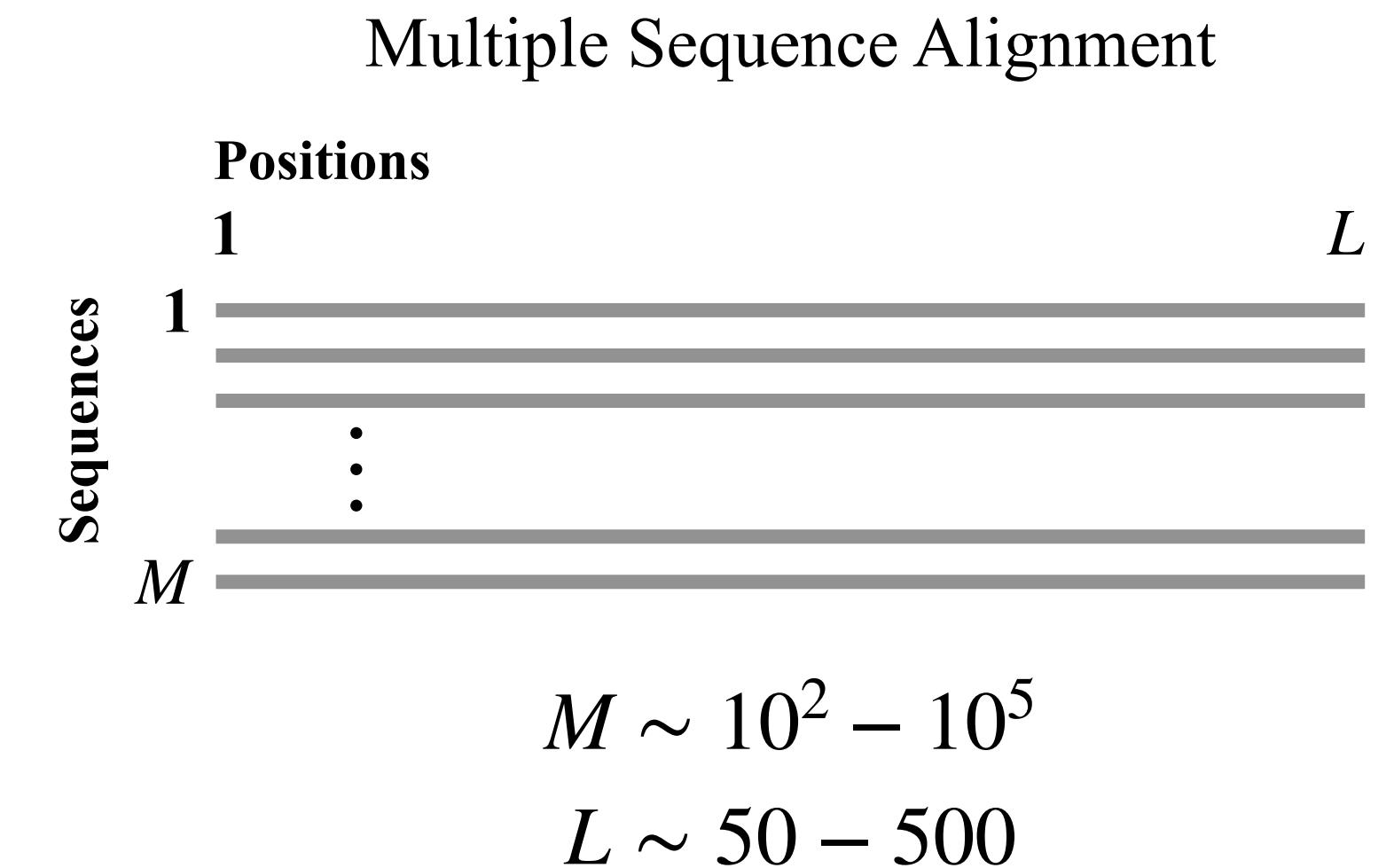
Watson L *et al.*, Nature, 2023

# Modeling a protein family with a Boltzmann Machine

# Modeling a protein family with a Boltzmann Machine (BM)

## Principle

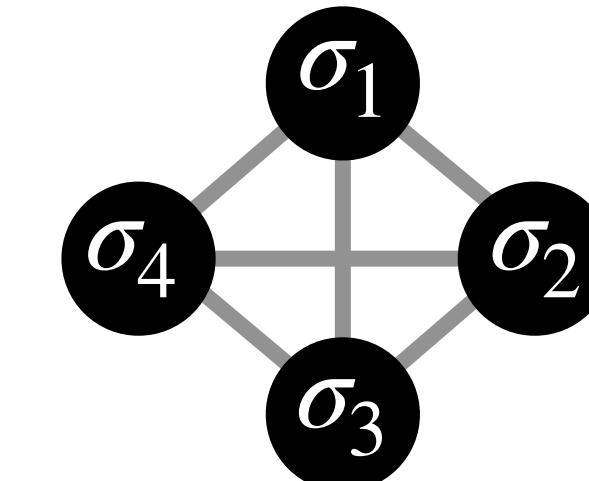
- Modeling of a protein family



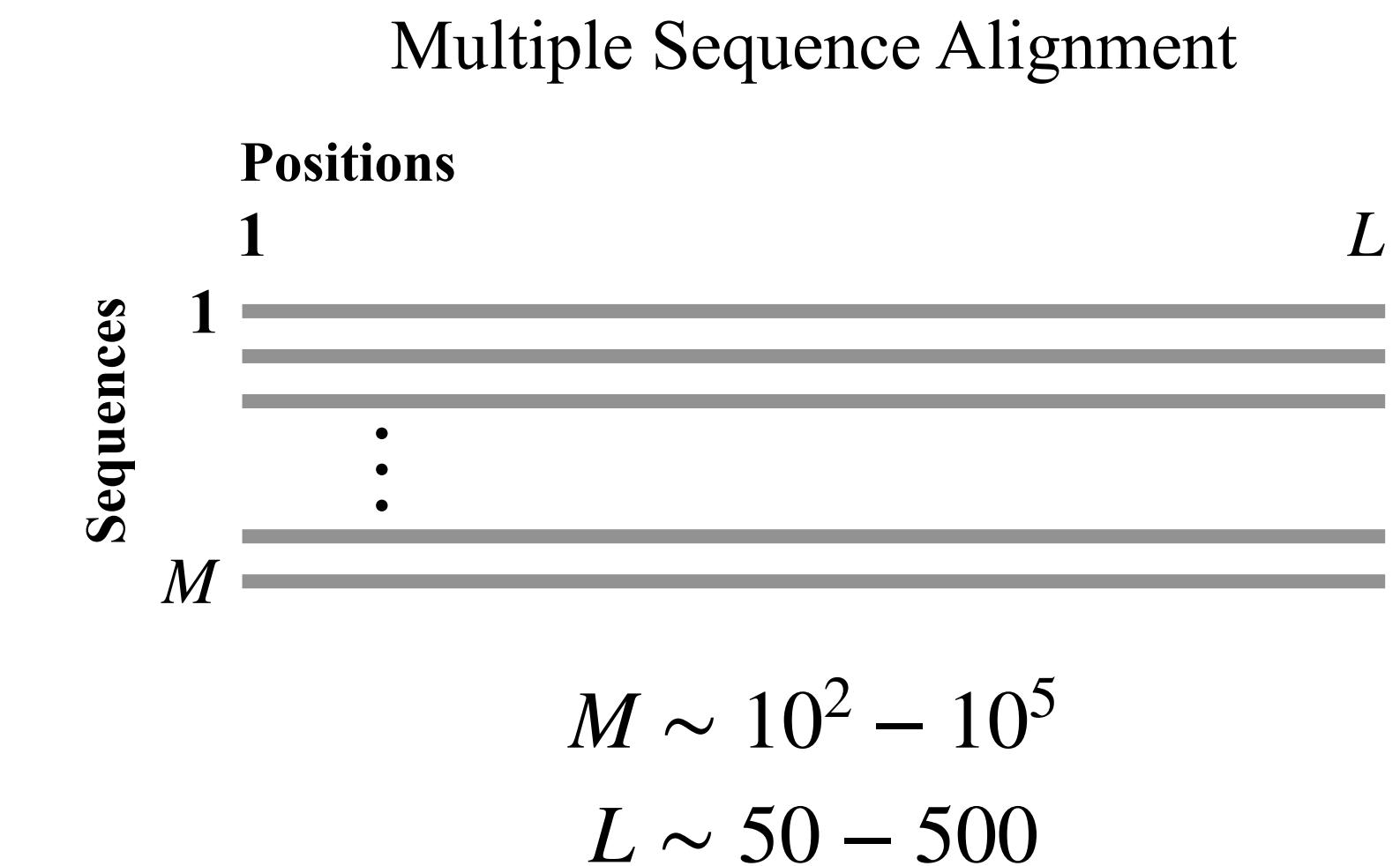
# Modeling a protein family with a Boltzmann Machine (BM)

## Principle

- ▶ Modeling of a protein family
- ▶ Graphical model: fully connected graph



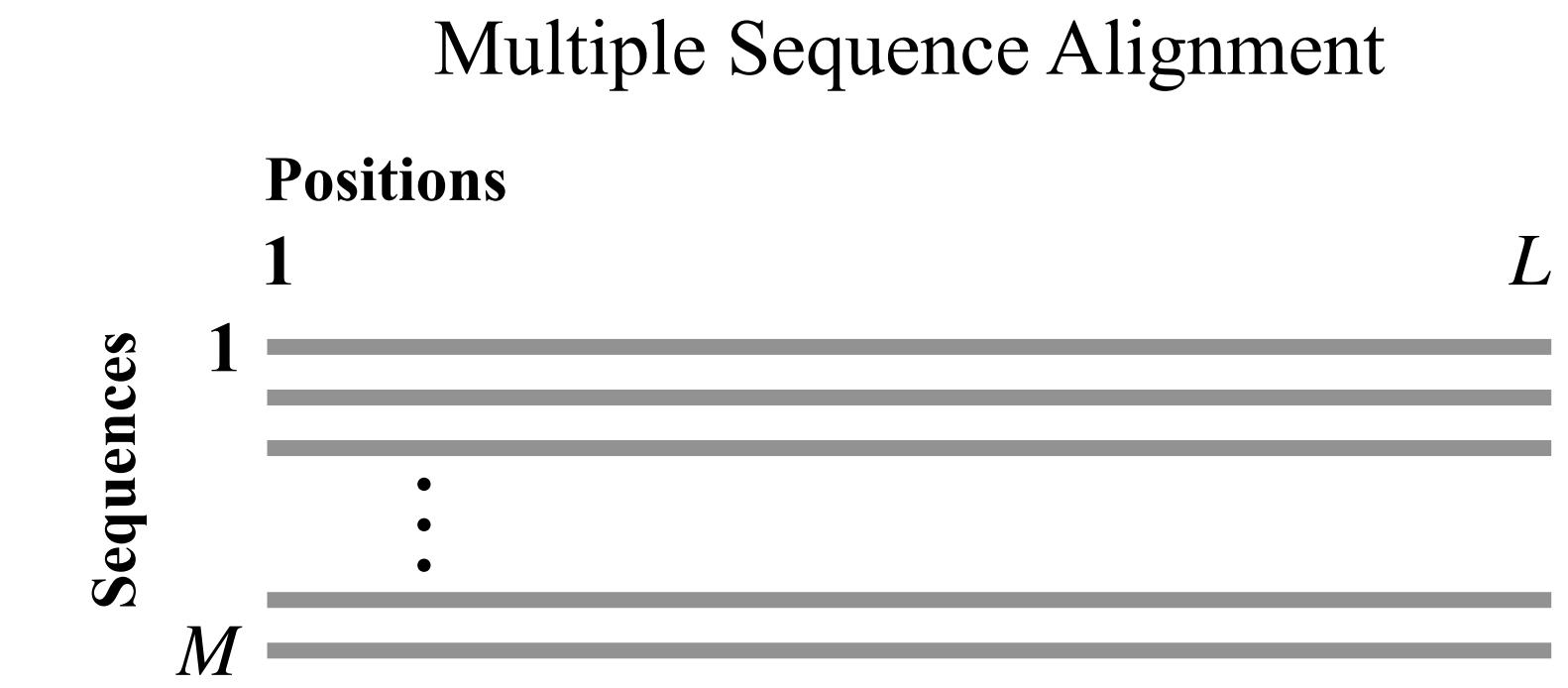
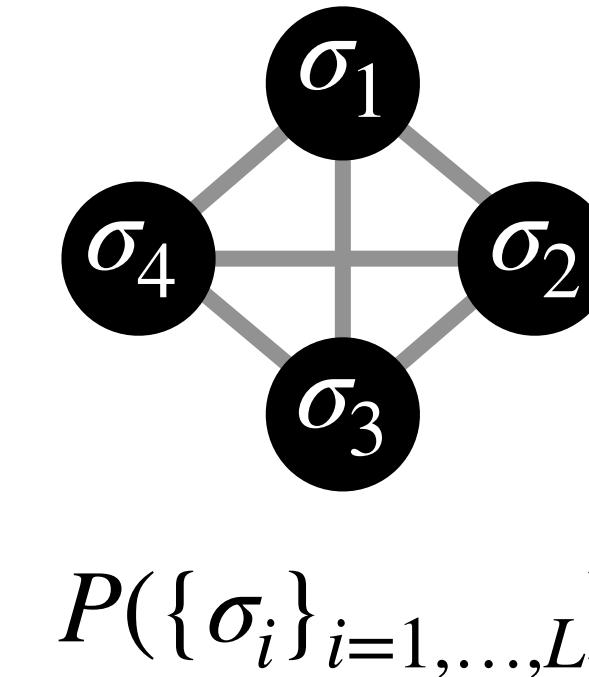
$$P(\{\sigma_i\}_{i=1,\dots,L})$$



# Modeling a protein family with a Boltzmann Machine (BM)

## Principle

- ▶ Modeling of a protein family
- ▶ Graphical model: fully connected graph
- ▶ Potts model

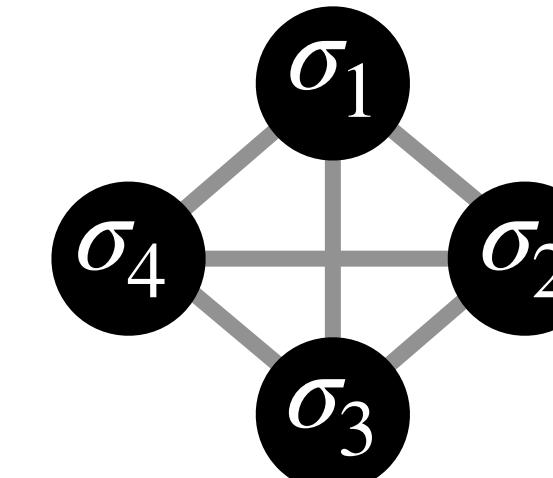


$$P(\{\sigma_i\}_{i=1,\dots,L}) = \frac{e^{\sum_{i=1}^L h_i(\sigma_i) + \sum_{i < j} J_{ij}(\sigma_i, \sigma_j)}}{Z}$$

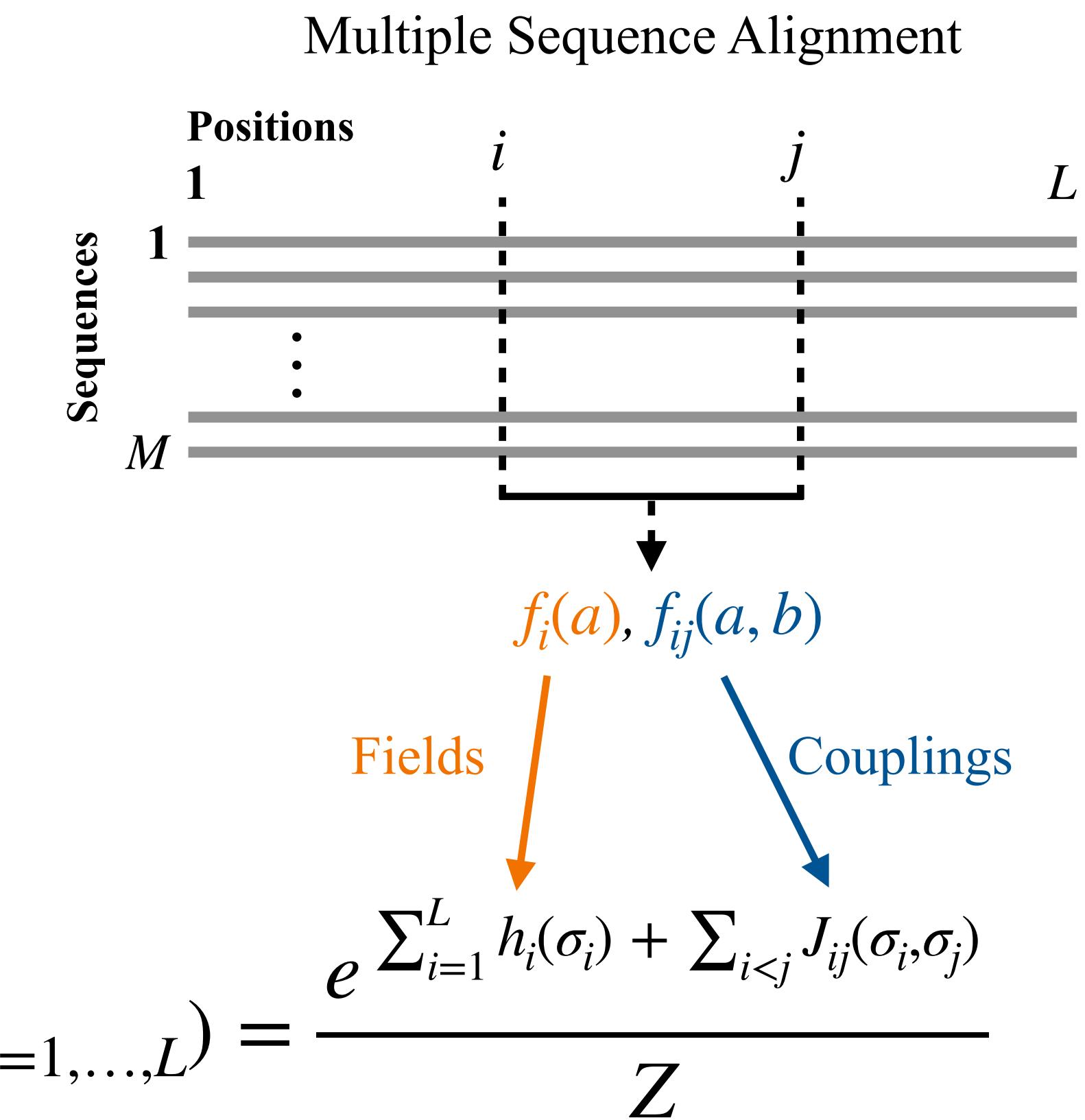
# Modeling a protein family with a Boltzmann Machine (BM)

## Principle

- ▶ Modeling of a protein family
- ▶ Graphical model: fully connected graph
- ▶ Potts model
- ▶ Trained to capture frequencies and pairwise frequencies (Maximum entropy approach)



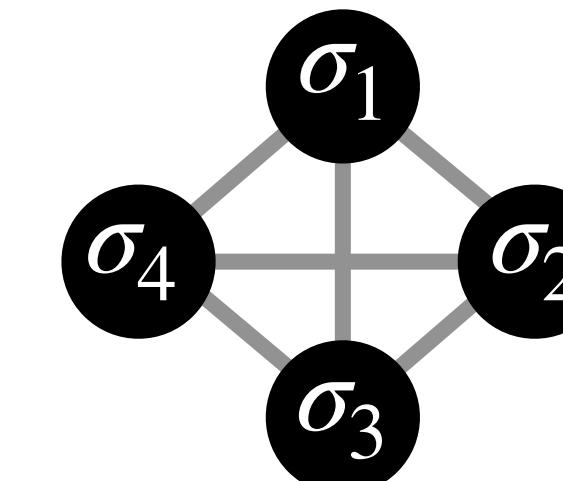
$$P(\{\sigma_i\}_{i=1,\dots,L})$$



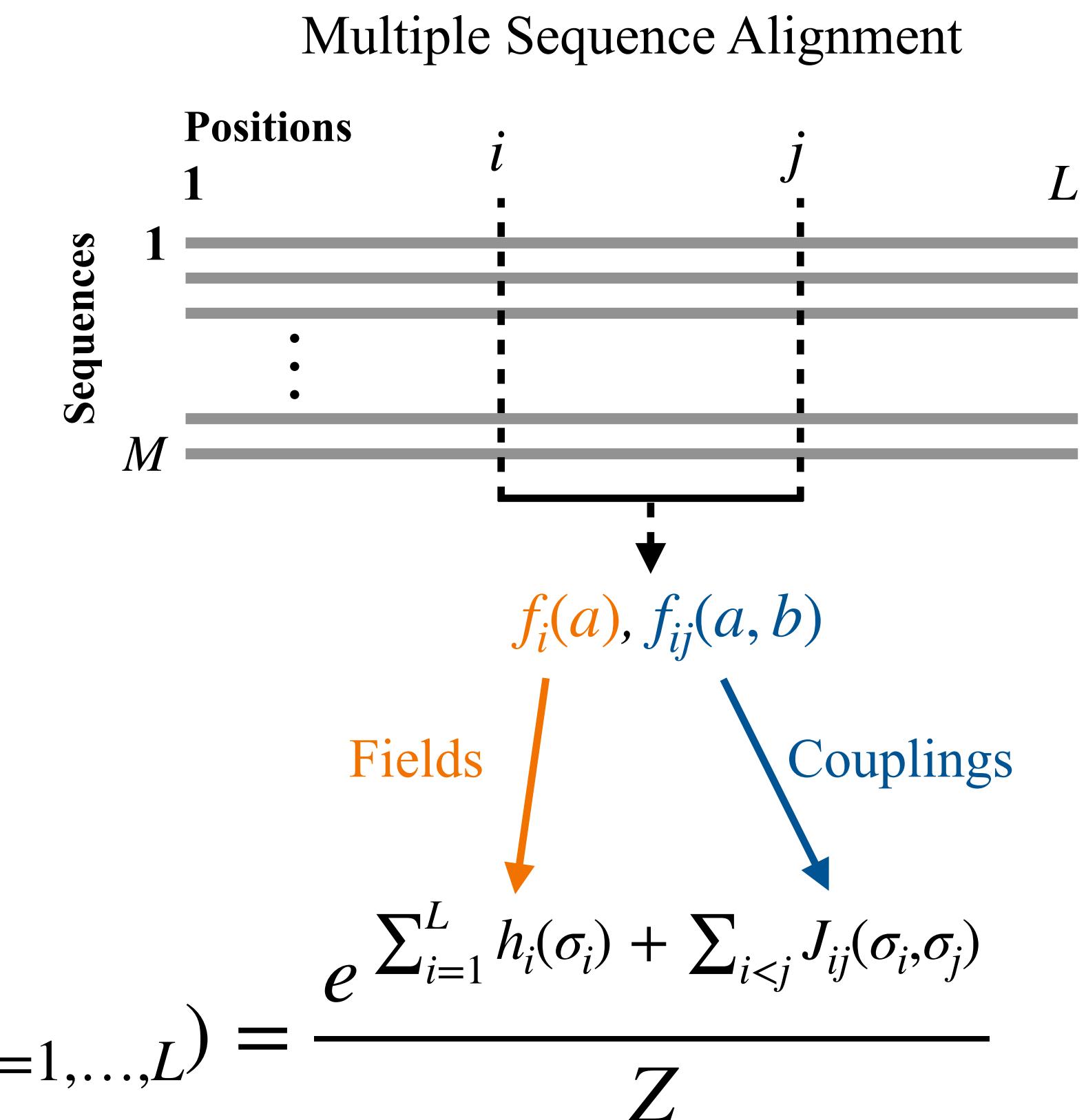
# Modeling a protein family with a Boltzmann Machine (BM)

## Principle

- ▶ Modeling of a protein family
- ▶ Graphical model: fully connected graph
- ▶ Potts model
- ▶ Trained to capture frequencies and pairwise frequencies (Maximum entropy approach)



$$P(\{\sigma_i\}_{i=1,\dots,L})$$



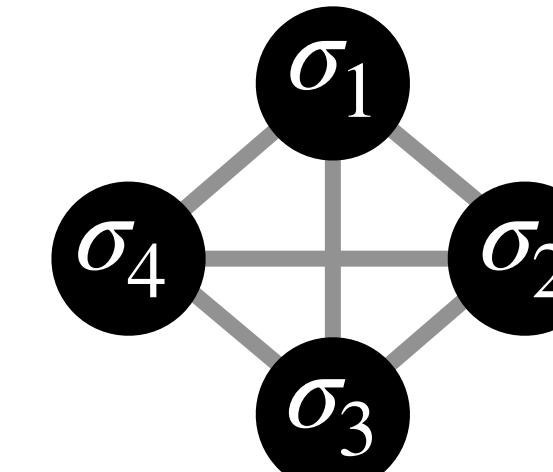
## Inference

- ▶ Parameters that maximize the probability of natural sequences (MLE)
- ▶ Other methods: Mean field, Pseudo-likelihood maximization, Autoregressive model...

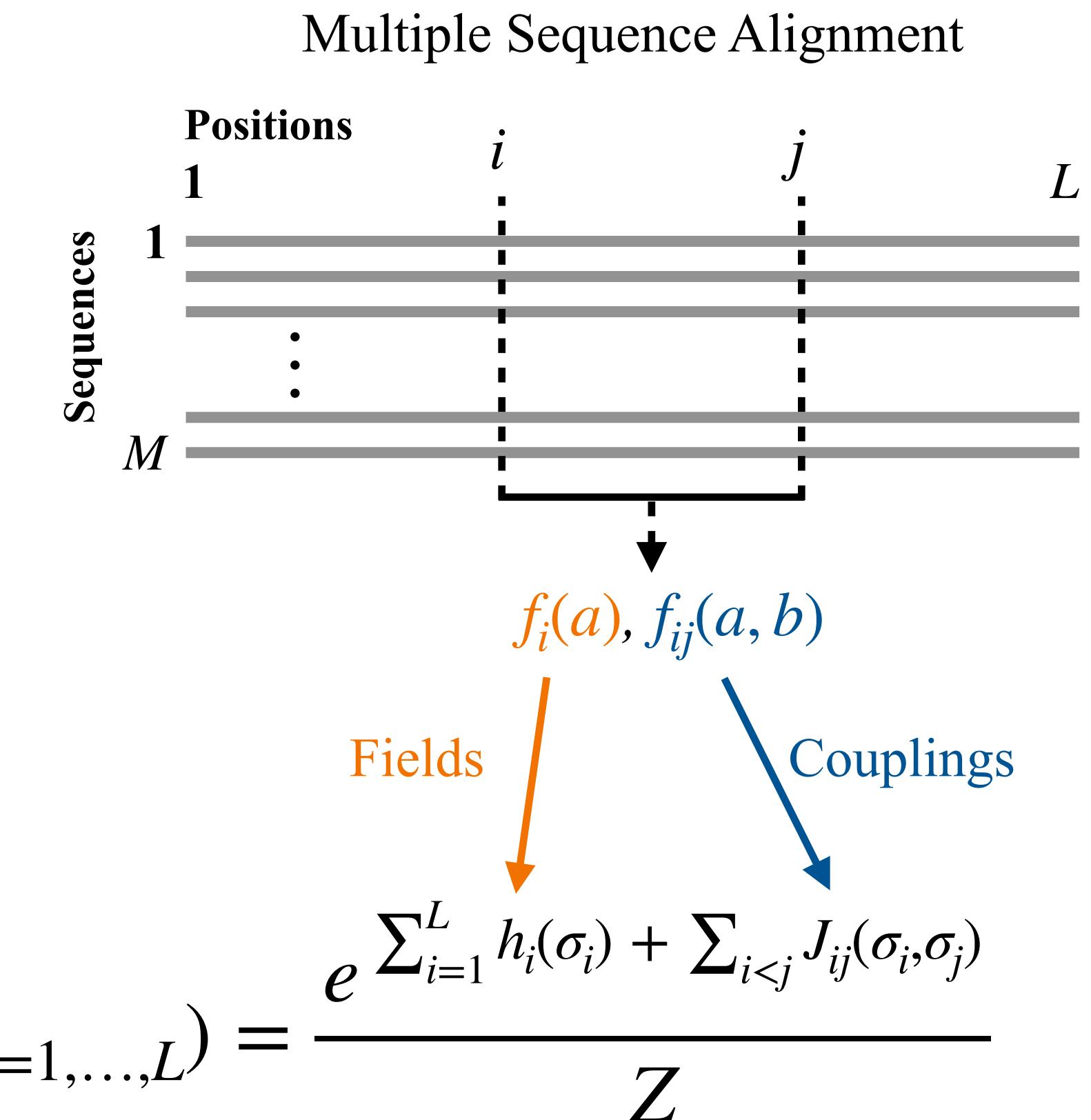
# Modeling a protein family with a Boltzmann Machine (BM)

## Principle

- ▶ Modeling of a protein family
- ▶ Graphical model: fully connected graph
- ▶ Potts model
- ▶ Trained to capture frequencies and pairwise frequencies (Maximum entropy approach)



$$P(\{\sigma_i\}_{i=1,\dots,L})$$



## Inference

- ▶ Parameters that maximize the probability of natural sequences (MLE)
- ▶ Other methods: Mean field, Pseudo-likelihood maximization, Autoregressive model...

## Results

- ▶ Predict structural contacts
- ▶ Predict mutational effect
- ▶ **Generative with low-temperature sampling** (Russ et al., 2020)

$$P(\{\sigma_i\}_{i=1,\dots,L}) \sim e^{-\frac{E(h, J)}{T}} \text{ with } T < 1$$

# The undersampling problem

## Problem

- ▶ # parameters  $\sim 10^5 - 10^7 \gg$  # sequences  $\sim 10^2 - 10^5$
- ▶ Extreme statistics from undersampling  $\rightarrow$  infinite parameters

## Multiple Sequence Alignment



# The undersampling problem

## Problem

- ▶ # parameters  $\sim 10^5 - 10^7 \gg$  # sequences  $\sim 10^2 - 10^5$
- ▶ Extreme statistics from undersampling  $\rightarrow$  infinite parameters

## Regularization methods

- ▶ Remove parameters (Pruning, Alphabet reduction...)
- ▶ Modify statistics (pseudo-counts)
- ▶ Constrain parameters during optimization ( $L_p$  norm...)

## Multiple Sequence Alignment



# The undersampling problem

## Problem

- ▶ # parameters  $\sim 10^5 - 10^7 \gg$  # sequences  $\sim 10^2 - 10^5$
- ▶ Extreme statistics from undersampling  $\rightarrow$  infinite parameters

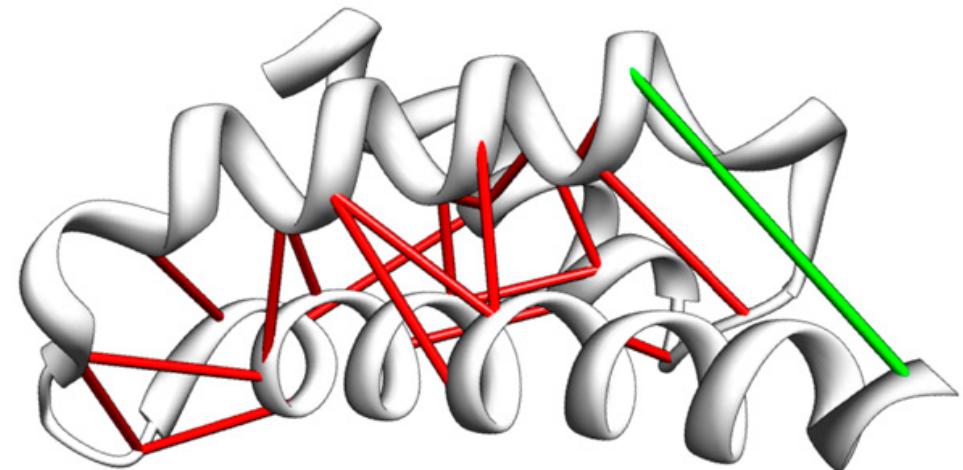
## Regularization methods

- ▶ Remove parameters (Pruning, Alphabet reduction...)
- ▶ Modify statistics (pseudo-counts)
- ▶ Constrain parameters during optimization ( $L_p$  norm...)

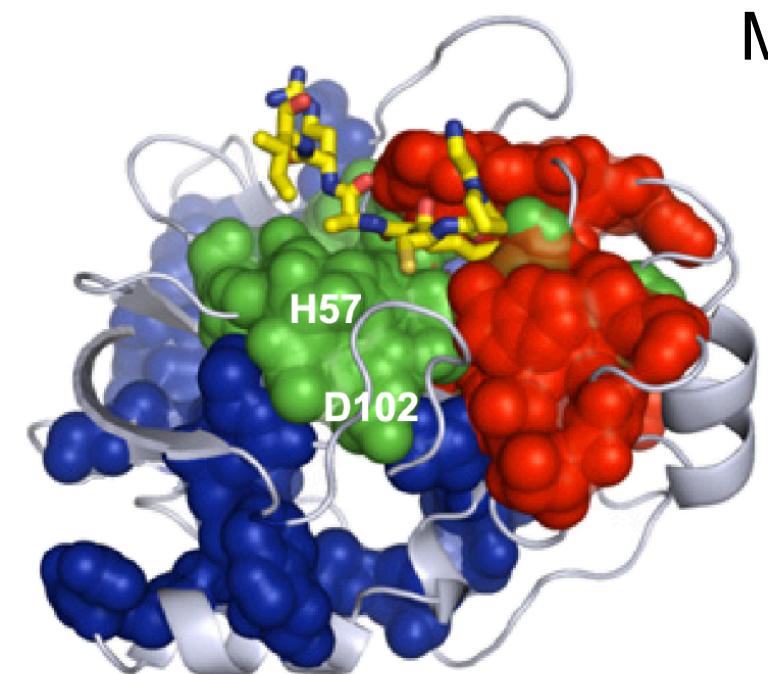
## Importance of data's statistical structure

- ▶ Real data often have rich statistical structure
- ▶ Proteins: Correlated units of different sizes, magnitude...
- ▶ Uneven impact of undersampling on different statistical signatures (Kleerorin *et al.* 2023)

## Multiple Sequence Alignment



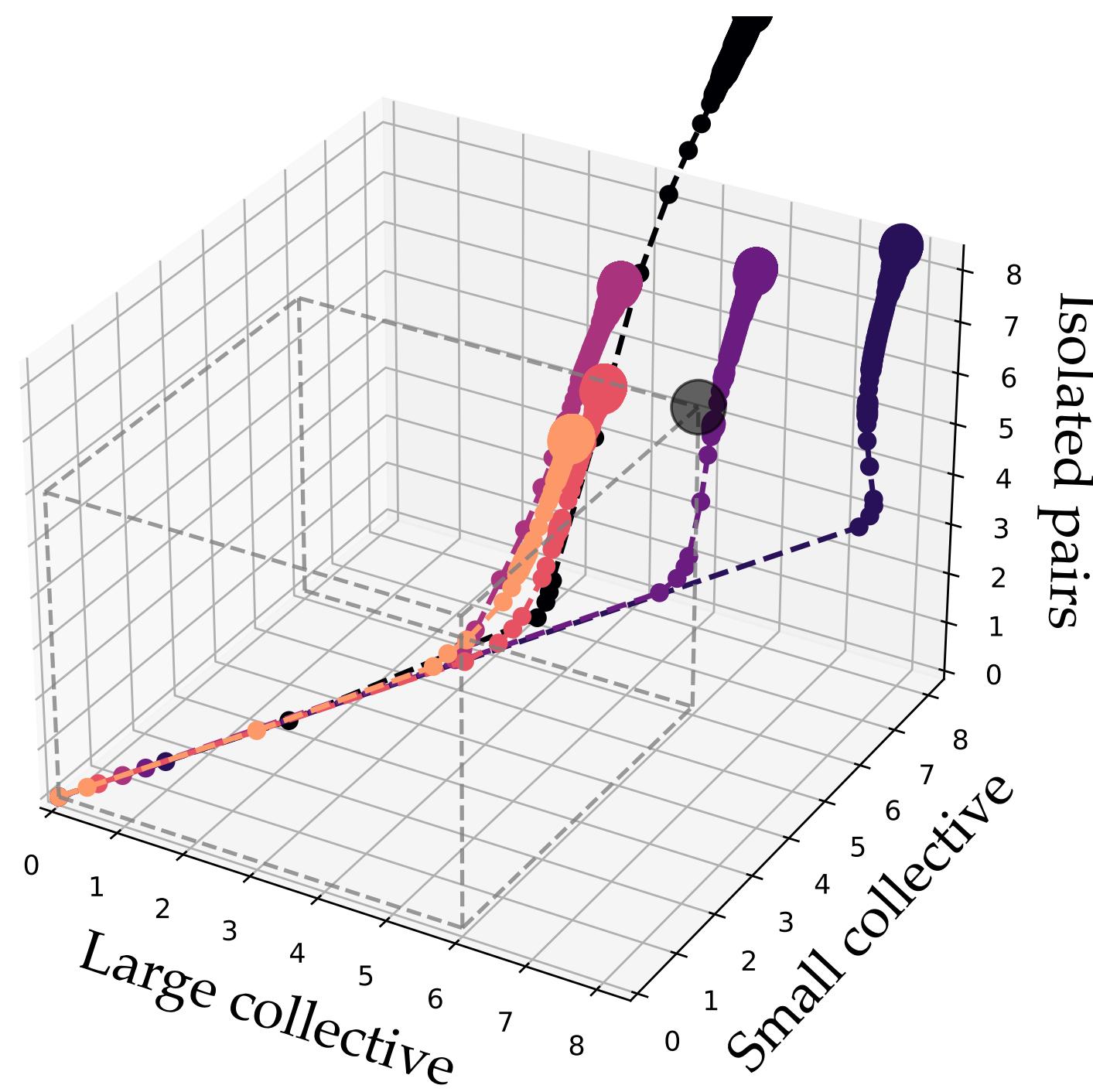
Morcos *et al.*, PNAS, 2011



Halabi *et al.*, Cell, 2009

## The undersampling problem

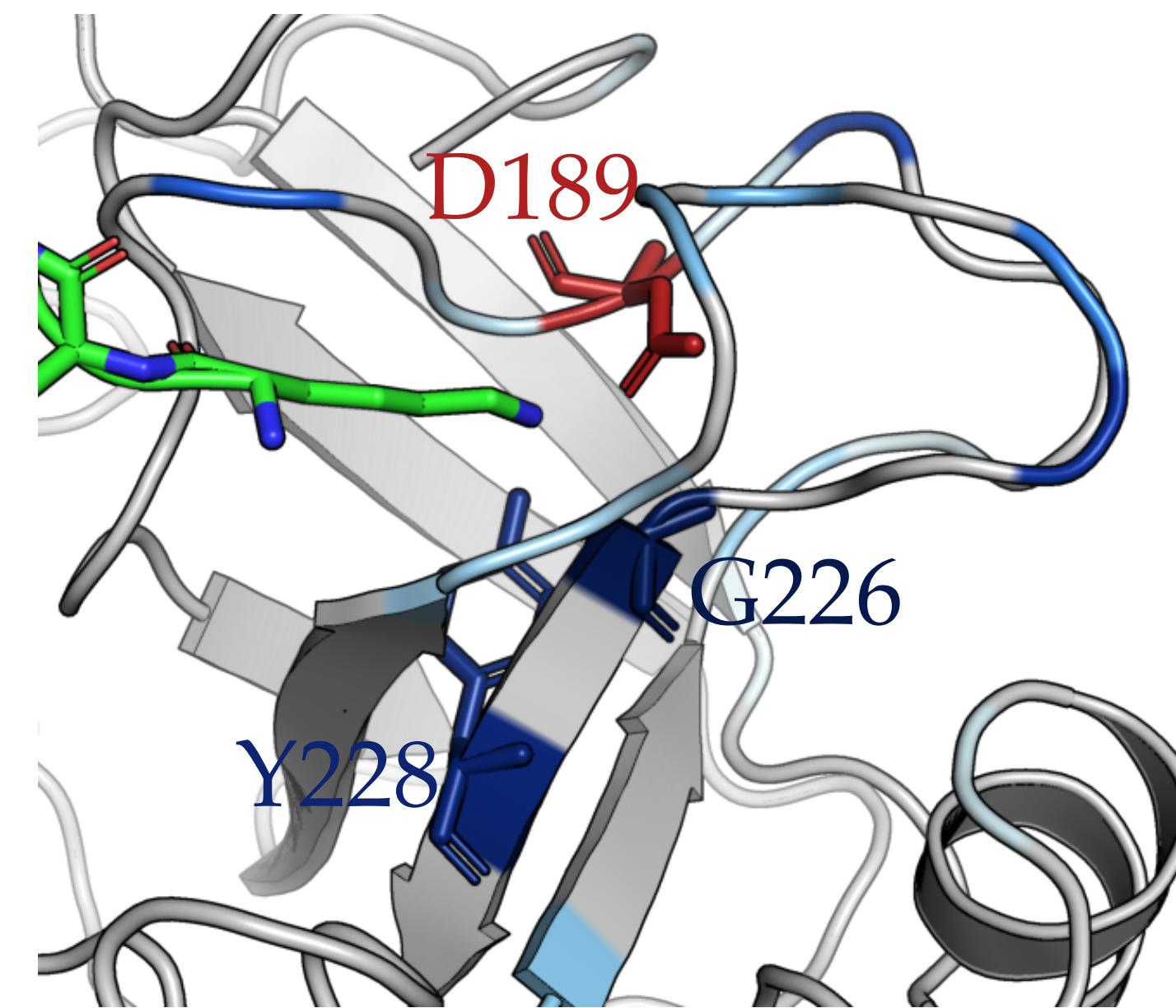
Overcoming undersampling induced biases



*In collaboration with Emily Hinds, Yaakov Kleeorin, Rama Ranganathan (University of Chicago, USA)*

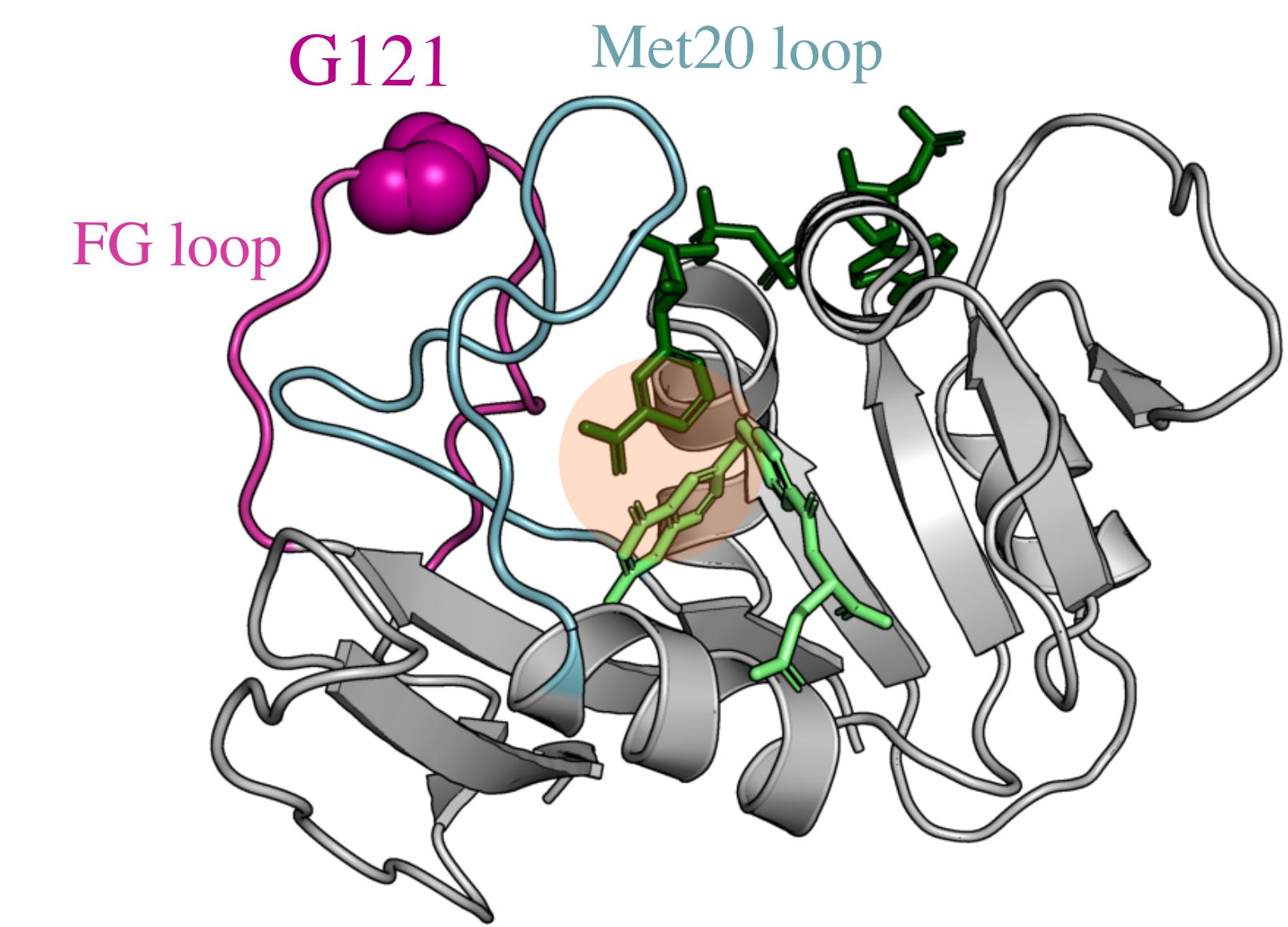
## Investigate protein properties with statistical learning

Specificity mechanism in S1A family



*In collaboration with Amaury Paveyranne, Timothé Lucas, Shoichi Yip, Clément Nizak (LJP, Sorbonne University, France)*

Allosteric network of *E. Coli* DHFR



*In collaboration with Paul Guenon, Damien Laage, Guillaume Stirnemann (ENS, France), Clément Nizak (LJP, France), Karolina Filipowska, Kim Reynolds (University of Texas, USA)*

# The Undersampling problem

*In collaboration with Emily Hinds, Yaakov Kleeeorin,  
Rama Ranganathan (University of Chicago, USA)*

# The undersampling problem

**I. Undersampling-induced biases (Kleedorin *et al.*, Cell system, 2023)**

**II. Generative Capacity of the Boltzmann Machine (Russ *et al.*, Science, 2020)**

**III. New inference method: Stochastic Boltzmann Machine**

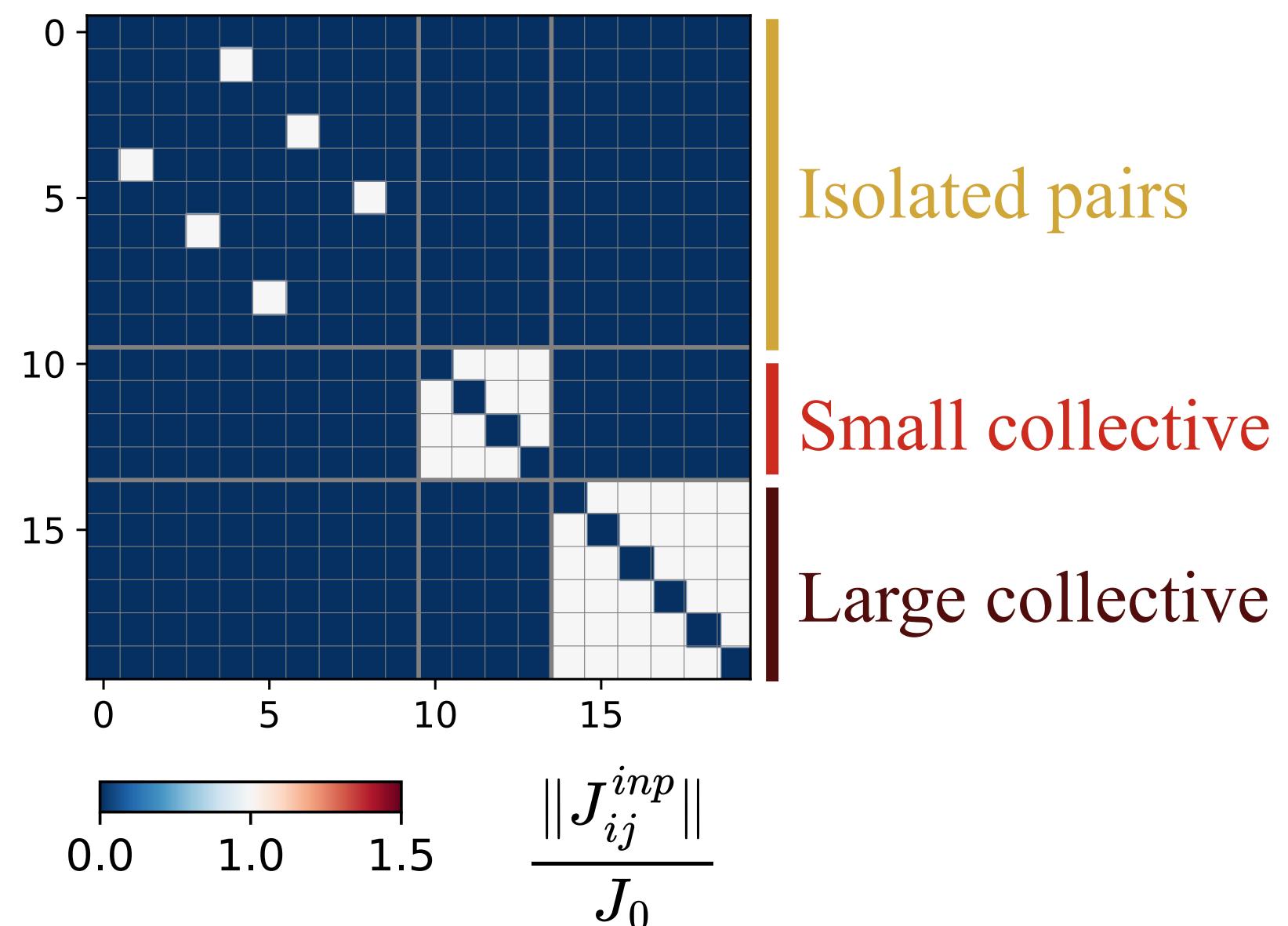
# I. Undersampling-induced biases (Kleeorin *et al.* 2023)

*Assess model performance with a toy model*

## Toy model features

- ▶ Boltzmann Machine model
- ▶ Correlated units of different sizes

$$P(\{\sigma_i\}_{i=1,\dots,L}) = \frac{e^{\sum_{i=1}^L h_i(\sigma_i) + \sum_{i < j} J_{ij}(\sigma_i, \sigma_j)}}{Z}$$



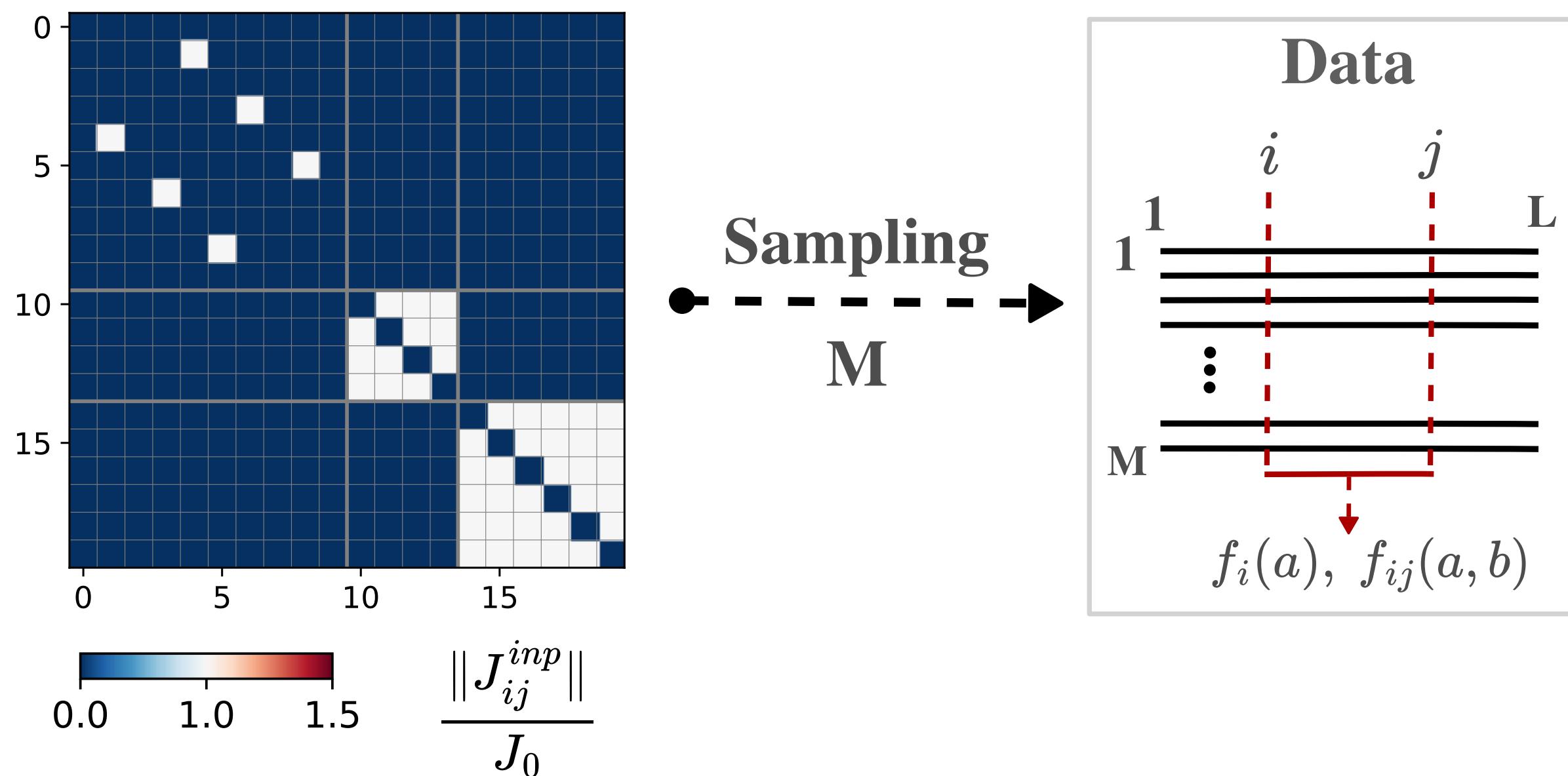
# I. Undersampling-induced biases (Kleeorin *et al.* 2023)

*Assess model performance with a toy model*

## Toy model features

- ▶ Boltzmann Machine model
- ▶ Correlated units of different sizes

$$P(\{\sigma_i\}_{i=1,\dots,L}) = \frac{e^{\sum_{i=1}^L h_i(\sigma_i) + \sum_{i < j} J_{ij}(\sigma_i, \sigma_j)}}{Z}$$



# I. Undersampling-induced biases (Kleeorin *et al.* 2023)

*Assess model performance with a toy model*

## Toy model features

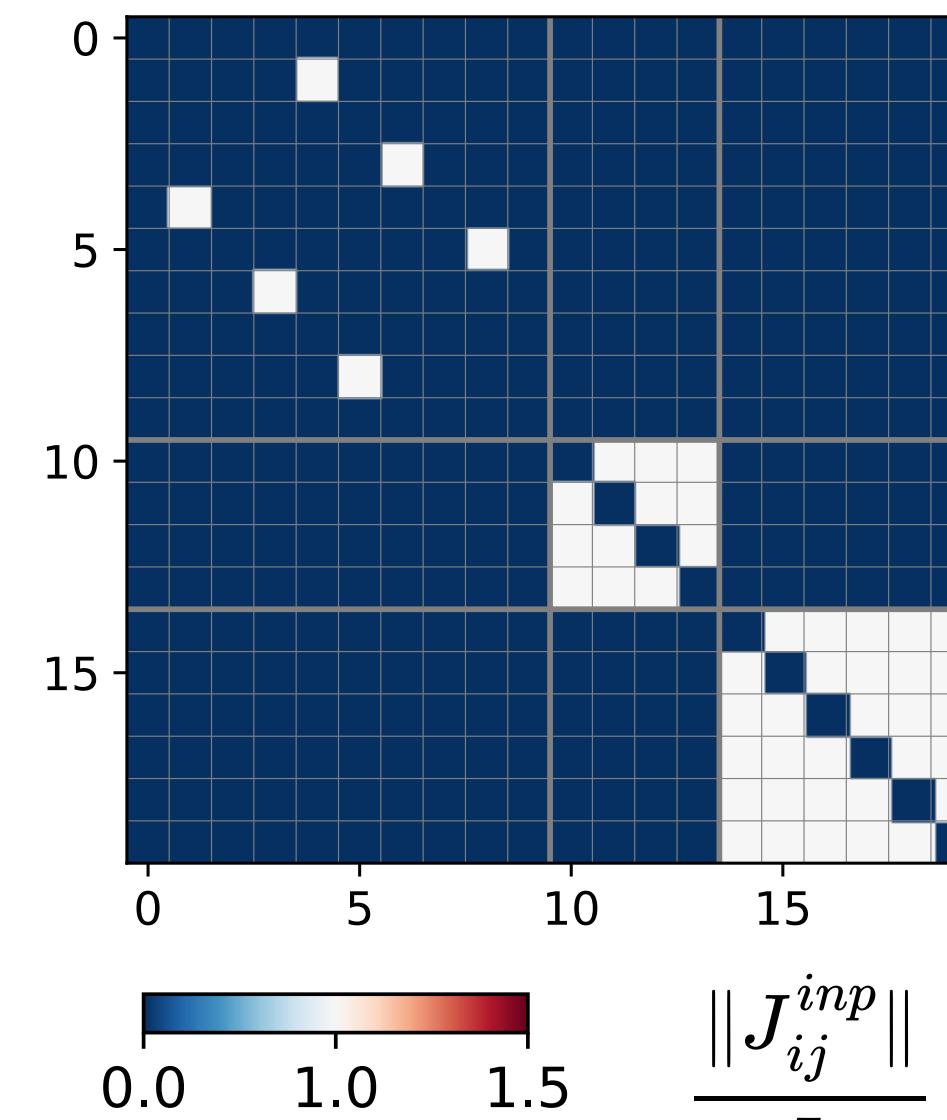
- Boltzmann Machine model
- Correlated units of different sizes

$$P(\{\sigma_i\}_{i=1,\dots,L}) = \frac{e^{\sum_{i=1}^L h_i(\sigma_i) + \sum_{i < j} J_{ij}(\sigma_i, \sigma_j)}}{Z}$$

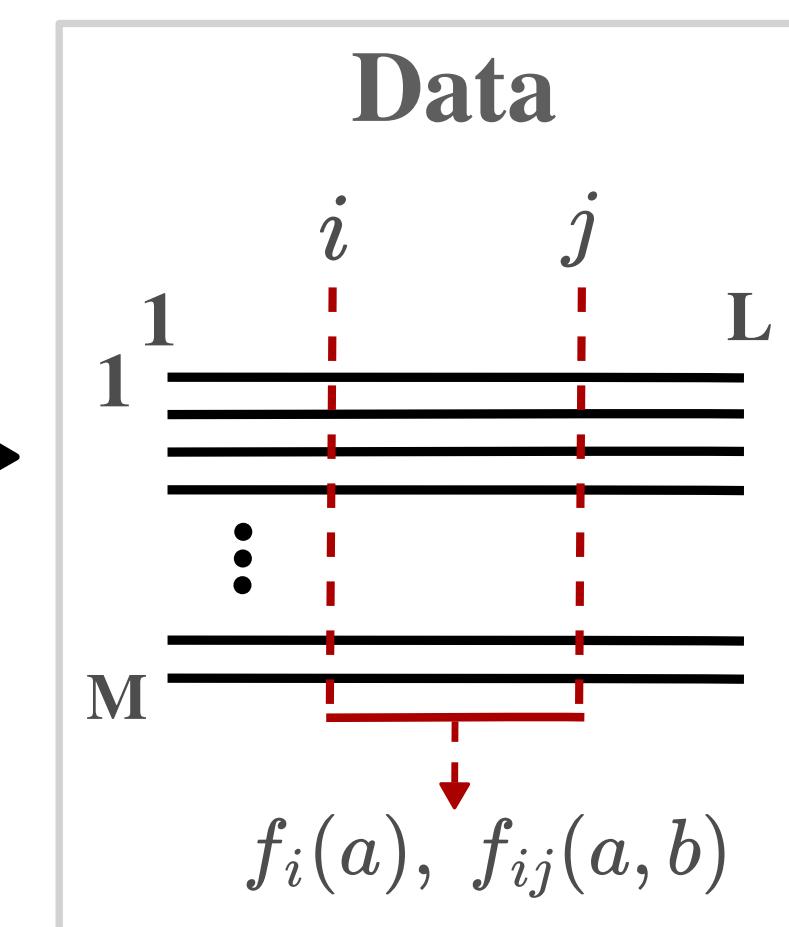
## Inference:

- Undersampling regime
- Log-likelihood maximization with  $L_2$  regularization

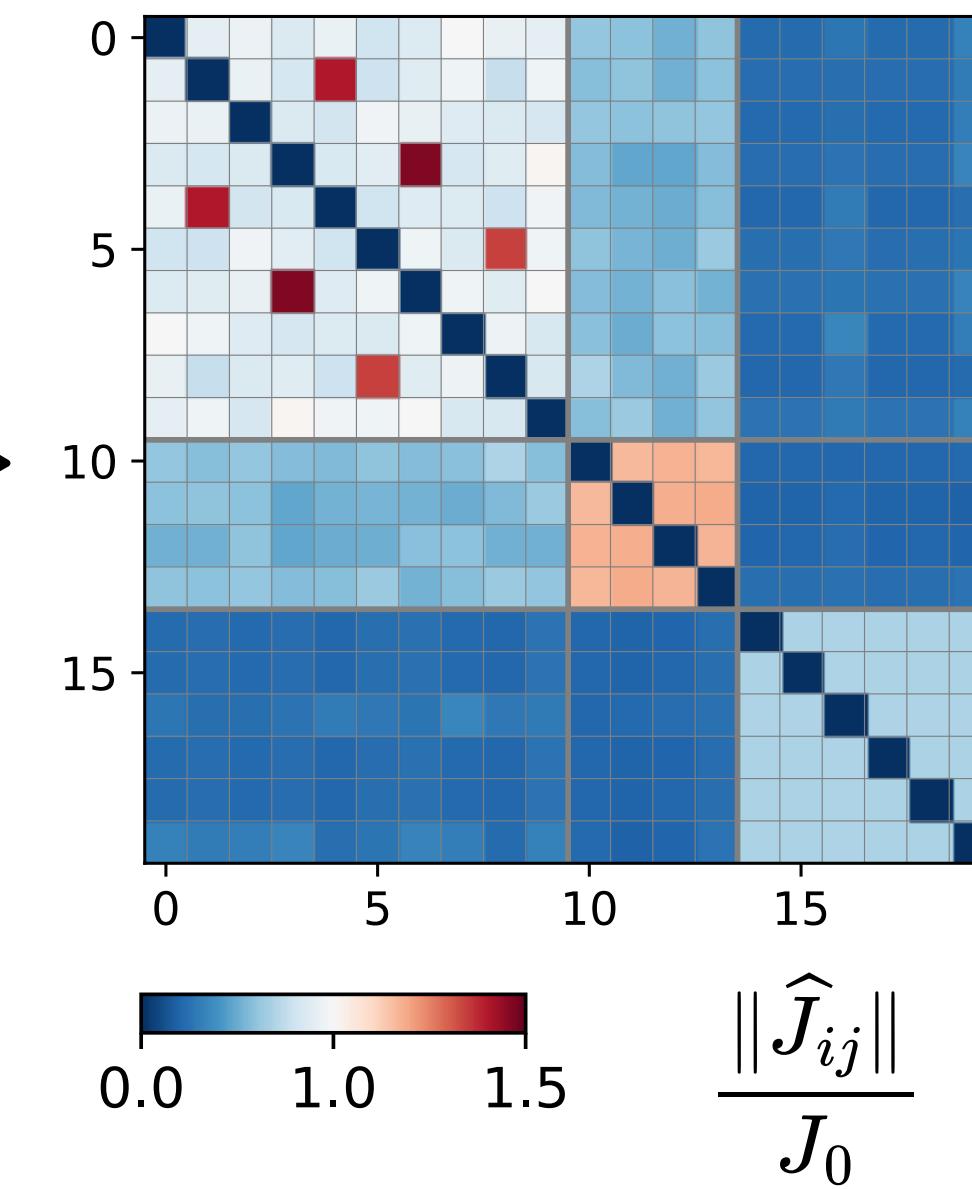
$$\theta^* = \arg \max_{\theta} \left[ \frac{1}{M} \sum_{m=1}^M \log P(\sigma^{(m)} | \theta) - \lambda_J \|J\|^2 - \lambda_h \|h\|^2 \right]$$



Sampling  
M



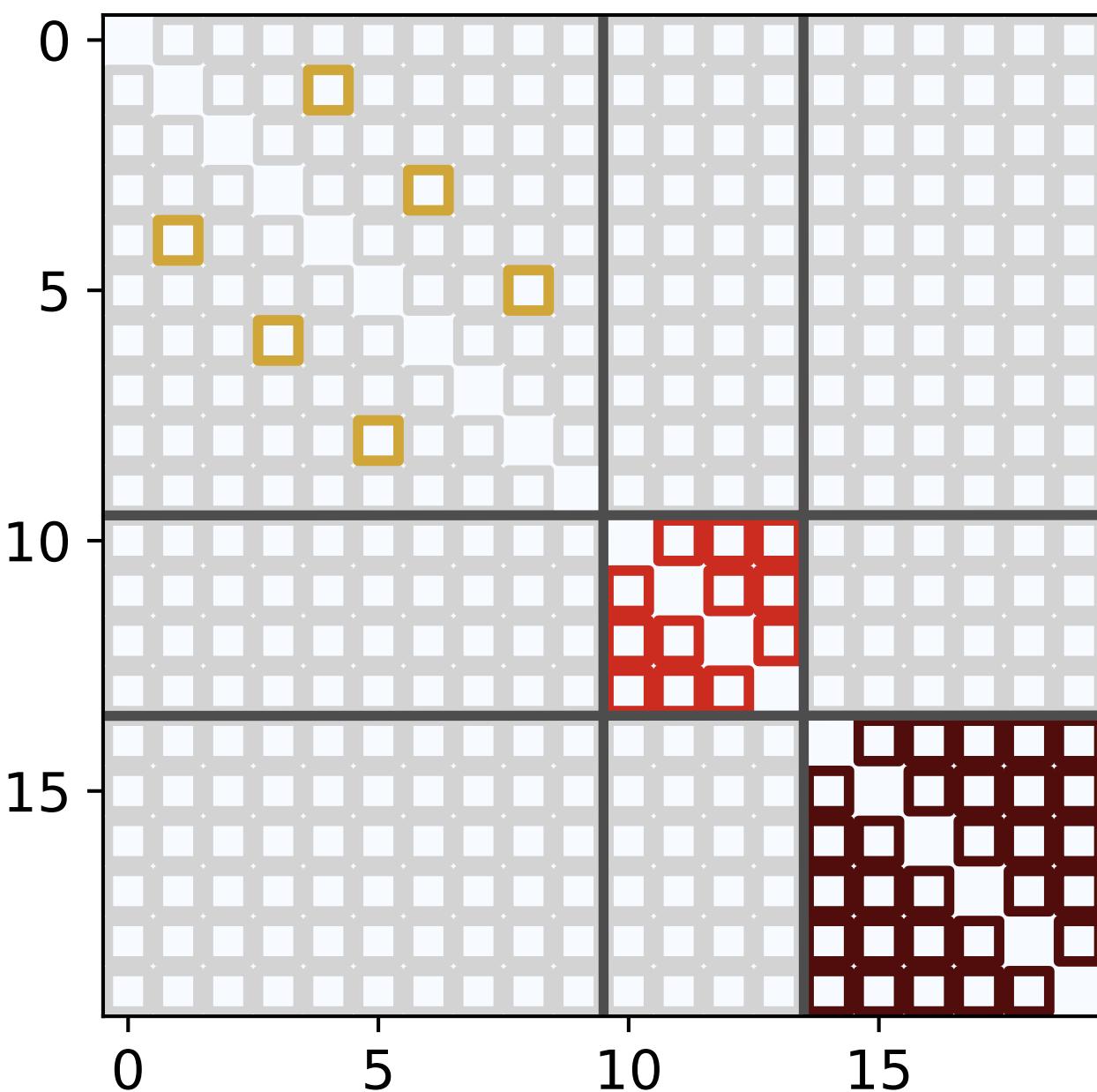
Inference



# I. Undersampling-induced biases (Kleeorin *et al.* 2023)

*Inference as function of regularization strength*

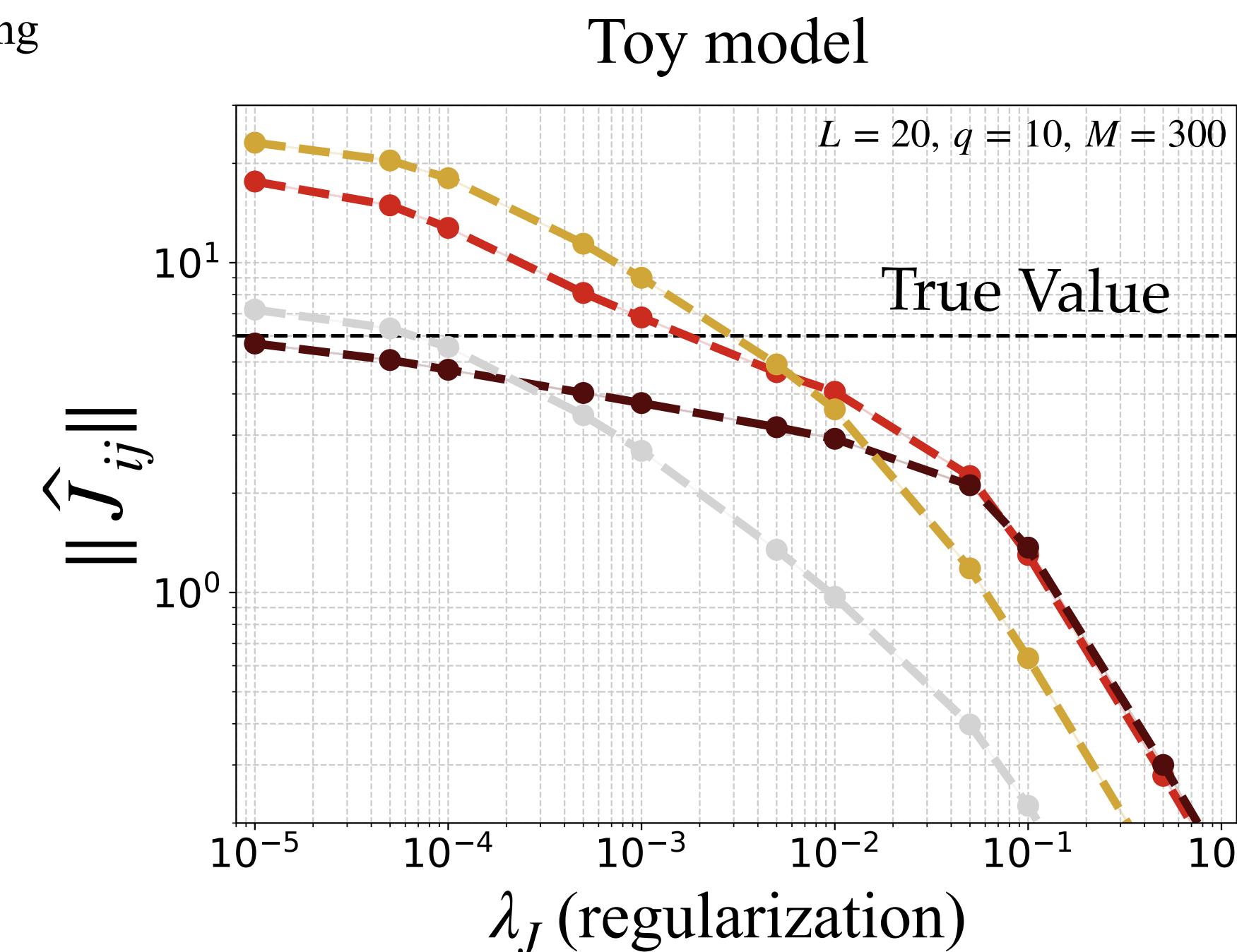
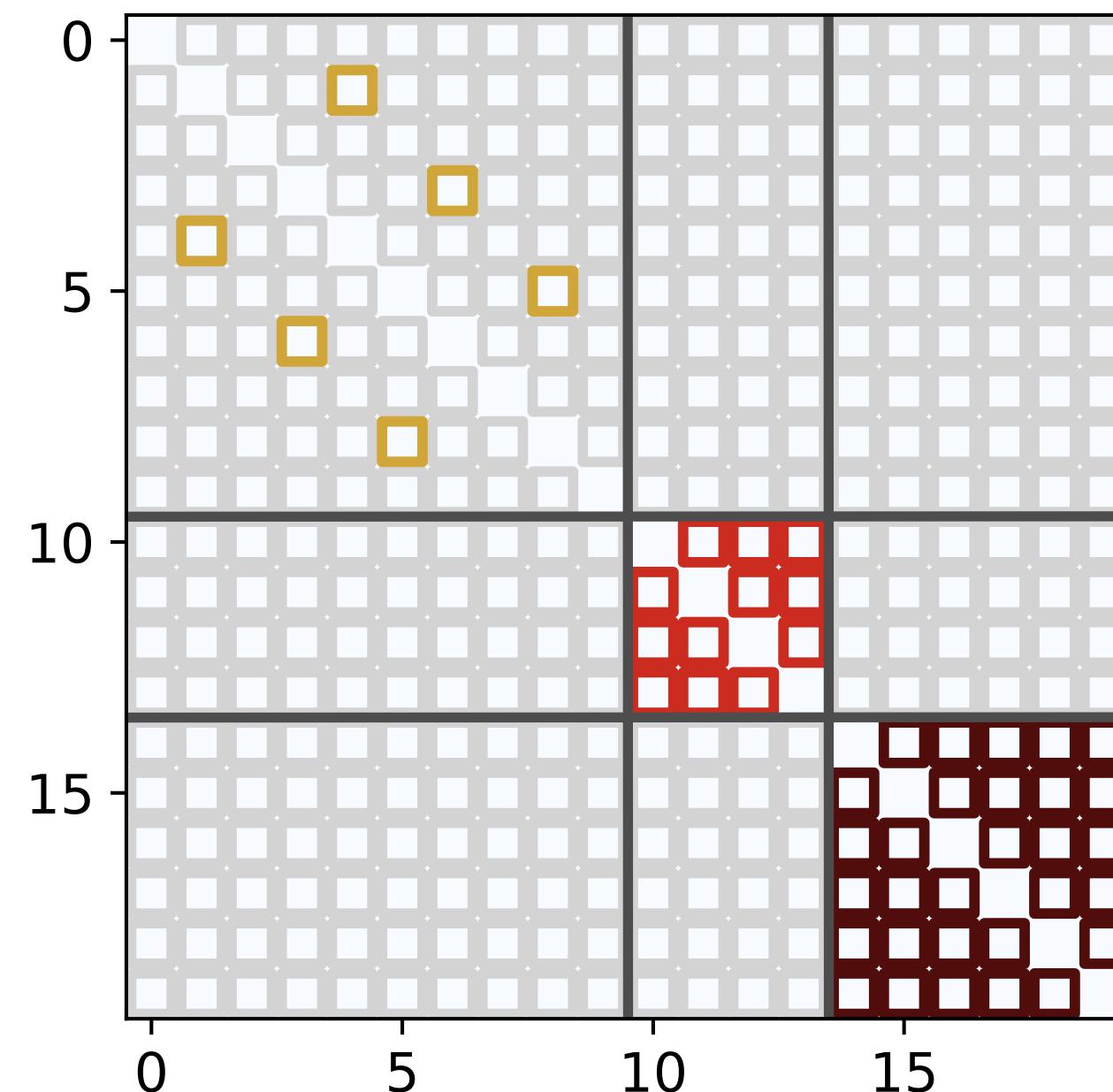
Isolated pairs      Large collective  
Small collective      Non interacting



# I. Undersampling-induced biases (Kleeorin *et al.* 2023)

*Inference as function of regularization strength*

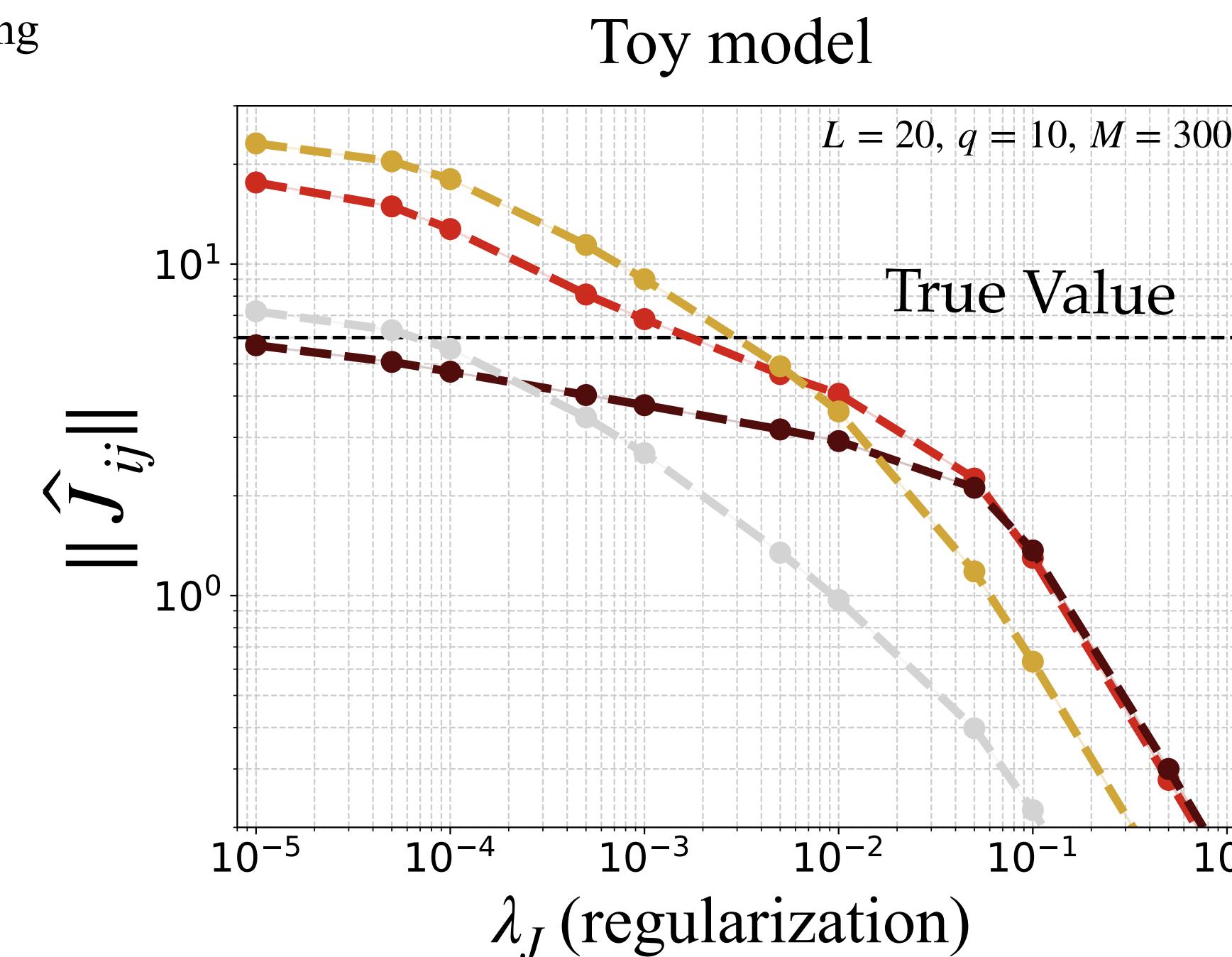
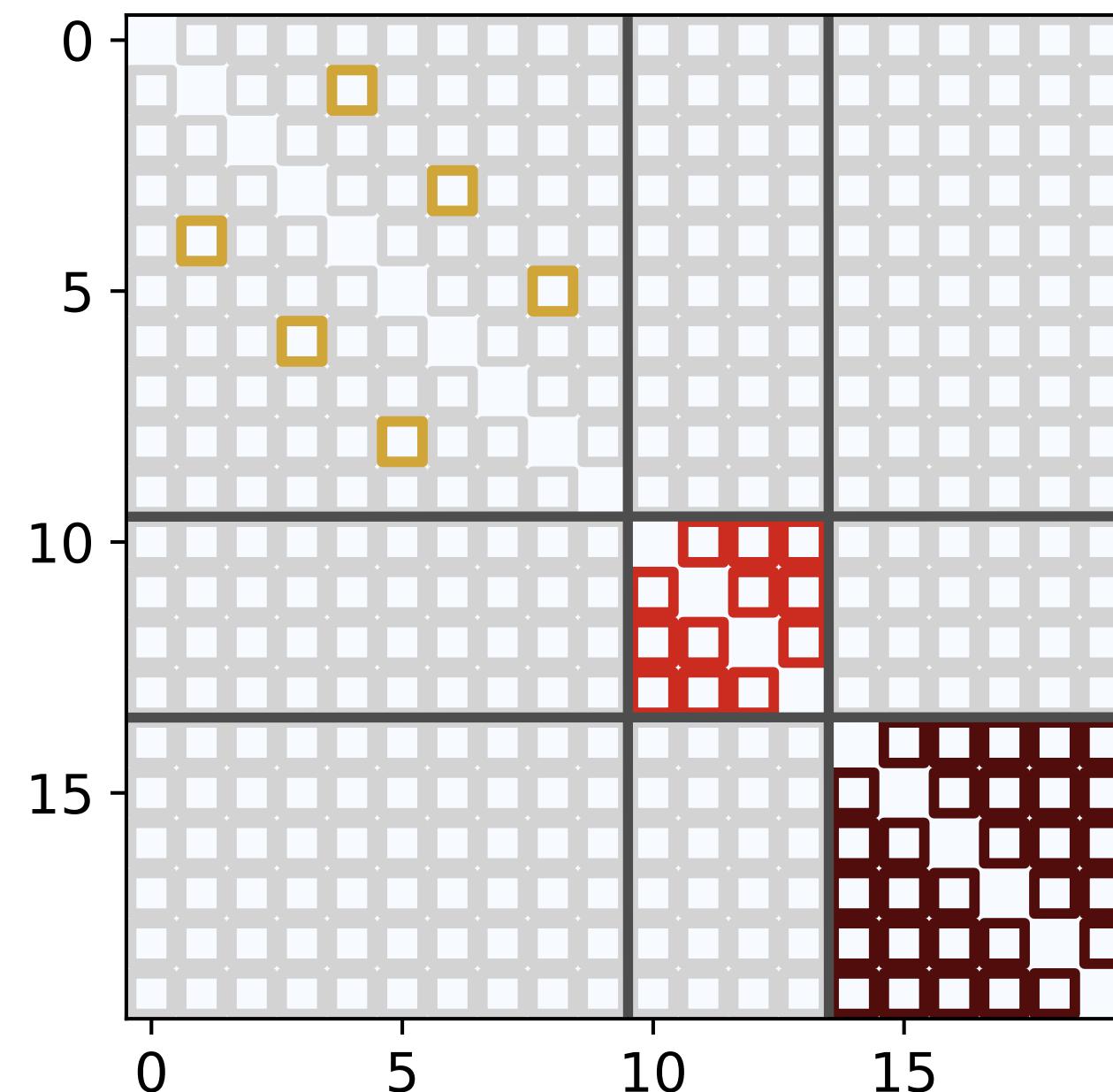
█ Isolated pairs      █ Large collective  
█ Small collective      █ Non interacting



# I. Undersampling-induced biases (Kleeorin *et al.* 2023)

*Inference as function of regularization strength*

█ Isolated pairs      █ Large collective  
█ Small collective      █ Non interacting

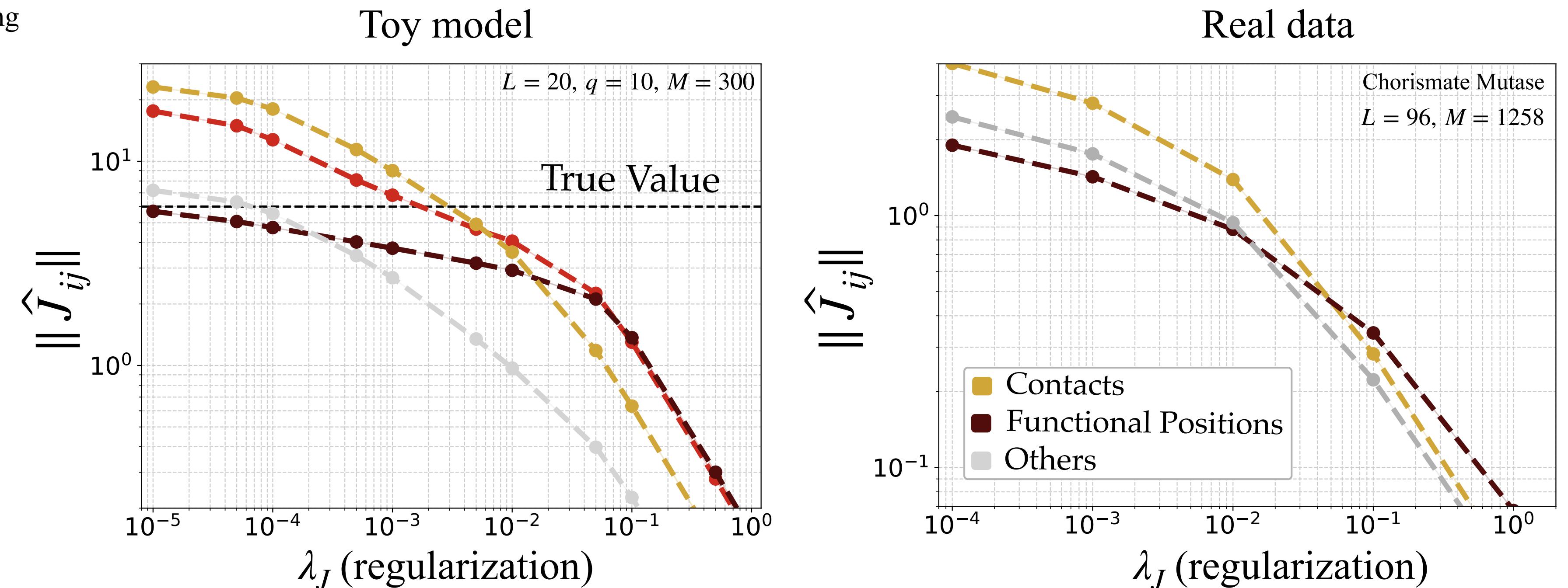
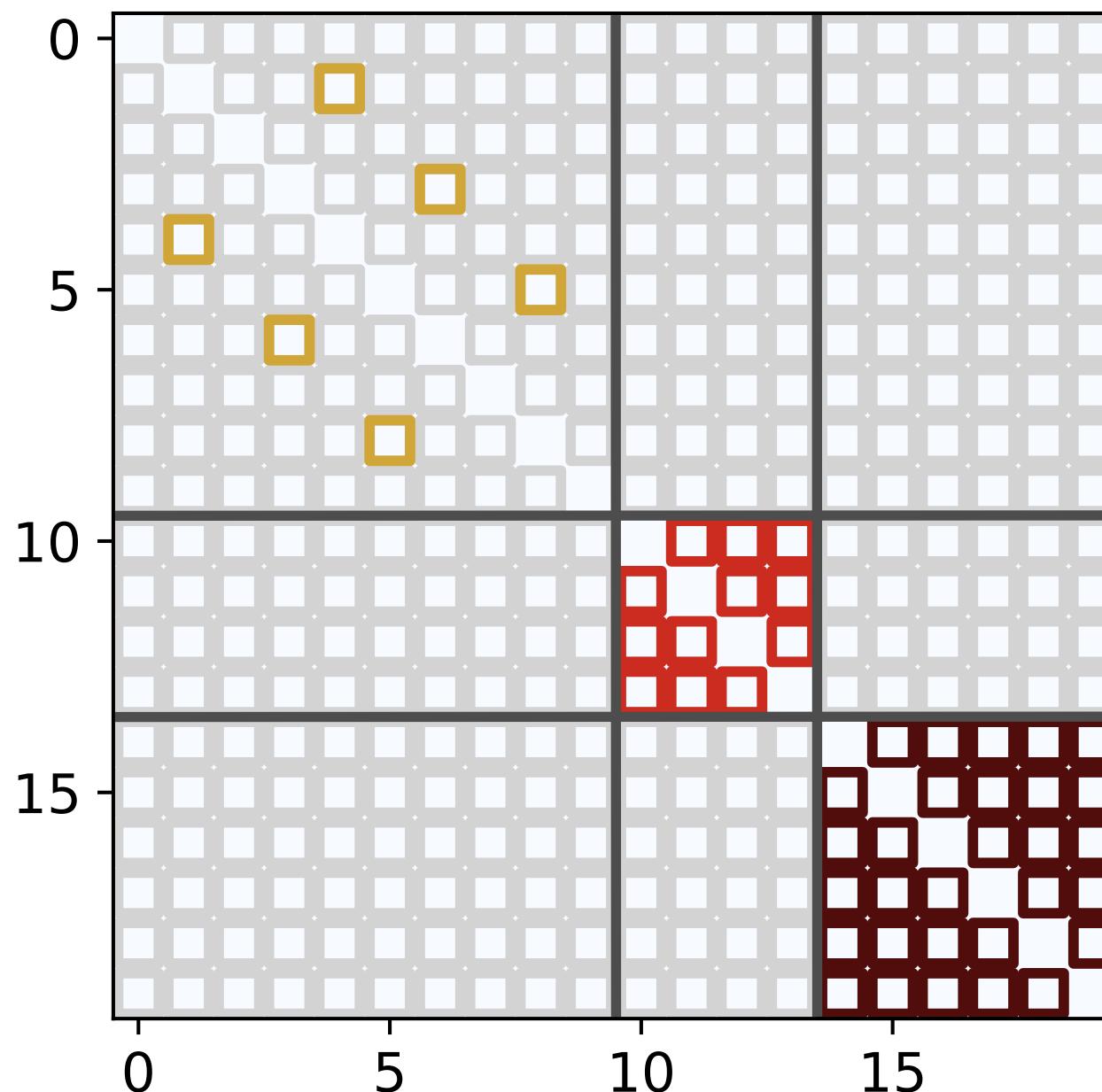


- Systematic bias between the estimation of collective modes and isolated interactions

# I. Undersampling-induced biases (Kleeorin *et al.* 2023)

*Inference as function of regularization strength*

█ Isolated pairs      █ Large collective  
█ Small collective      █ Non interacting



- ▶ Systematic bias between the estimation of collective modes and isolated interactions
- ▶ Relevance for real protein data

# Generative Capacity of the Boltzmann Machine (Russ *et al.* 2020)

*Application to the Chorismate Mutase family*

## Fidelity

*Do the artificial proteins share key properties with those observed in the training data?*

## Novelty

*How much do the artificial sequences differ from the training data?  
→ Generalization capacity*

## Diversity

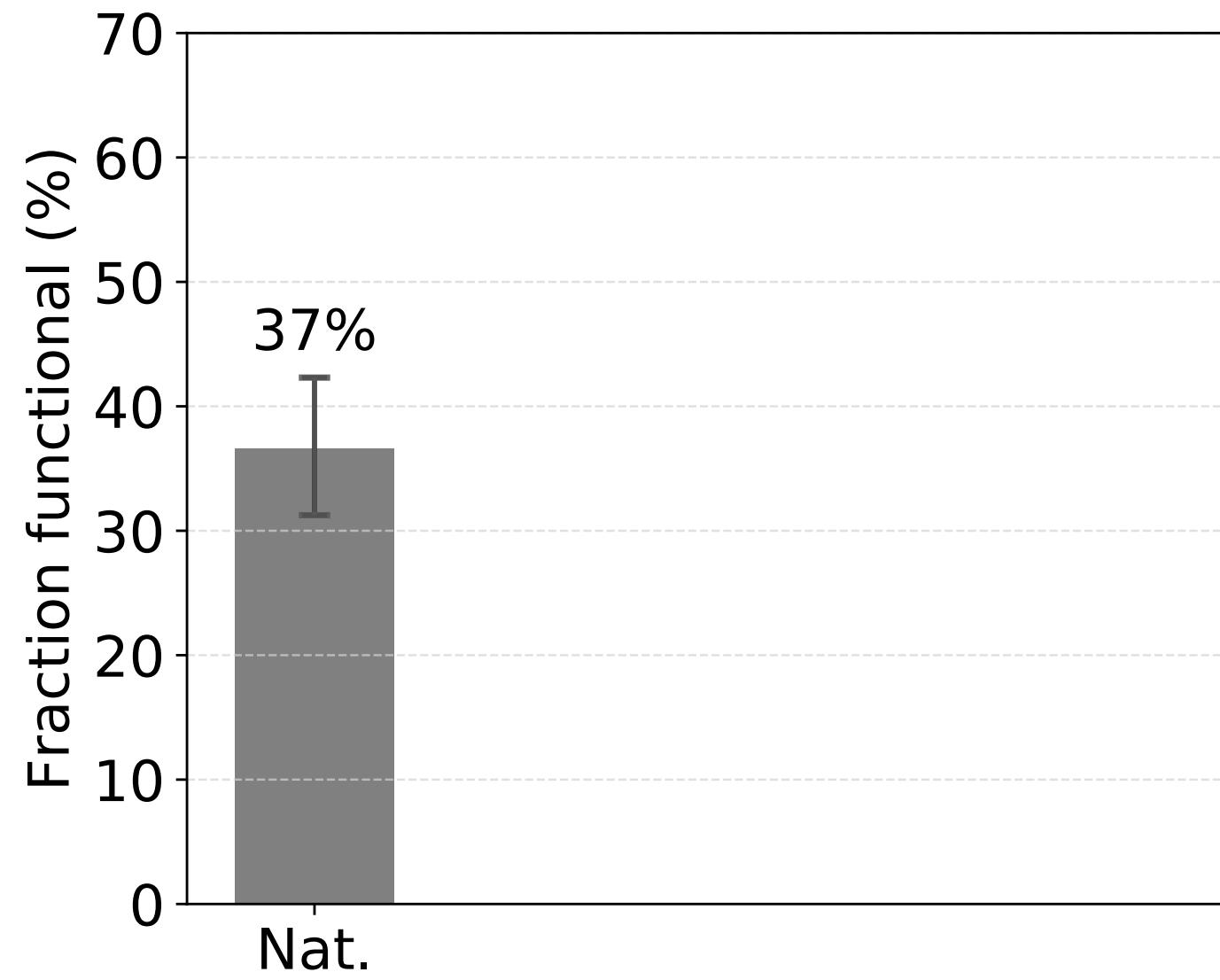
*Do the artificial sequences span the same range of variability as the training data?*

## II. Generative Capacity of the Boltzmann Machine (Russ *et al.* 2020)

*Application to the Chorismate Mutase family*

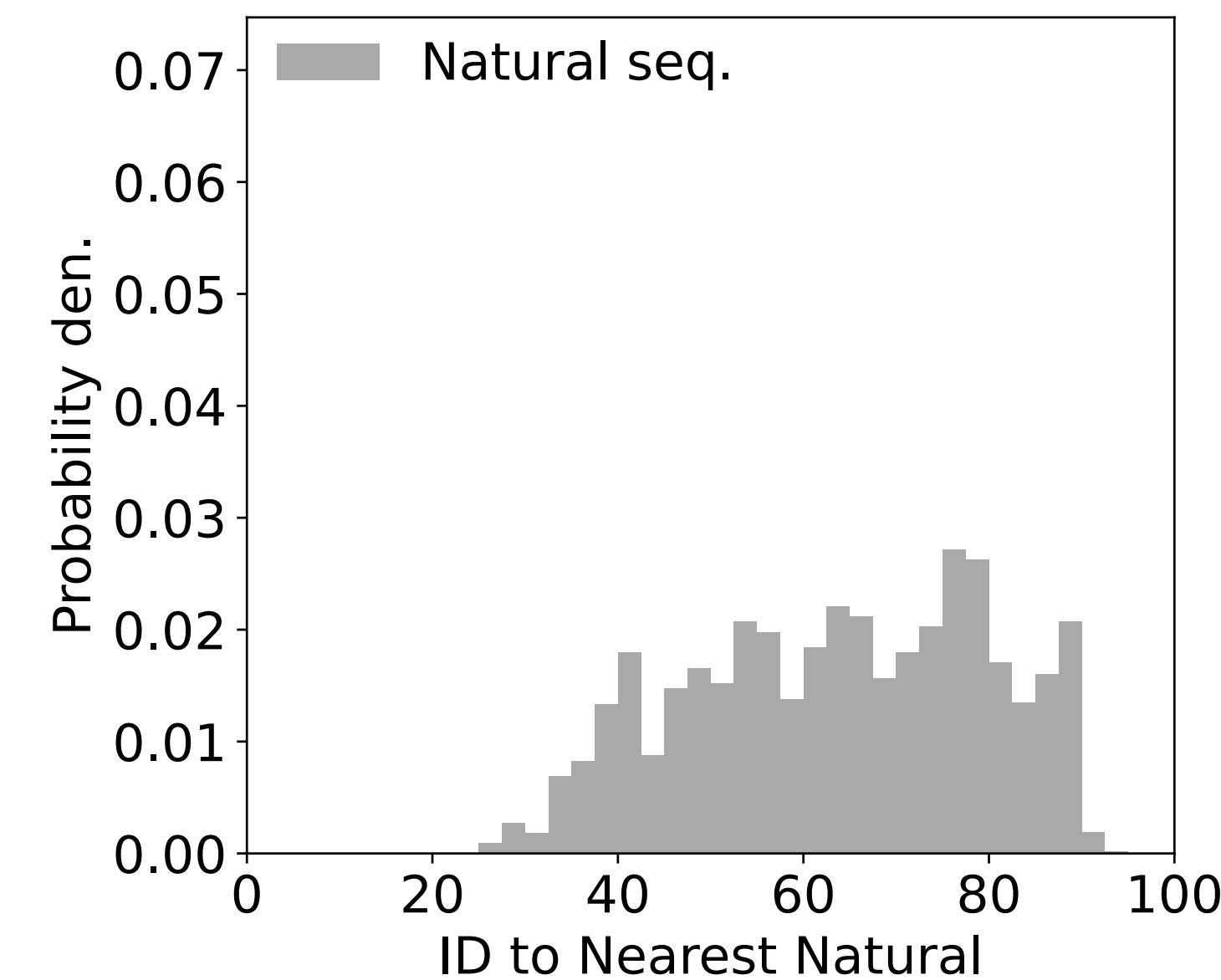
### Fidelity

*Fraction of functional sequences*



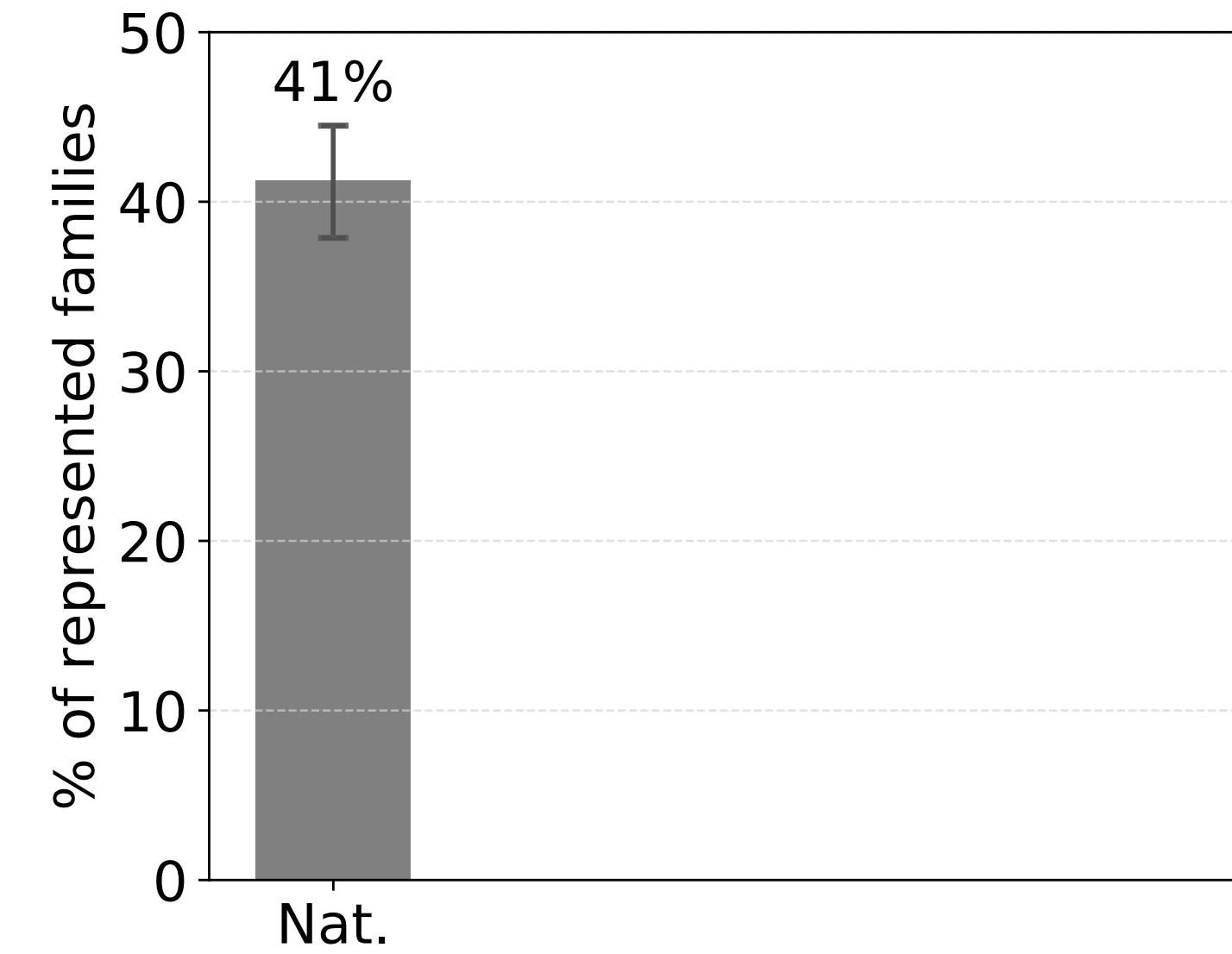
### Novelty

*Distribution of identity to the nearest natural sequence*



### Diversity

*% of taxonomic families represented*



$L_2$  regularization:  $\lambda = 0.01$

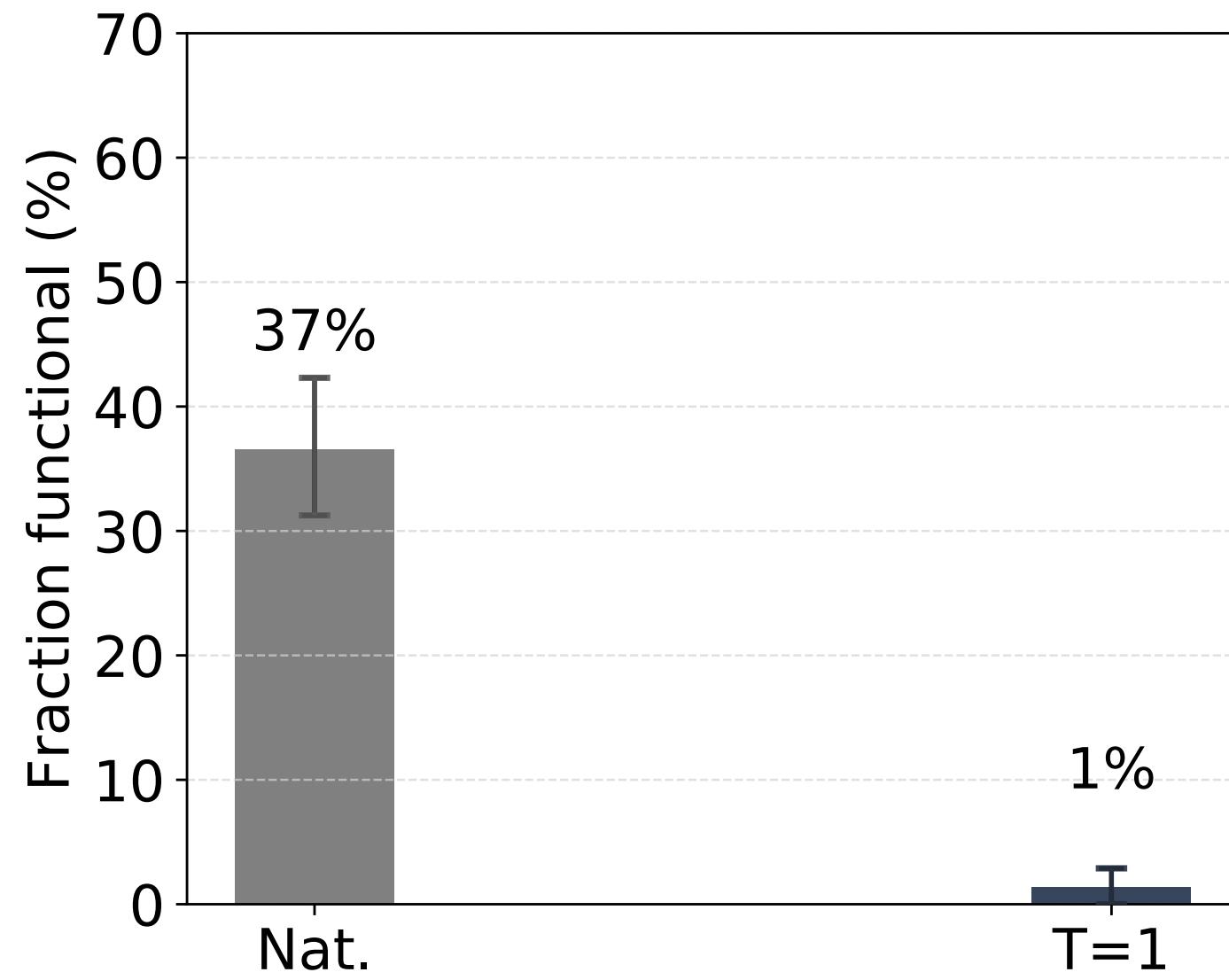
Figures made from Russ et al., *Science* 2020 data (sequences tested in *E. Coli*)

## II. Generative Capacity of the Boltzmann Machine (Russ *et al.* 2020)

*Application to the Chorismate Mutase family*

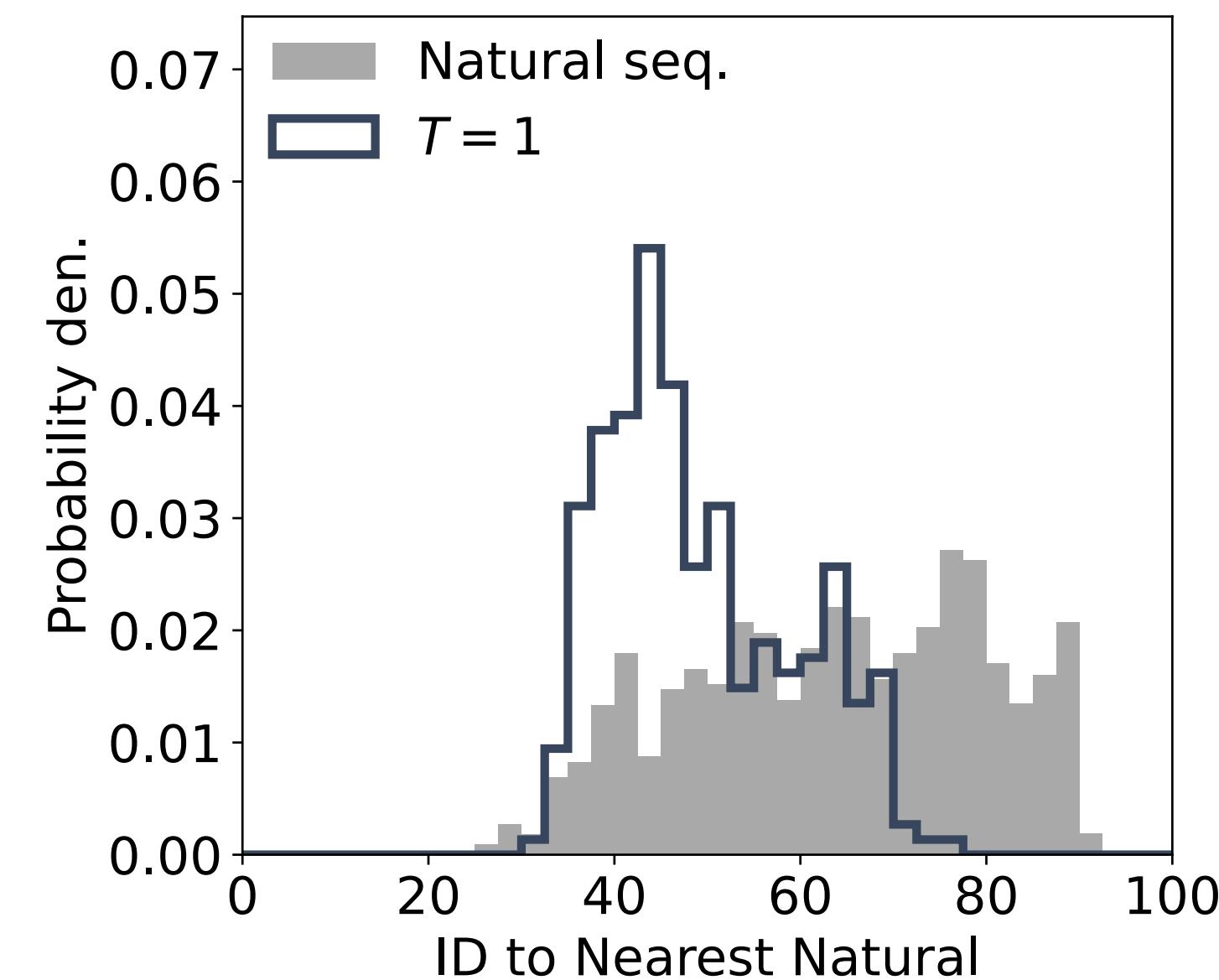
### Fidelity

*Fraction of functional sequences*



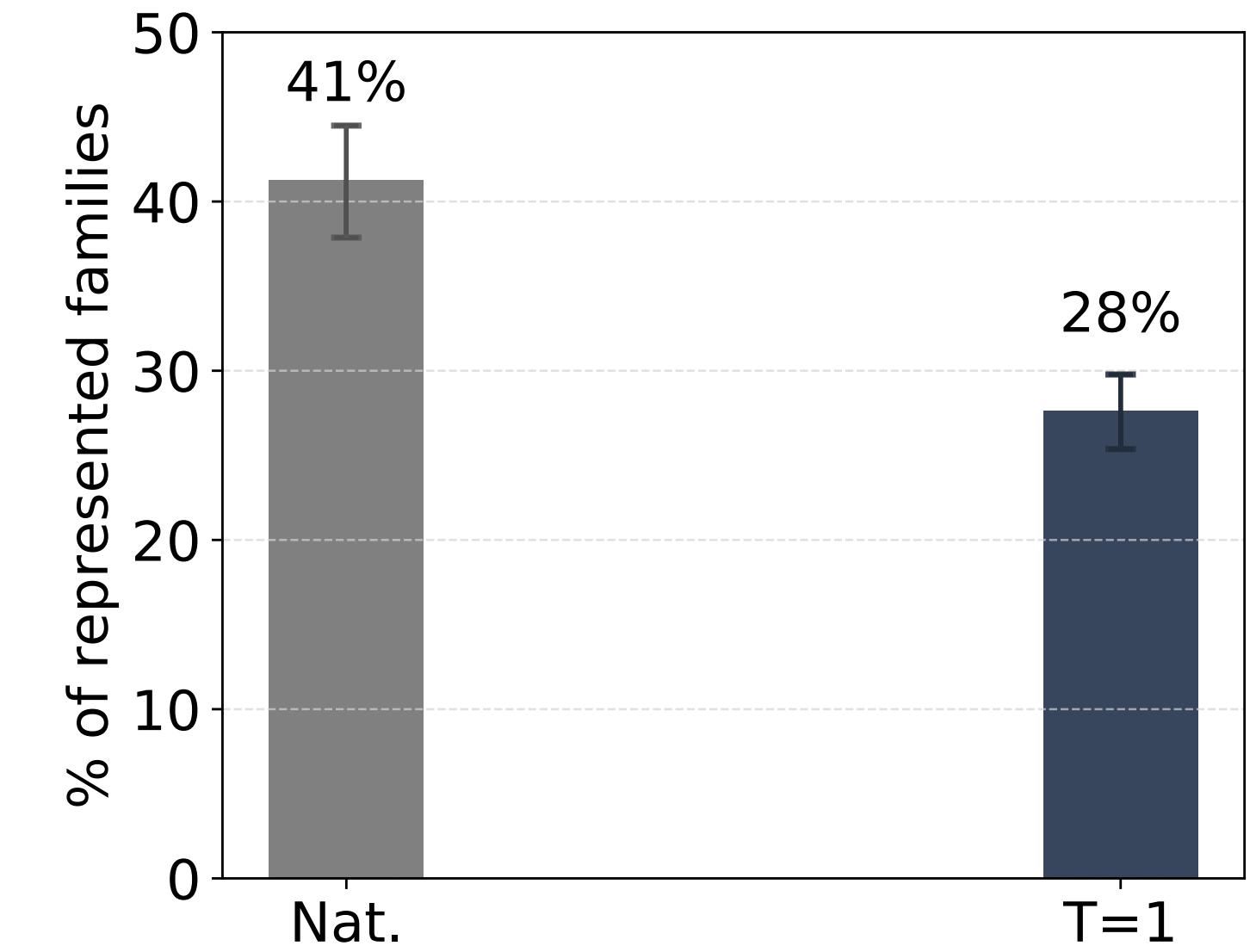
### Novelty

*Distribution of identity to the nearest natural sequence*



### Diversity

*% of taxonomic families represented*



$L_2$  regularization:  $\lambda = 0.01$

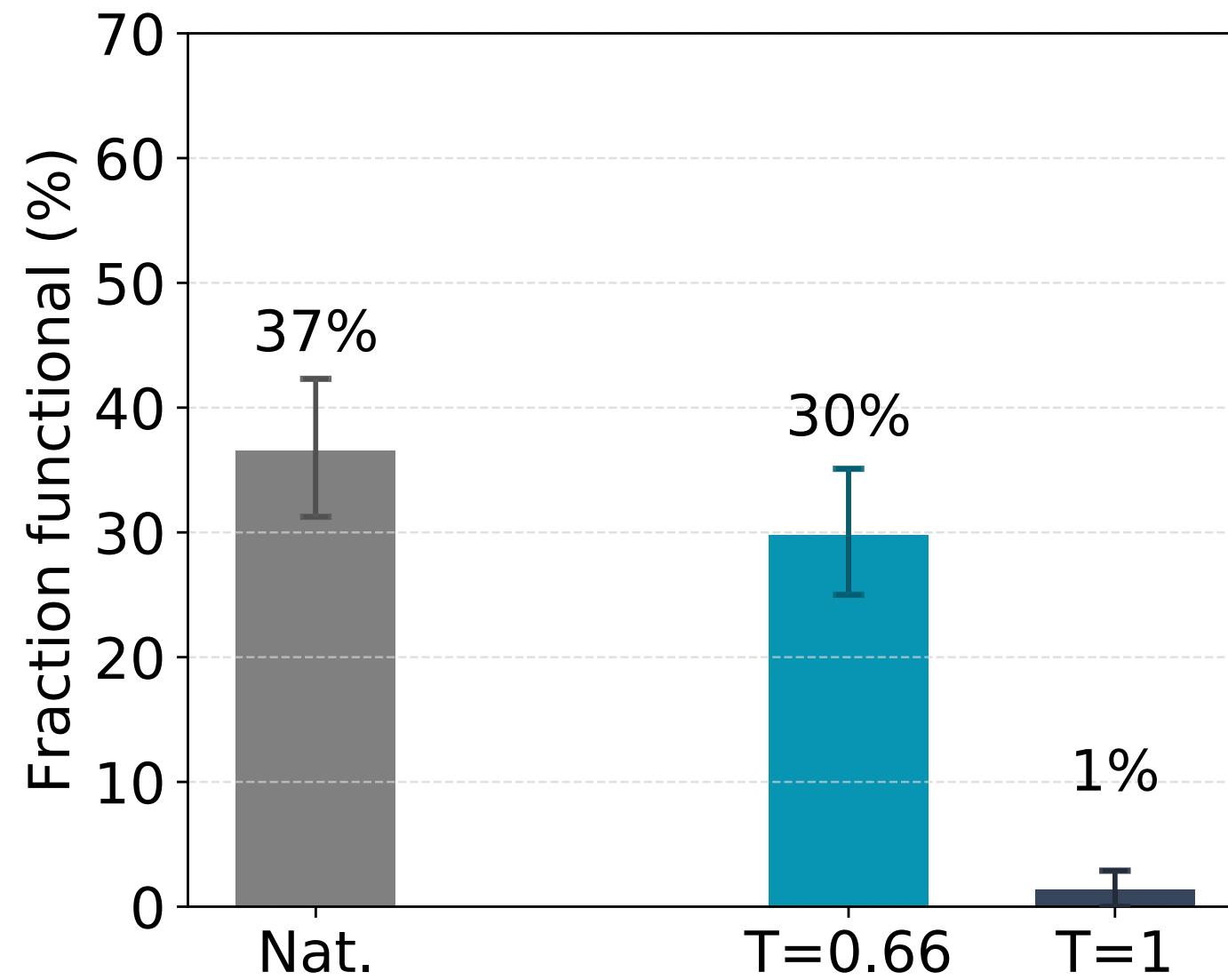
Figures made from Russ et al., *Science* 2020 data (sequences tested in *E. Coli*)

## II. Generative Capacity of the Boltzmann Machine (Russ *et al.* 2020)

*Application to the Chorismate Mutase family*

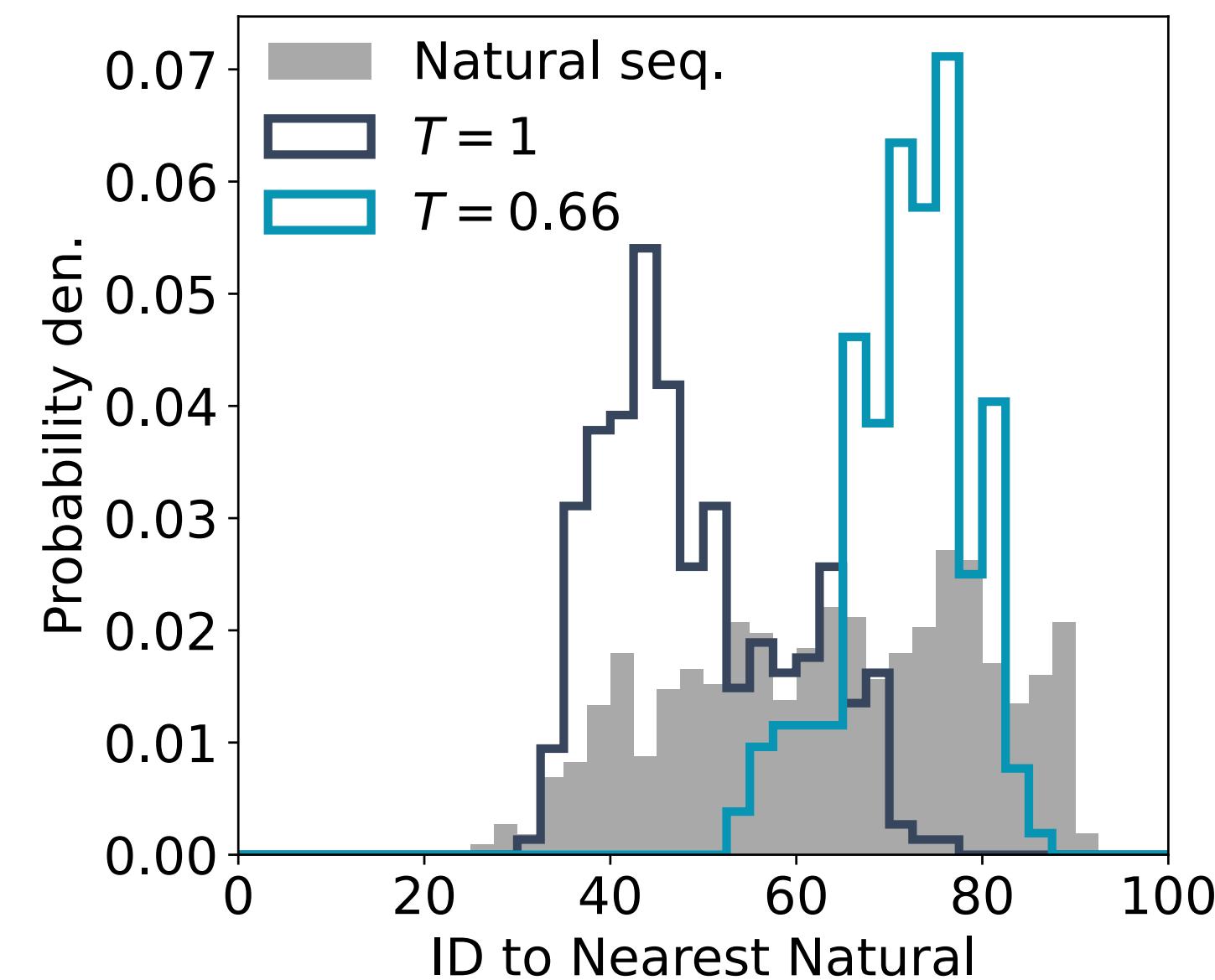
### Fidelity

*Fraction of functional sequences*



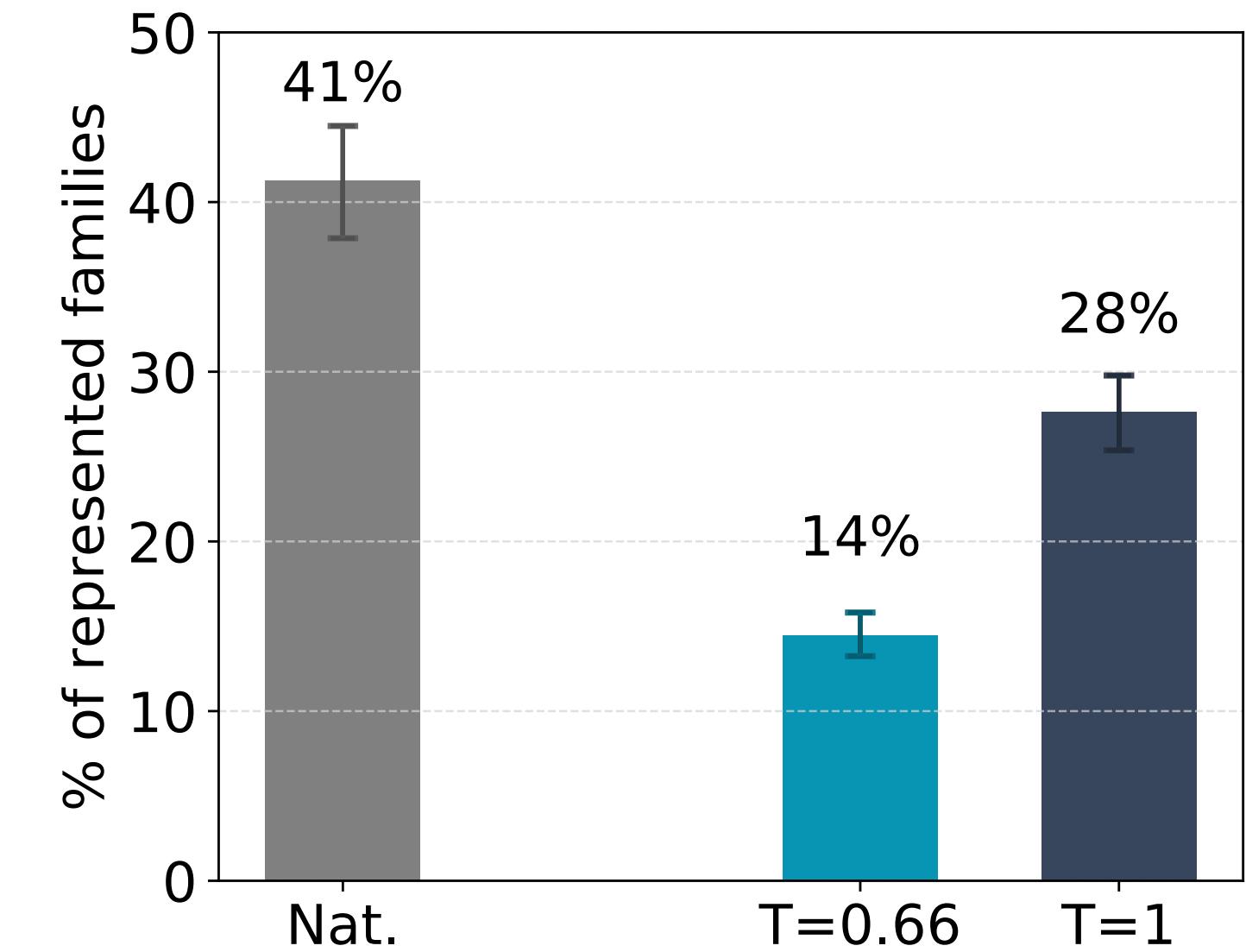
### Novelty

*Distribution of identity to the nearest natural sequence*



### Diversity

*% of taxonomic families represented*



$L_2$  regularization:  $\lambda = 0.01$

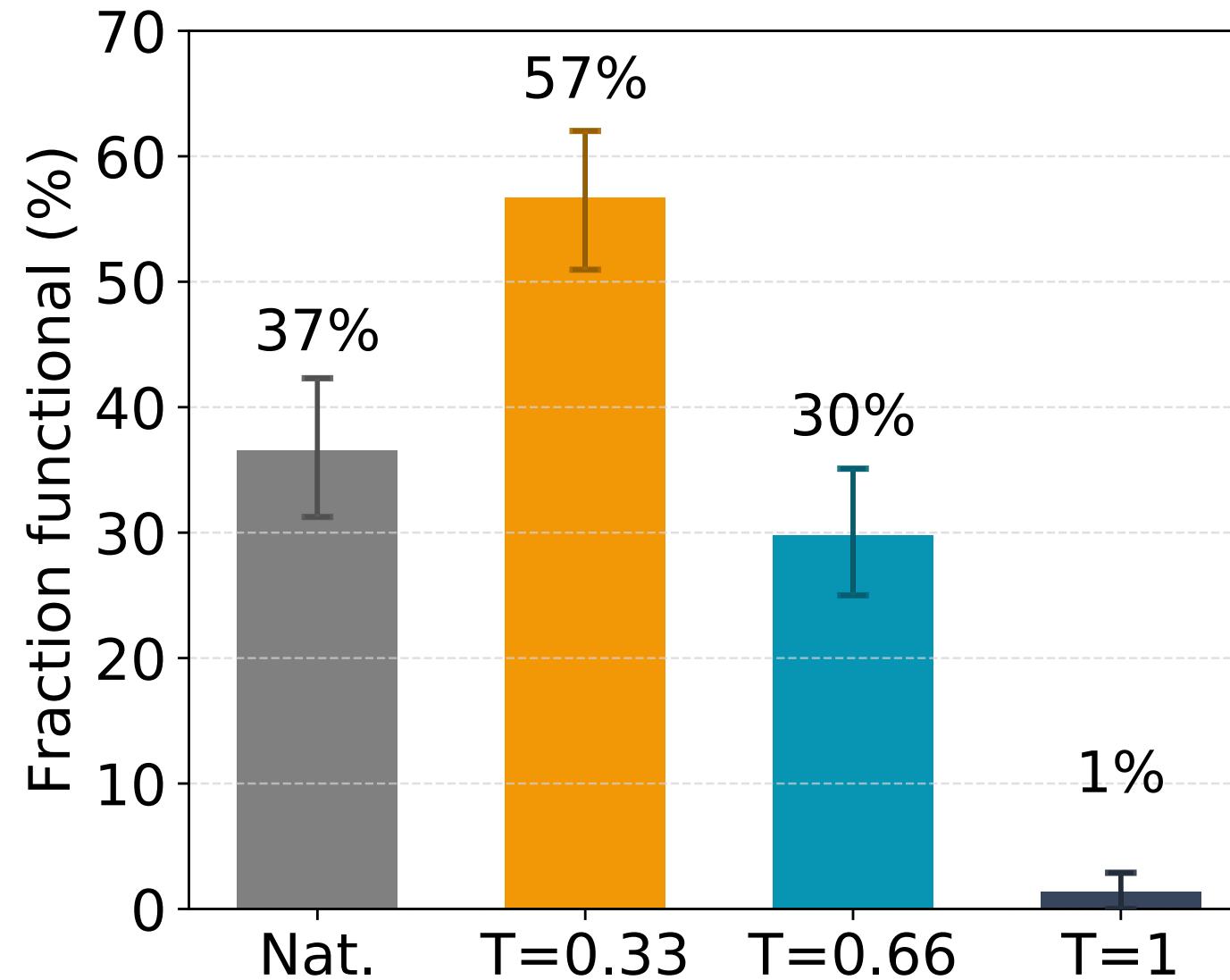
Figures made from Russ et al., *Science* 2020 data (sequences tested in *E. Coli*)

## II. Generative Capacity of the Boltzmann Machine (Russ *et al.* 2020)

*Application to the Chorismate Mutase family*

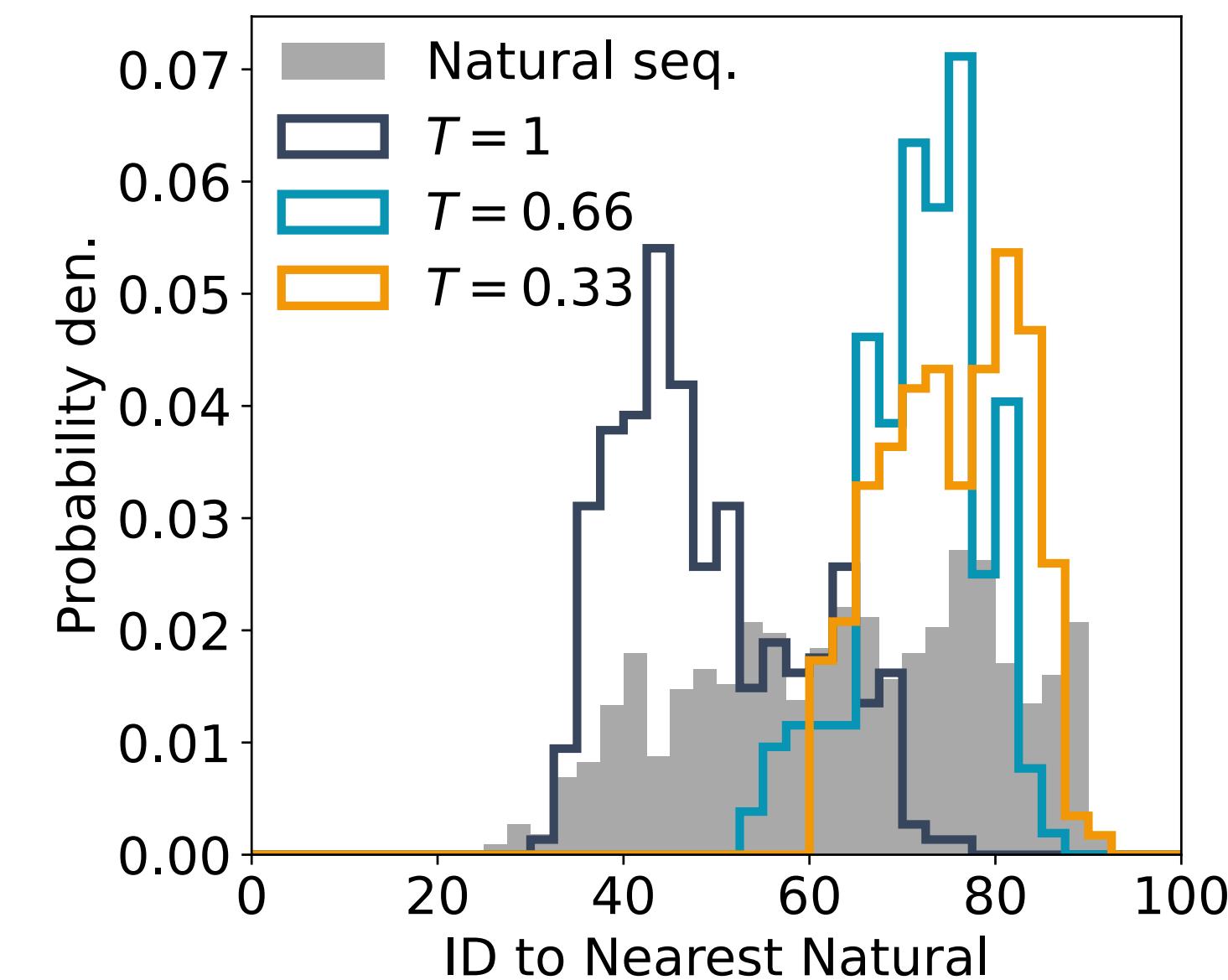
### Fidelity

*Fraction of functional sequences*



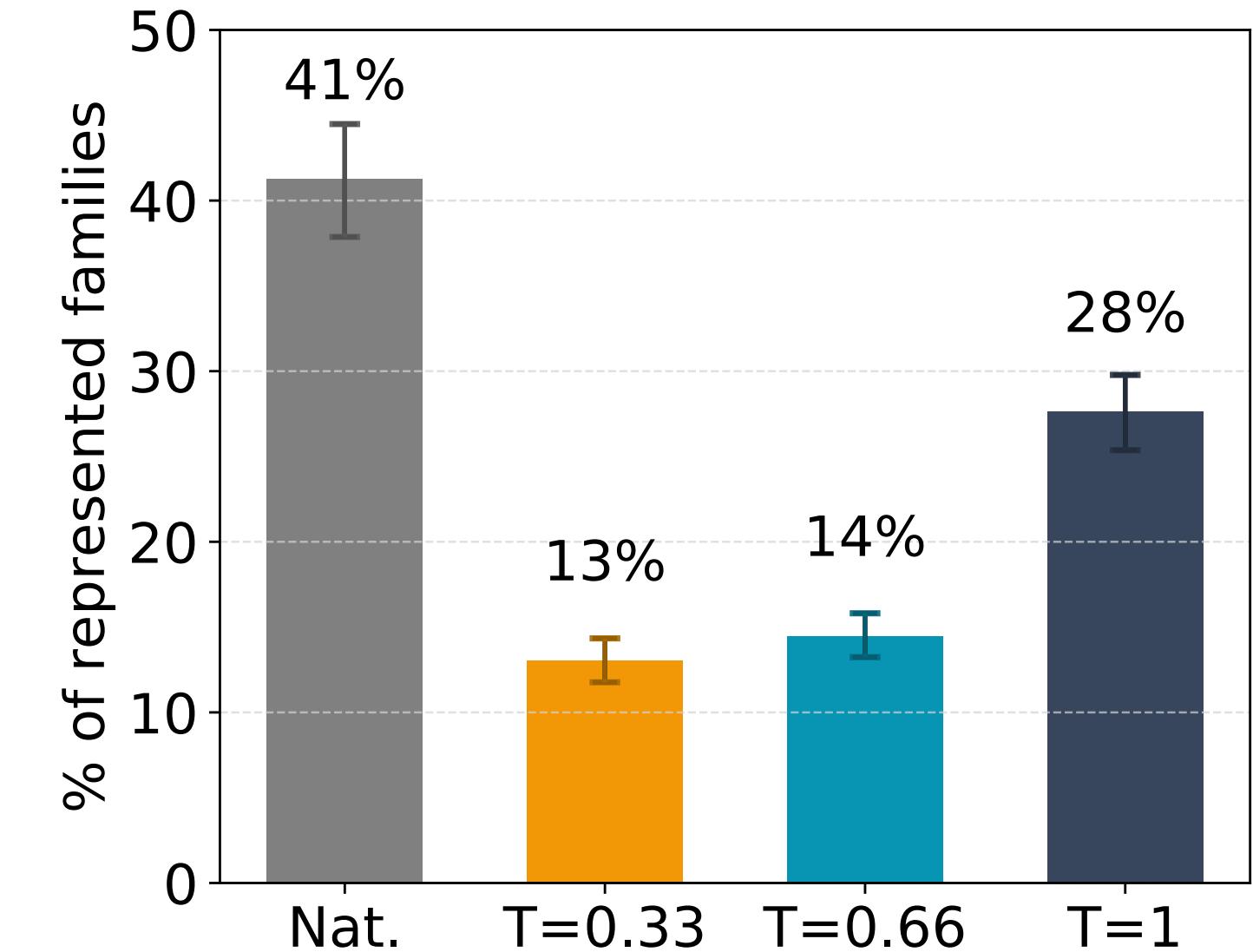
### Novelty

*Distribution of identity to the nearest natural sequence*



### Diversity

*% of taxonomic families represented*



- ▶ Functional sequences requires low-temperature sampling:  $P(\{\sigma_i\}_{i=1,\dots,L}) \sim e^{-\frac{E(\mathbf{h}, \mathbf{J})}{T}}$  with  $T < 1$
- ▶ The gain in functionality comes at the cost of **reduced diversity** and **novelty**

$L_2$  regularization:  $\lambda = 0.01$

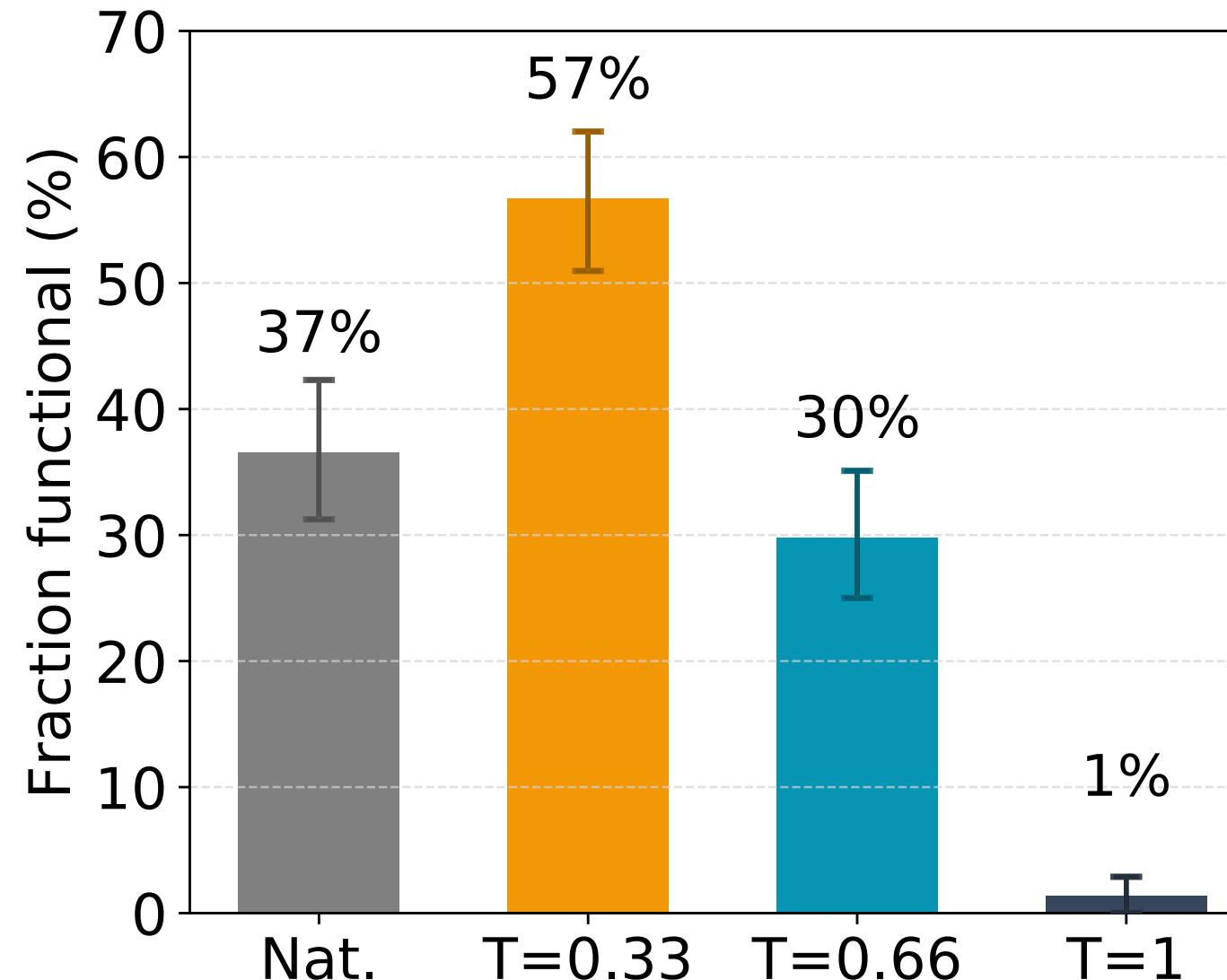
Figures made from Russ et al., *Science* 2020 data (sequences tested in *E. Coli*)

## II. Generative Capacity of the Boltzmann Machine (Russ *et al.* 2020)

*Application to the Chorismate Mutase family*

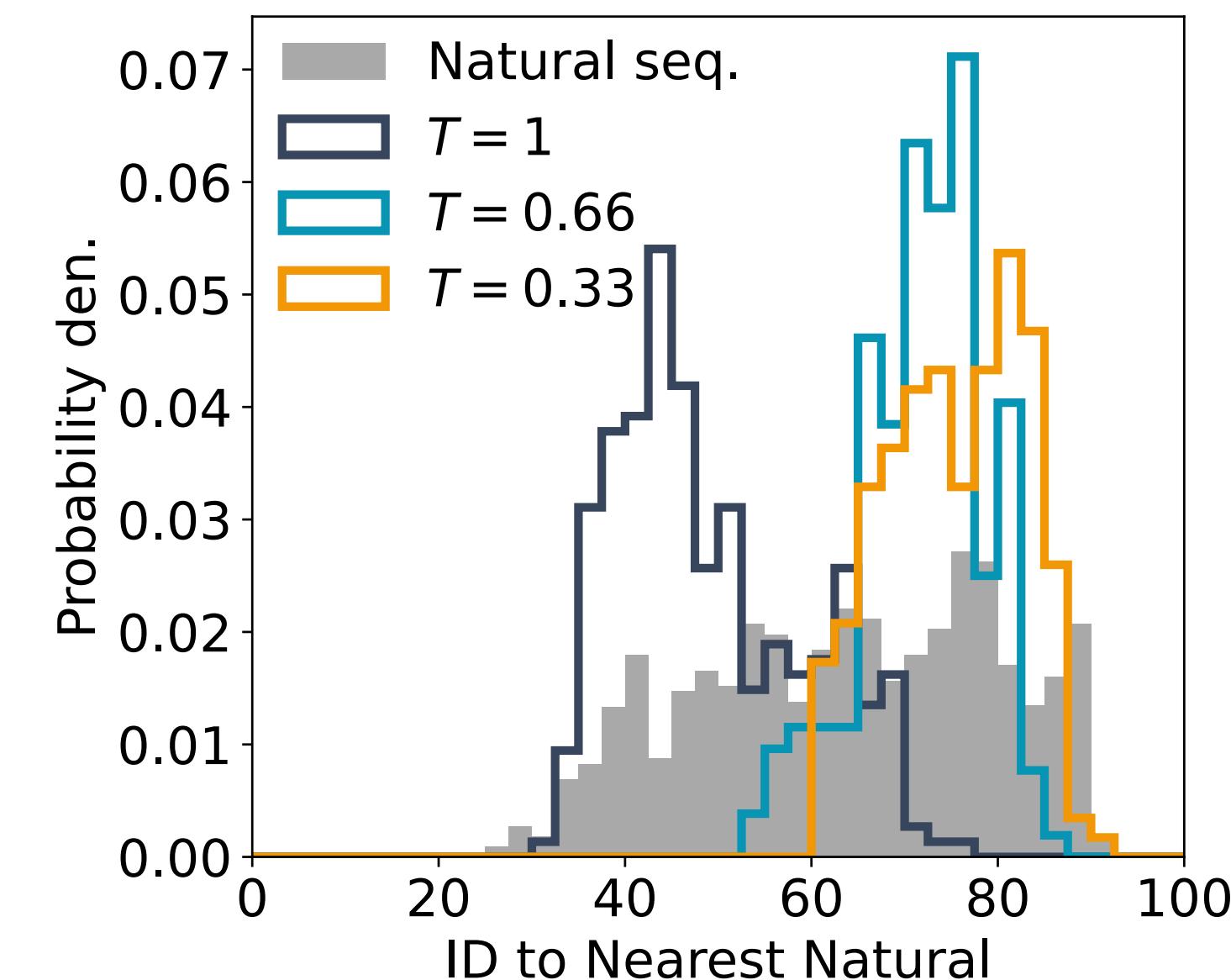
### Fidelity

*Fraction of functional sequences*



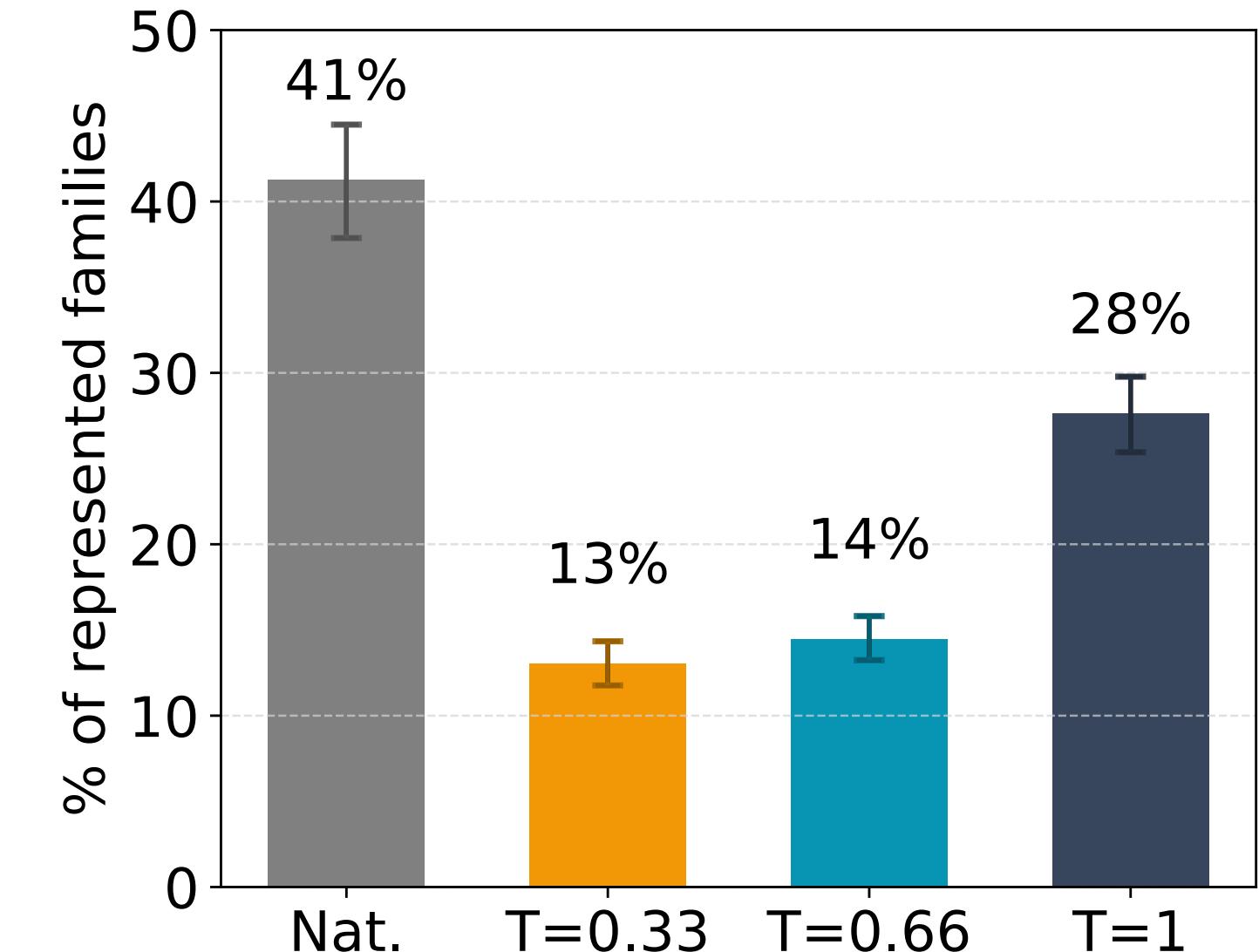
### Novelty

*Distribution of identity to the nearest natural sequence*



### Diversity

*% of taxonomic families represented*



- ▶ Functional sequences requires low-temperature sampling:  $P(\{\sigma_i\}_{i=1,\dots,L}) \sim e^{-\frac{E(\mathbf{h}, \mathbf{J})}{T}}$  with  $T < 1$
- ▶ The gain in functionality comes at the cost of **reduced diversity** and **novelty**
- ▶ Success may depend on the experimental assay

$L_2$  regularization:  $\lambda = 0.01$

Figures made from Russ et al., *Science* 2020 data (sequences tested in *E. Coli*)

### III. Stochastic Boltzmann Machine

*Assess model performance with a toy model*

$$\theta = \{J, h\} \quad f(\theta) = \frac{1}{M} \sum_{m=1}^M \log P(\sigma^{(m)} \mid \theta) - \lambda_J \|J\|^2 - \lambda_h \|h\|^2$$

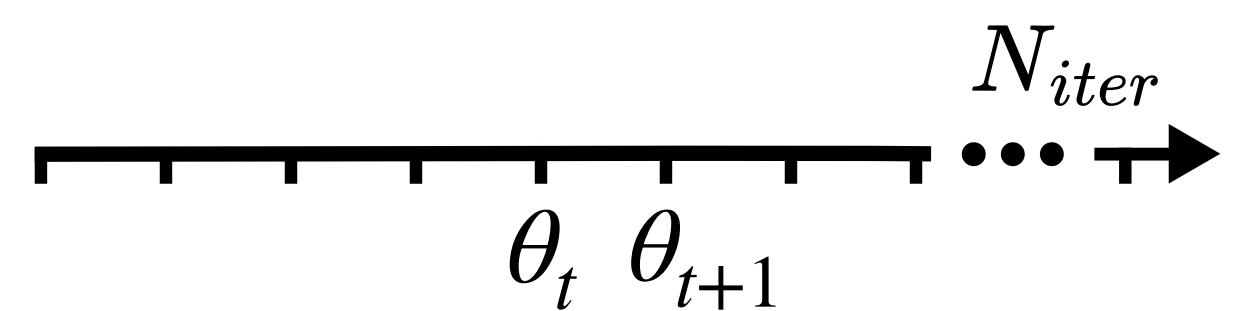
### III. Stochastic Boltzmann Machine

*Assess model performance with a toy model*

$$\theta = \{J, h\} \quad f(\theta) = \frac{1}{M} \sum_{m=1}^M \log P(\sigma^{(m)} \mid \theta) - \lambda_J \|J\|^2 - \lambda_h \|h\|^2$$

$$\theta_{t+1} = \theta_t - \eta_t p_t$$

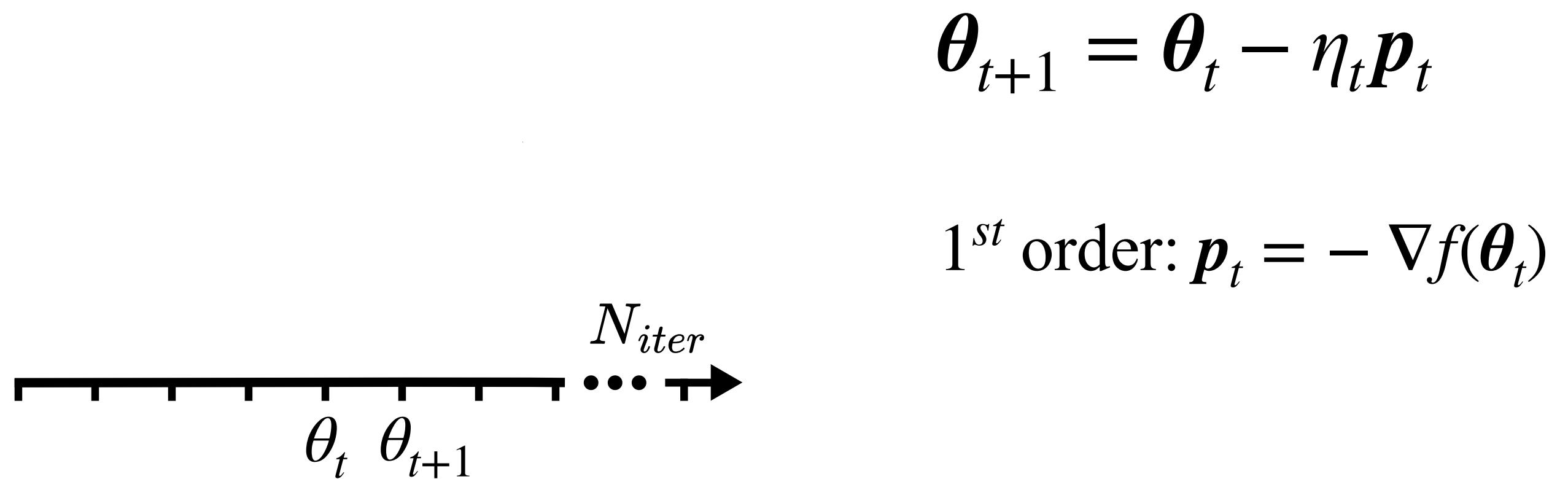
1<sup>st</sup> order:  $p_t = -\nabla f(\theta_t)$



### III. Stochastic Boltzmann Machine

*Assess model performance with a toy model*

$$\theta = \{J, h\} \quad f(\theta) = \frac{1}{M} \sum_{m=1}^M \log P(\sigma^{(m)} \mid \theta) - \lambda_J \|J\|^2 - \lambda_h \|h\|^2$$

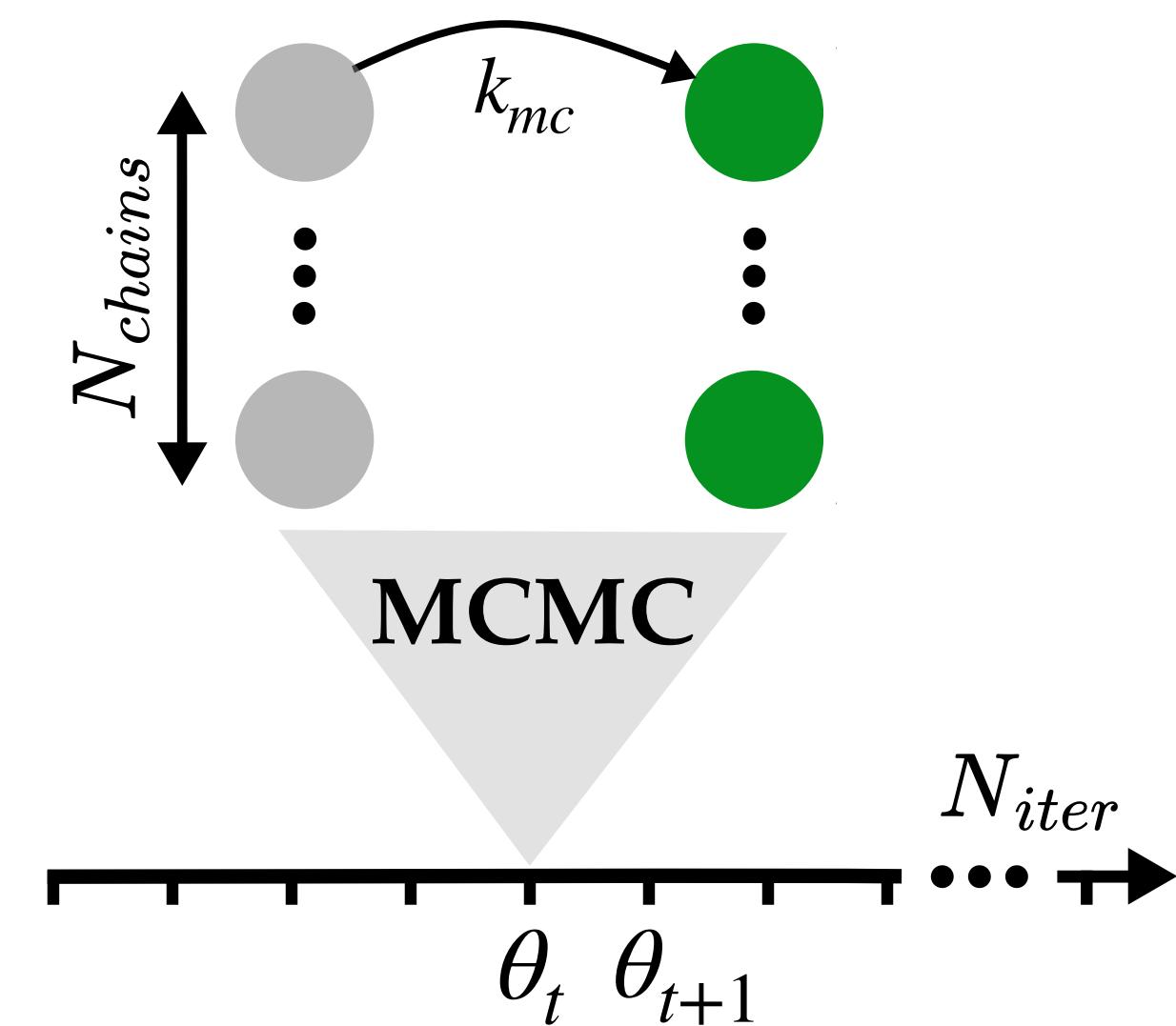


$$\frac{\partial f(\boldsymbol{\theta})}{\partial J_{ij}(a, b)} = \underbrace{f_{ij}(a, b)}_{\text{Empirical}} - \underbrace{\langle \delta(\sigma_i, a) \delta(\sigma_j, b) \rangle}_{\text{Model}} + \lambda_J J_{ij}(a, b)$$

### III. Stochastic Boltzmann Machine

*Assess model performance with a toy model*

$$\theta = \{J, h\} \quad f(\theta) = \frac{1}{M} \sum_{m=1}^M \log P(\sigma^{(m)} | \theta) - \lambda_J \|J\|^2 - \lambda_h \|h\|^2$$



$$\theta_{t+1} = \theta_t - \eta_t p_t$$

$$1^{st} \text{ order: } p_t = - \nabla f(\theta_t)$$

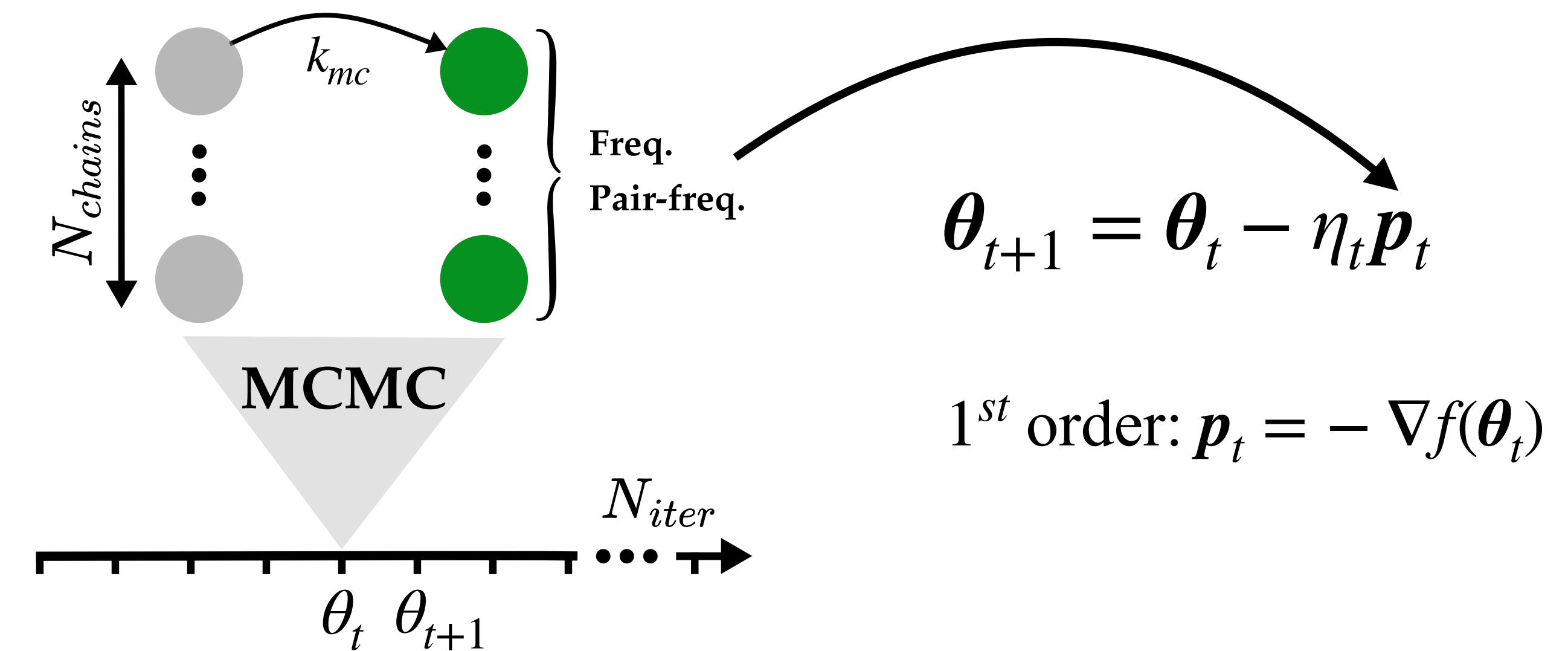
$$\frac{\partial f(\theta)}{\partial J_{ij}(a, b)} = f_{ij}(a, b) - \langle \delta(\sigma_i, a) \delta(\sigma_j, b) \rangle + \lambda_J J_{ij}(a, b)$$

Empirical                      Model

### III. Stochastic Boltzmann Machine

*Assess model performance with a toy model*

$$\theta = \{J, h\} \quad f(\theta) = \frac{1}{M} \sum_{m=1}^M \log P(\sigma^{(m)} | \theta) - \lambda_J \|J\|^2 - \lambda_h \|h\|^2$$

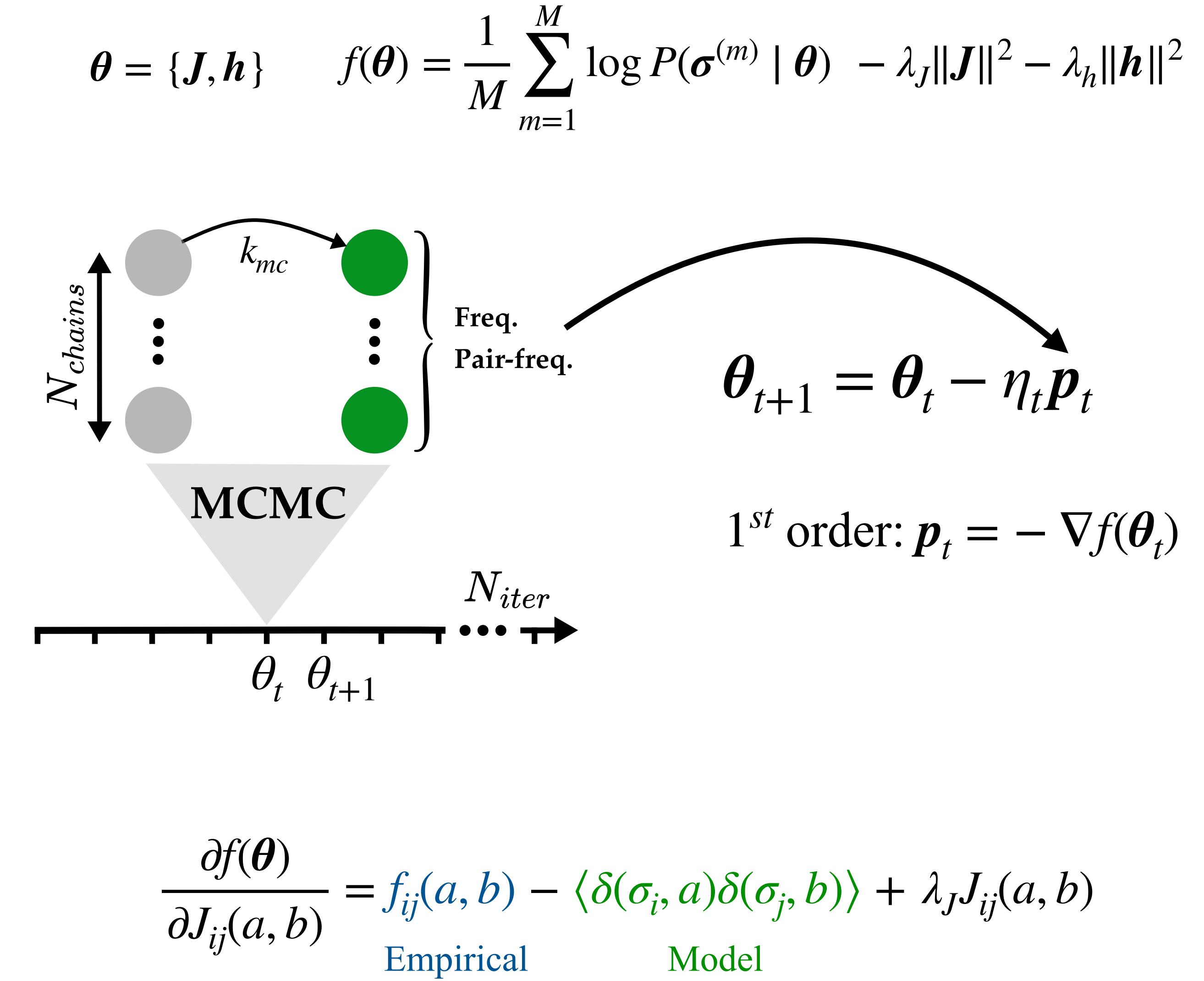
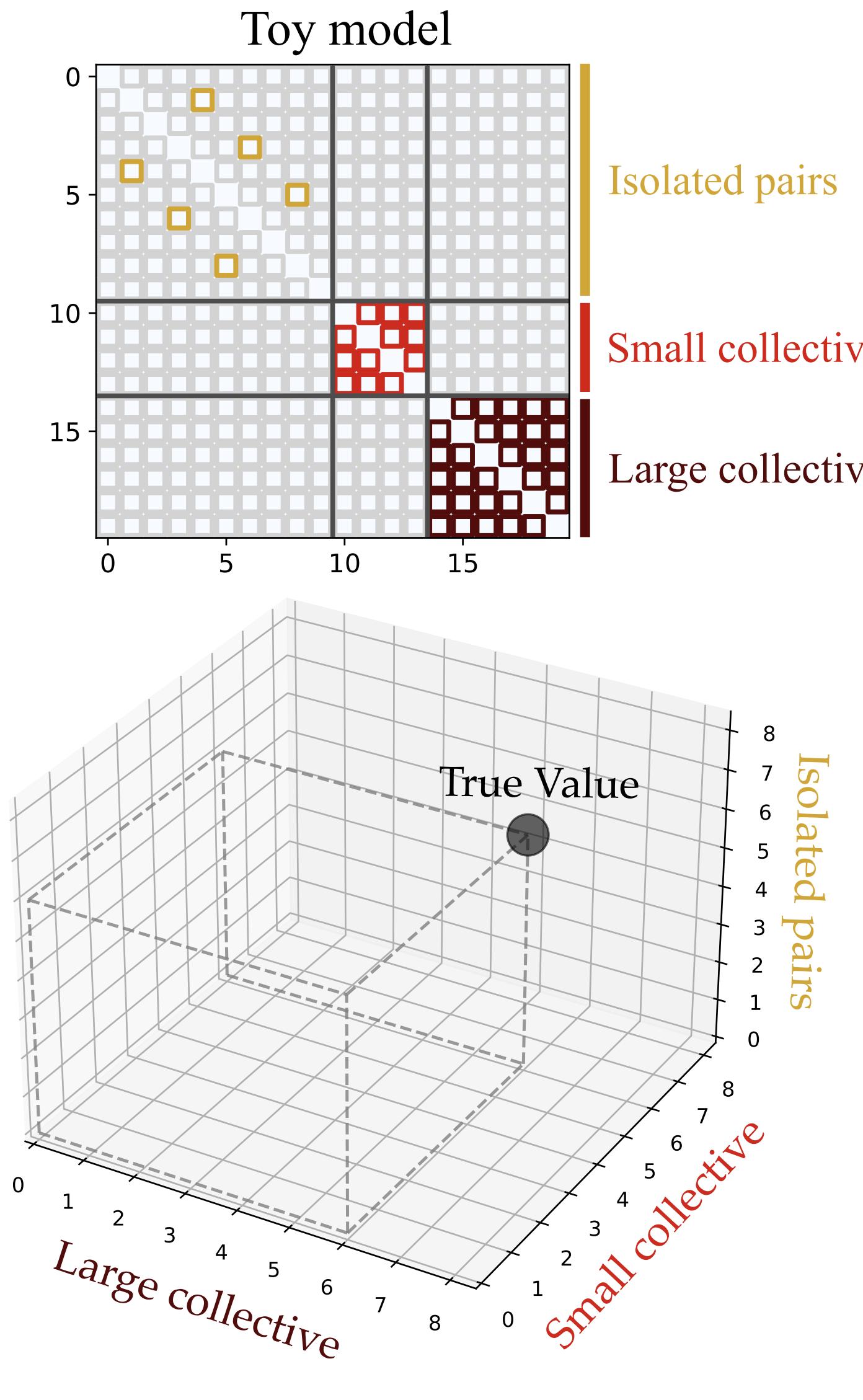


$$\frac{\partial f(\theta)}{\partial J_{ij}(a,b)} = f_{ij}(a,b) - \langle \delta(\sigma_i, a) \delta(\sigma_j, b) \rangle + \lambda_J J_{ij}(a,b)$$

Empirical                    Model

### III. Stochastic Boltzmann Machine

*Assess model performance with a toy model*



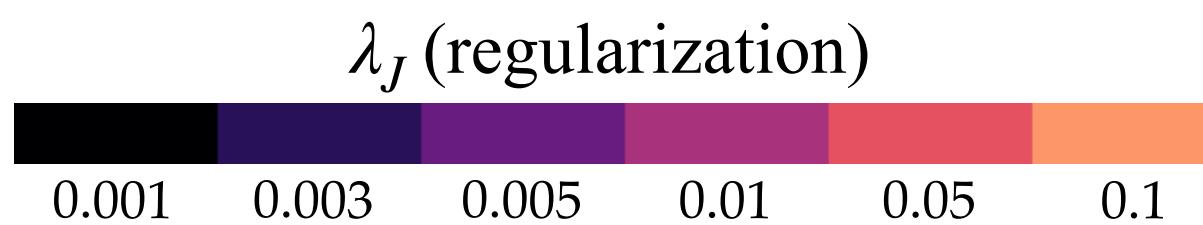
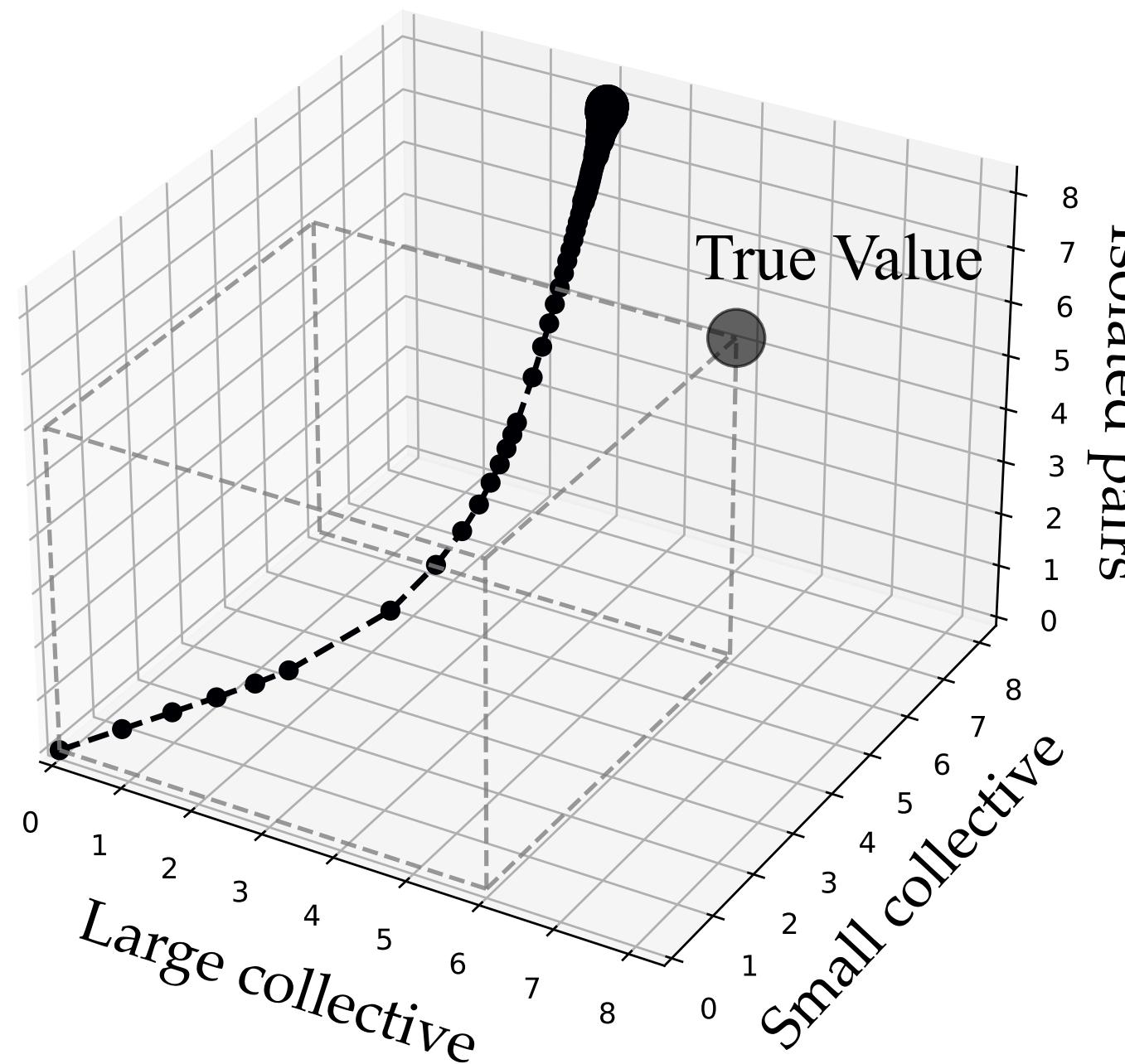
### III. Stochastic Boltzmann Machine

*Assess model performance with a toy model*

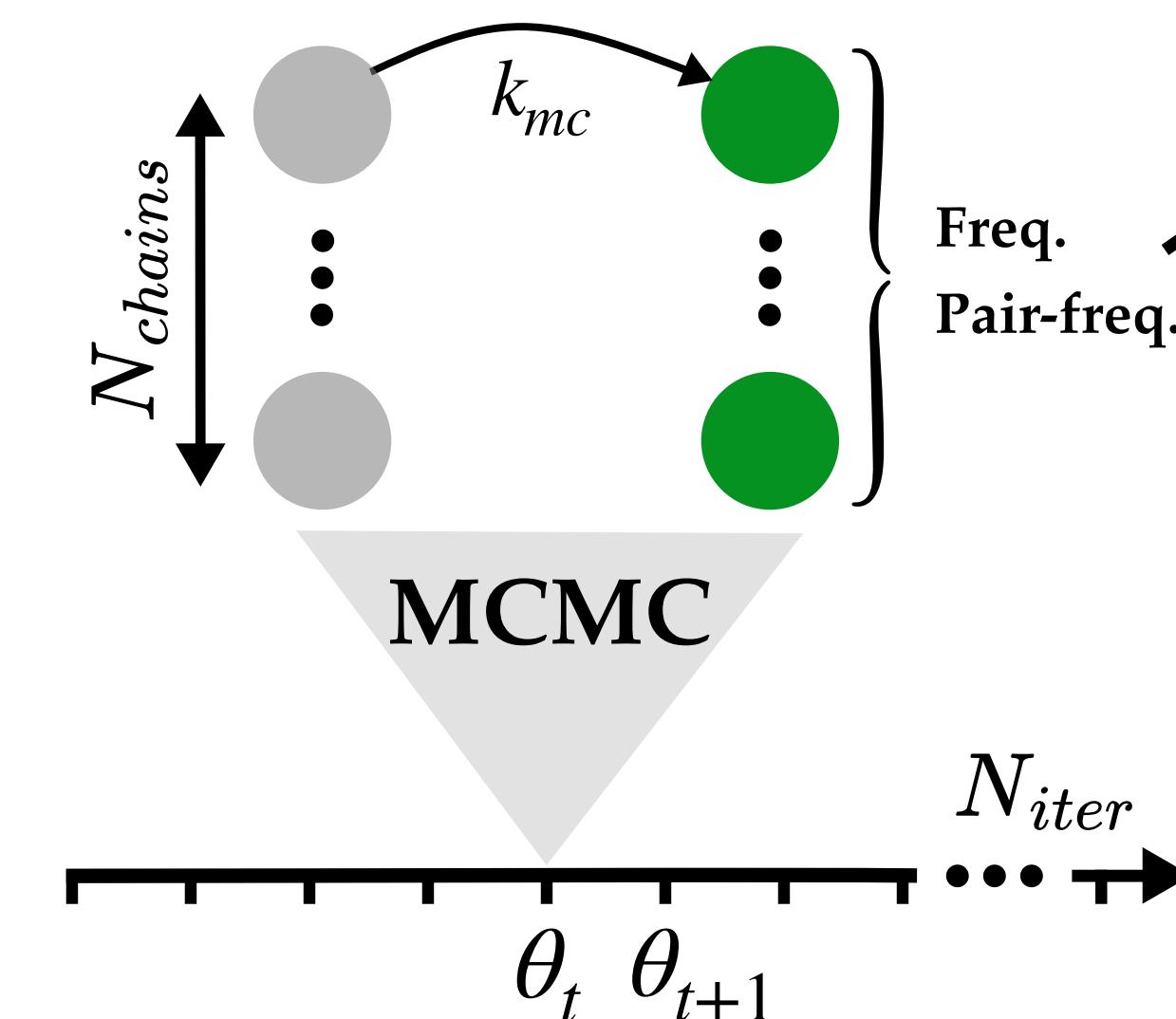
**BM**

$L_2$  regularization

$$1^{st} \text{ order: } \mathbf{p}_t = -\nabla f(\boldsymbol{\theta}_t)$$



$$\boldsymbol{\theta} = \{\mathbf{J}, \mathbf{h}\} \quad f(\boldsymbol{\theta}) = \frac{1}{M} \sum_{m=1}^M \log P(\boldsymbol{\sigma}^{(m)} | \boldsymbol{\theta}) - \lambda_J \|\mathbf{J}\|^2 - \lambda_h \|\mathbf{h}\|^2$$



$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \mathbf{p}_t$$

$$1^{st} \text{ order: } \mathbf{p}_t = -\nabla f(\boldsymbol{\theta}_t)$$

$$\frac{\partial f(\boldsymbol{\theta})}{\partial J_{ij}(a,b)} = f_{ij}(a,b) - \langle \delta(\sigma_i, a) \delta(\sigma_j, b) \rangle + \lambda_J J_{ij}(a,b)$$

Empirical                            Model

All models are inferred with  $k_{mc} = 10^5$ ,  $N_{iter} = 5000$

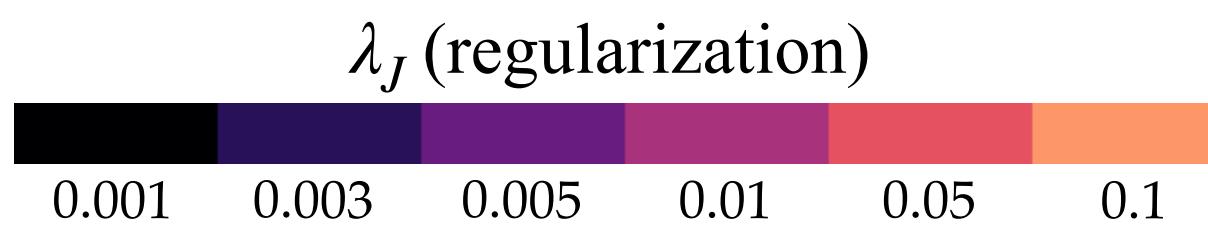
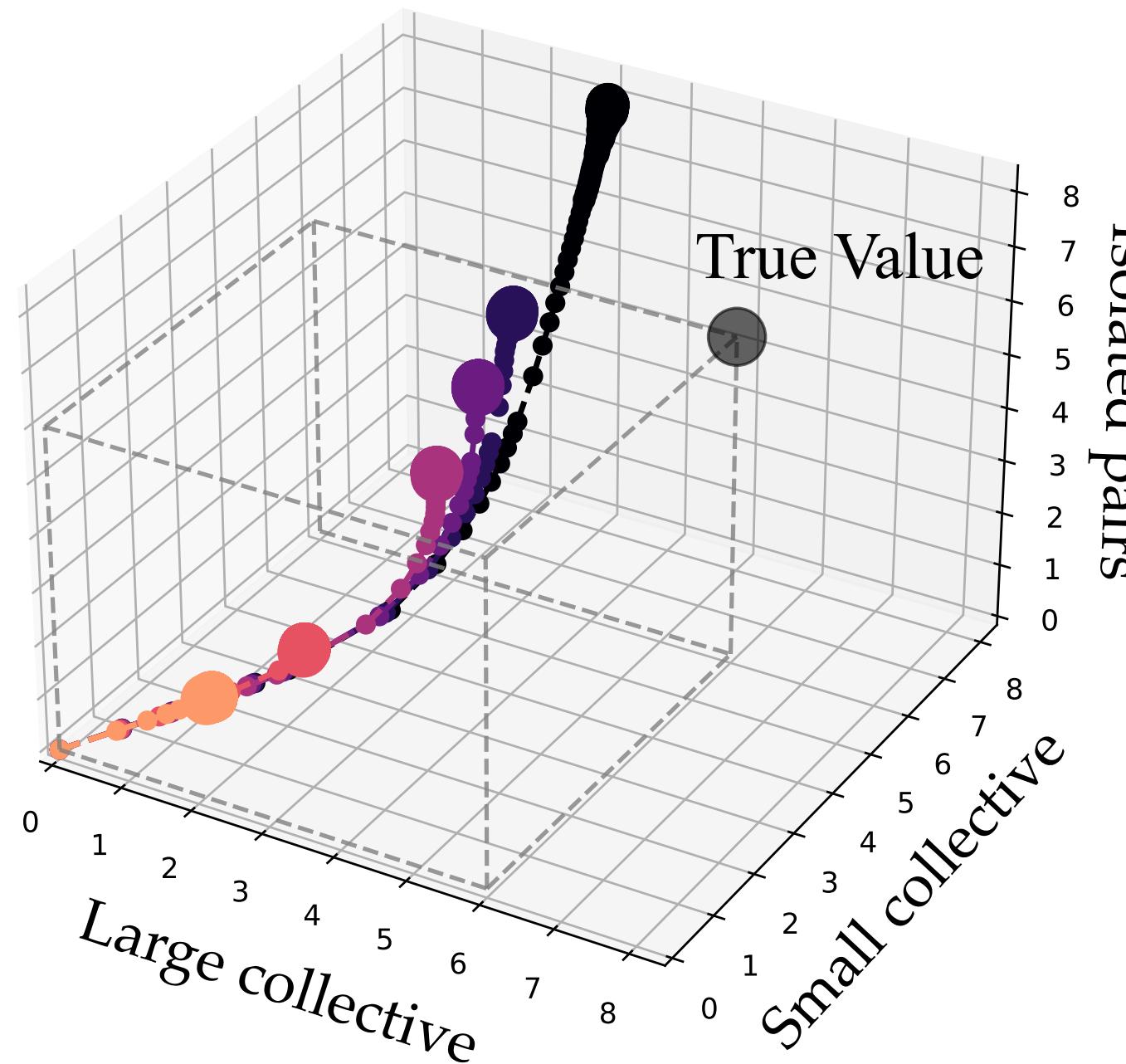
### III. Stochastic Boltzmann Machine

*Assess model performance with a toy model*

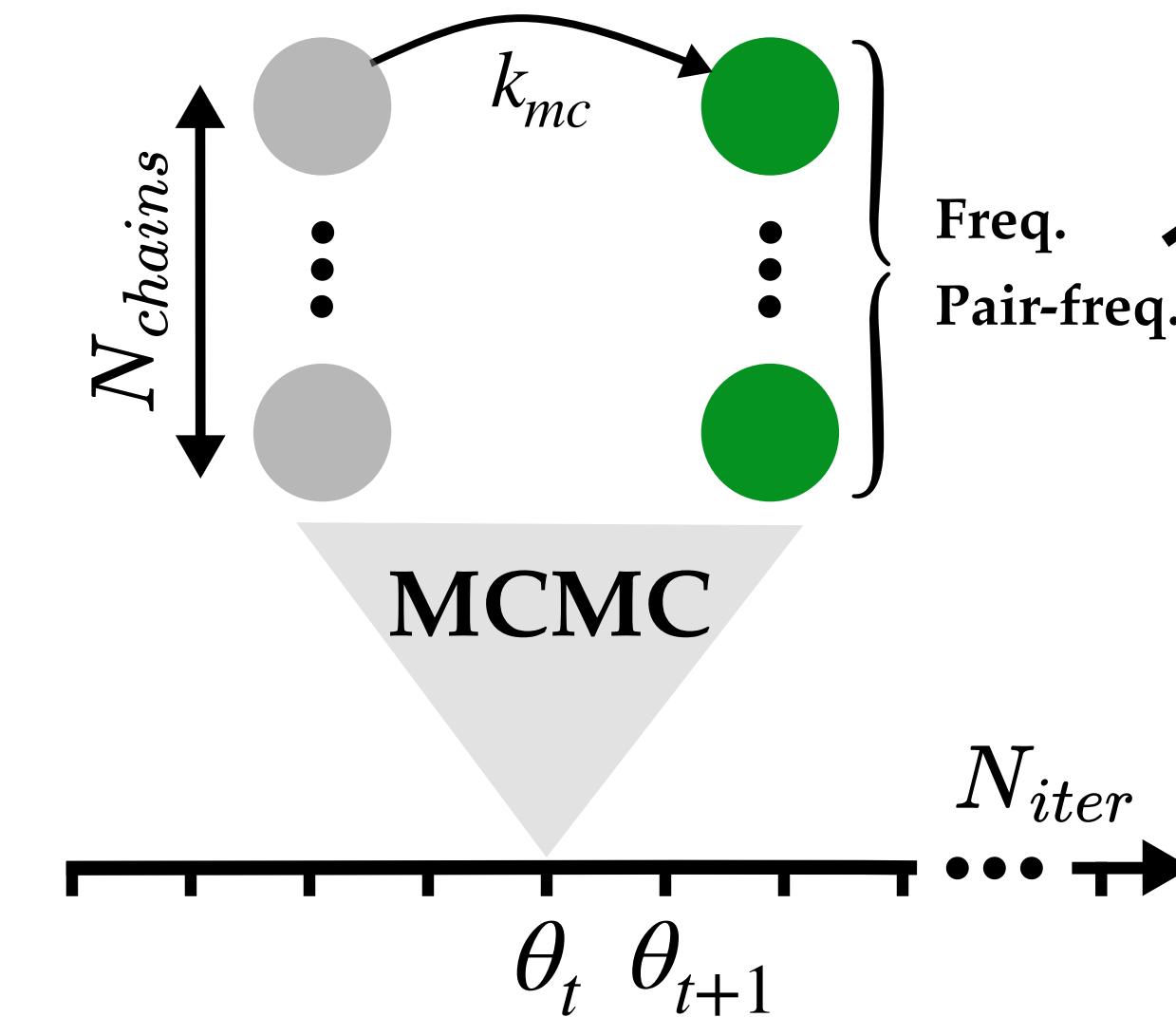
**BM**

$L_2$  regularization

1<sup>st</sup> order:  $\mathbf{p}_t = -\nabla f(\boldsymbol{\theta}_t)$



$$\boldsymbol{\theta} = \{\mathbf{J}, \mathbf{h}\} \quad f(\boldsymbol{\theta}) = \frac{1}{M} \sum_{m=1}^M \log P(\sigma^{(m)} | \boldsymbol{\theta}) - \lambda_J \|\mathbf{J}\|^2 - \lambda_h \|\mathbf{h}\|^2$$



$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \mathbf{p}_t$$

$$1^{st} \text{ order: } \mathbf{p}_t = -\nabla f(\boldsymbol{\theta}_t)$$

$$\frac{\partial f(\boldsymbol{\theta})}{\partial J_{ij}(a, b)} = f_{ij}(a, b) - \langle \delta(\sigma_i, a) \delta(\sigma_j, b) \rangle + \lambda_J J_{ij}(a, b)$$

Empirical                      Model

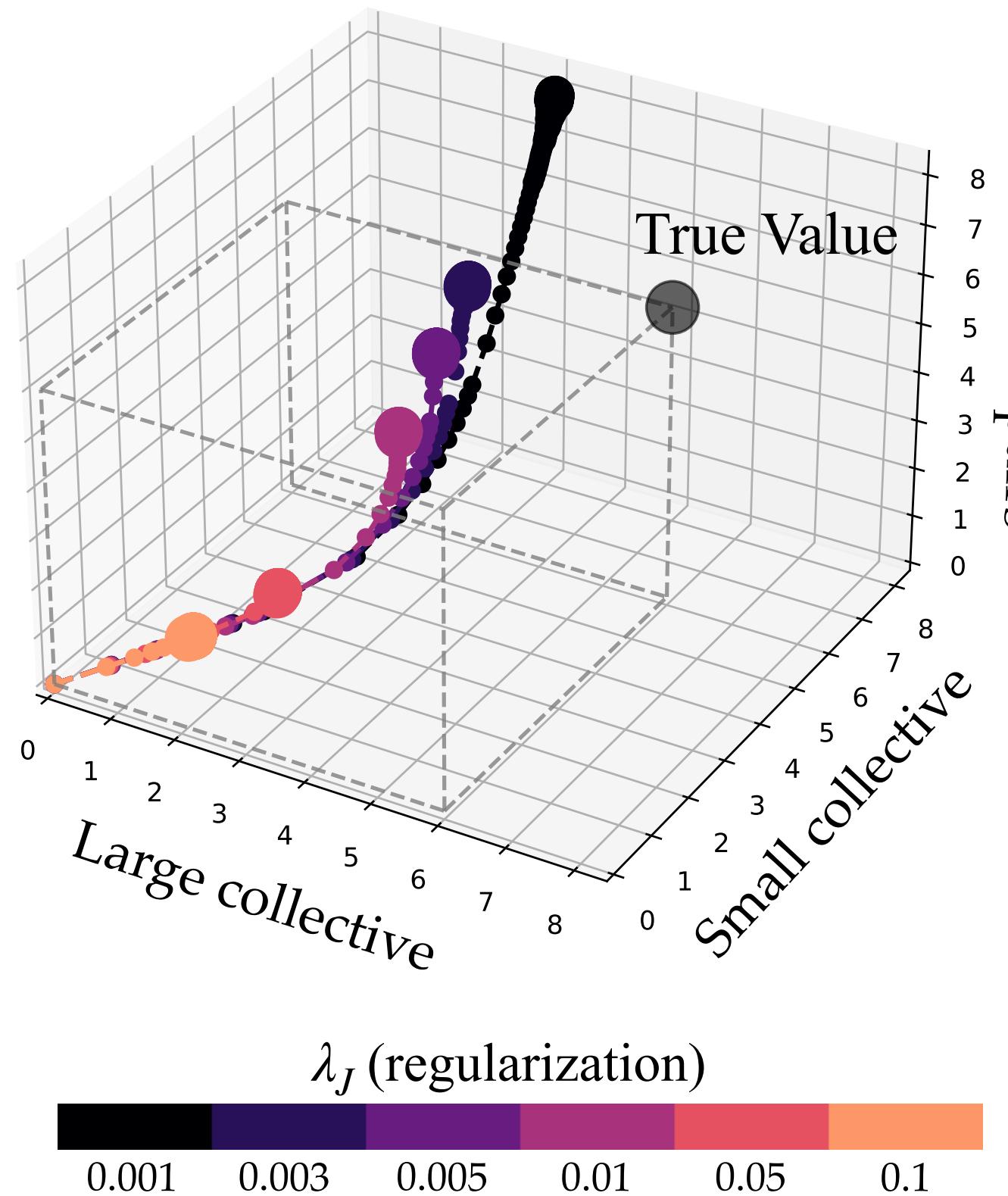
### III. Stochastic Boltzmann Machine

**BM**

$L_2$  regularization  
 $1^{st}$  order:  $\mathbf{p}_t = -\nabla f(\theta_t)$

**BM**

$L_2$  regularization  
 $2^{nd}$  order:  $\mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\theta_t)$



#### L-BFGS algorithm

**Why?**

- High dimensional space  $\sim 10^4$
- $1^{st}$  order suboptimal because of anisotropic curvature

**How?**

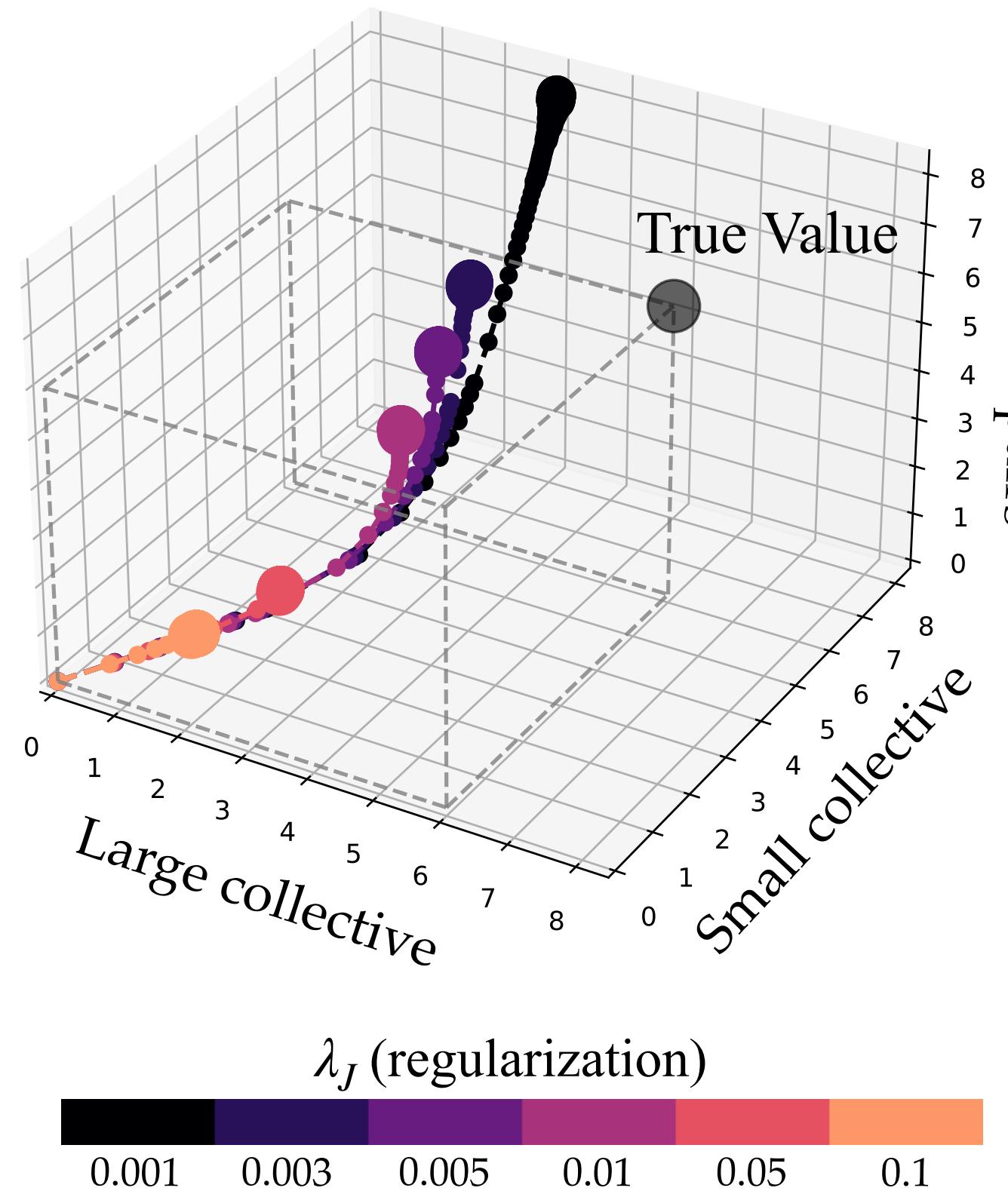
- Curvature information from the  $m$  most recent iterations

$$\mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\theta_t)$$

### III. Stochastic Boltzmann Machine

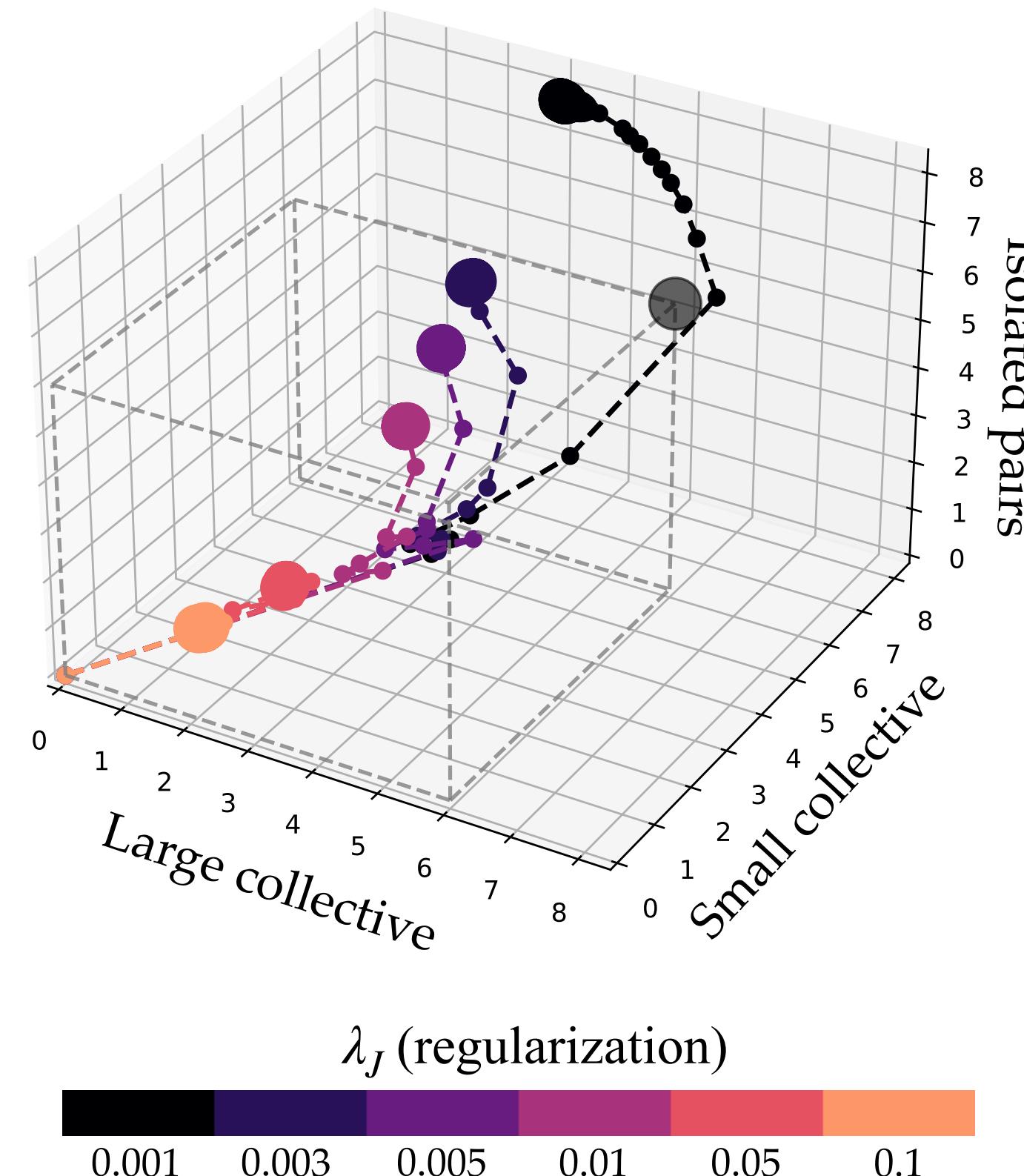
**BM**

$L_2$  regularization  
 $1^{st}$  order:  $\mathbf{p}_t = -\nabla f(\theta_t)$



**BM**

$L_2$  regularization  
 $2^{nd}$  order:  $\mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\theta_t)$



#### L-BFGS algorithm

**Why?**

- High dimensional space  $\sim 10^4$
- $1^{st}$  order suboptimal because of anisotropic curvature

**How?**

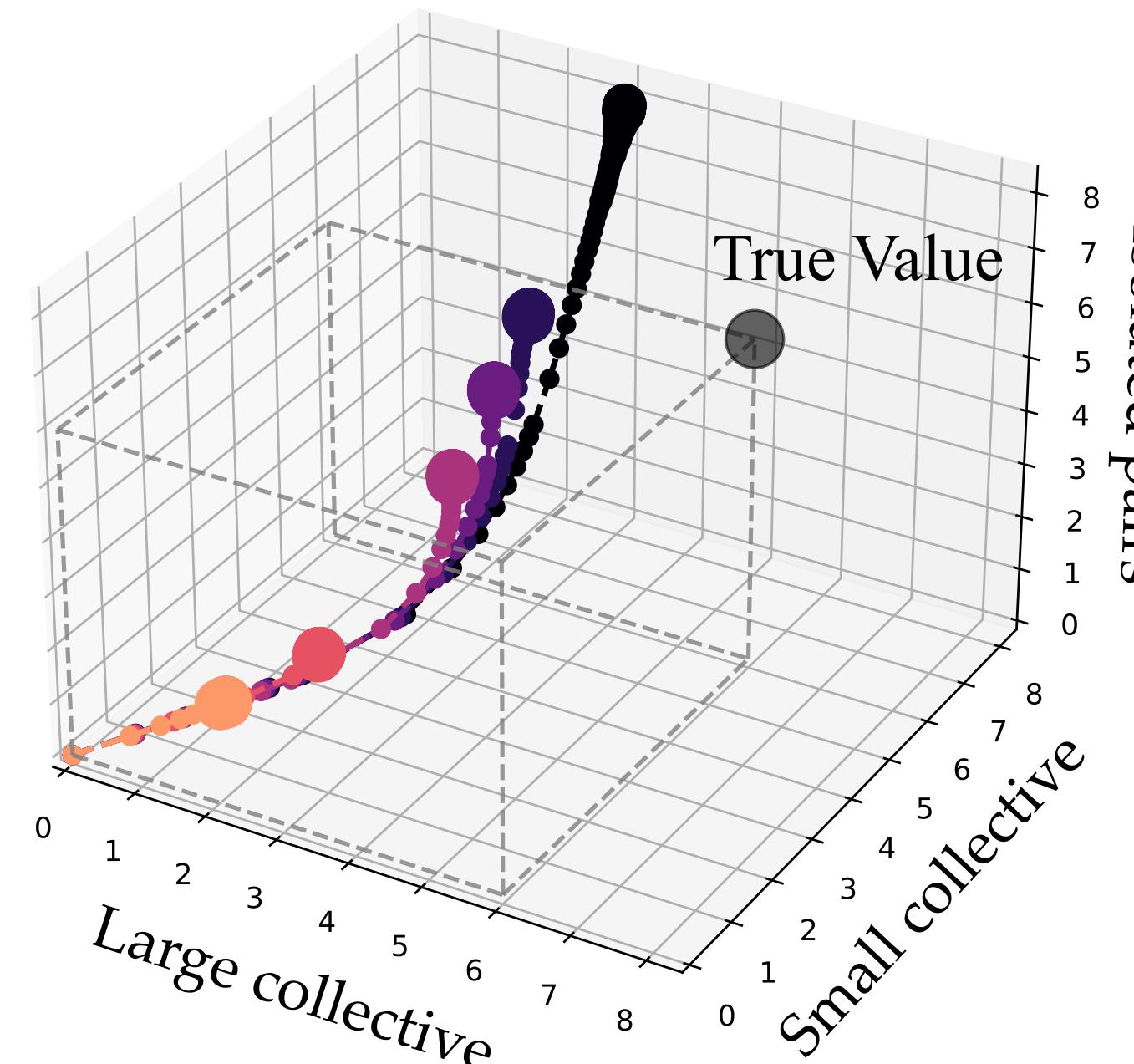
- Curvature information from the  $m$  most recent iterations

$$\mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\theta_t)$$

### III. Stochastic Boltzmann Machine

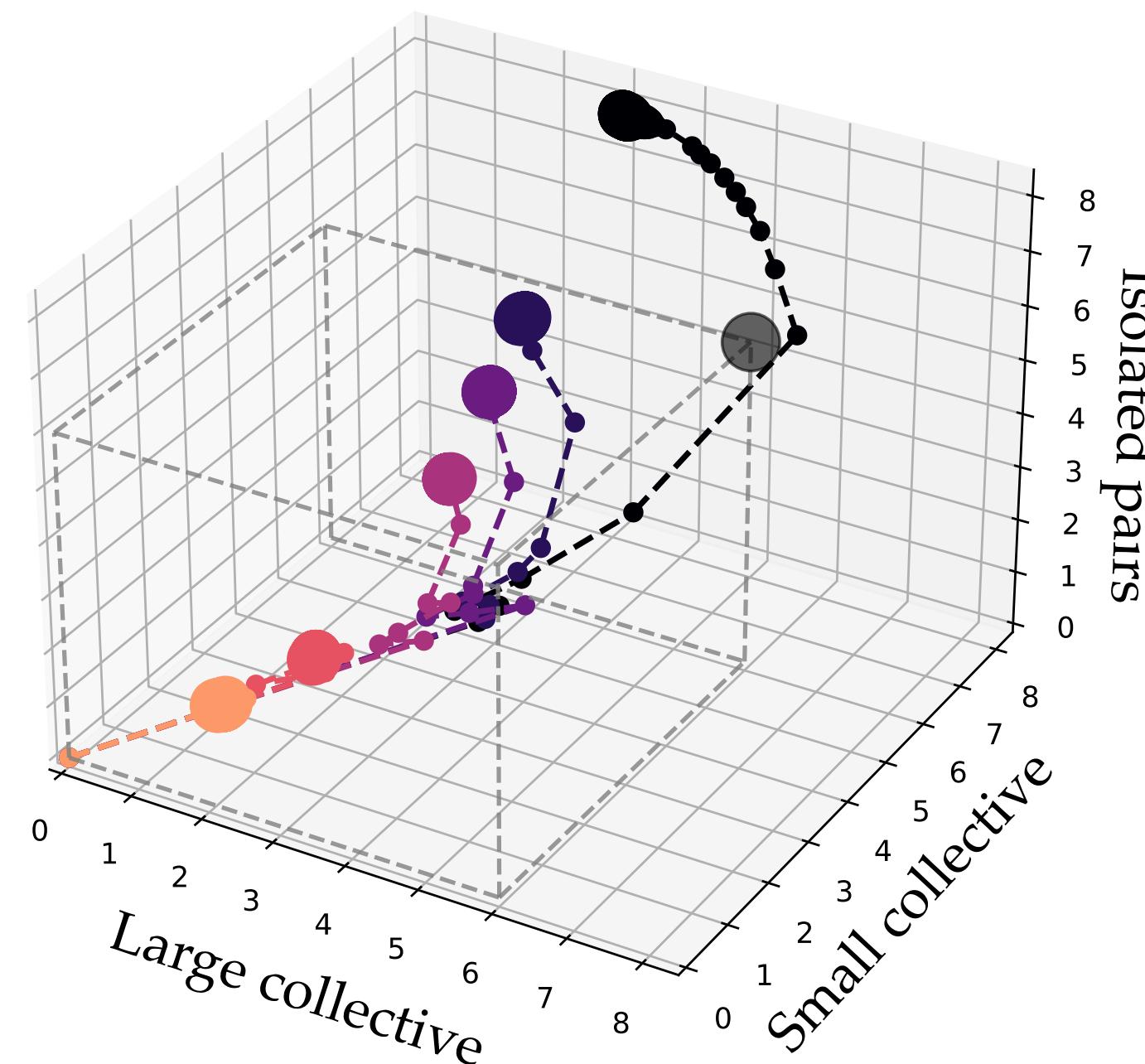
**BM**

$L_2$  regularization  
 $1^{st}$  order:  $\mathbf{p}_t = -\nabla f(\theta_t)$



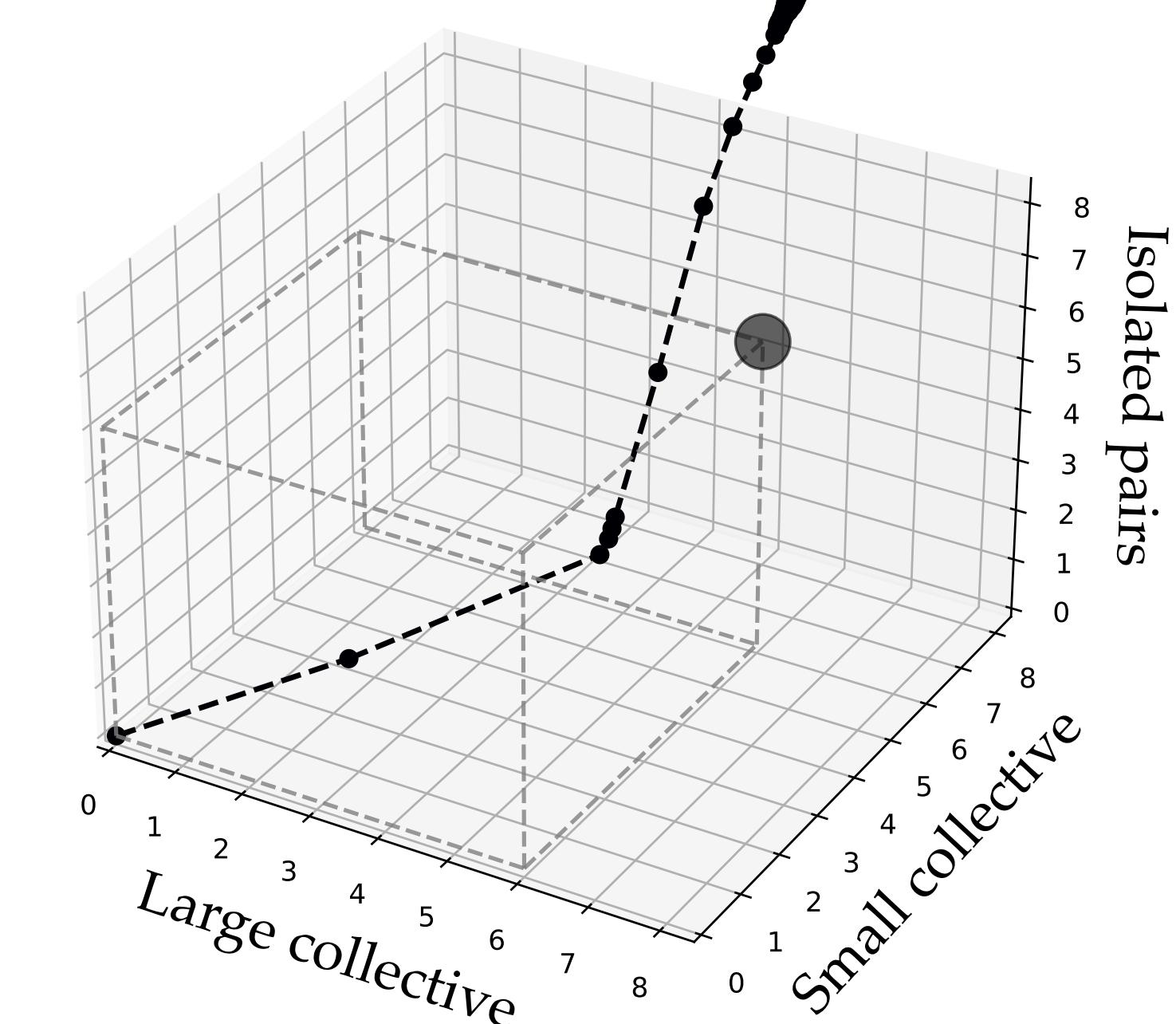
**BM**

$L_2$  regularization  
 $2^{nd}$  order:  $\mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\theta_t)$



**Stochastic Boltzmann Machine**

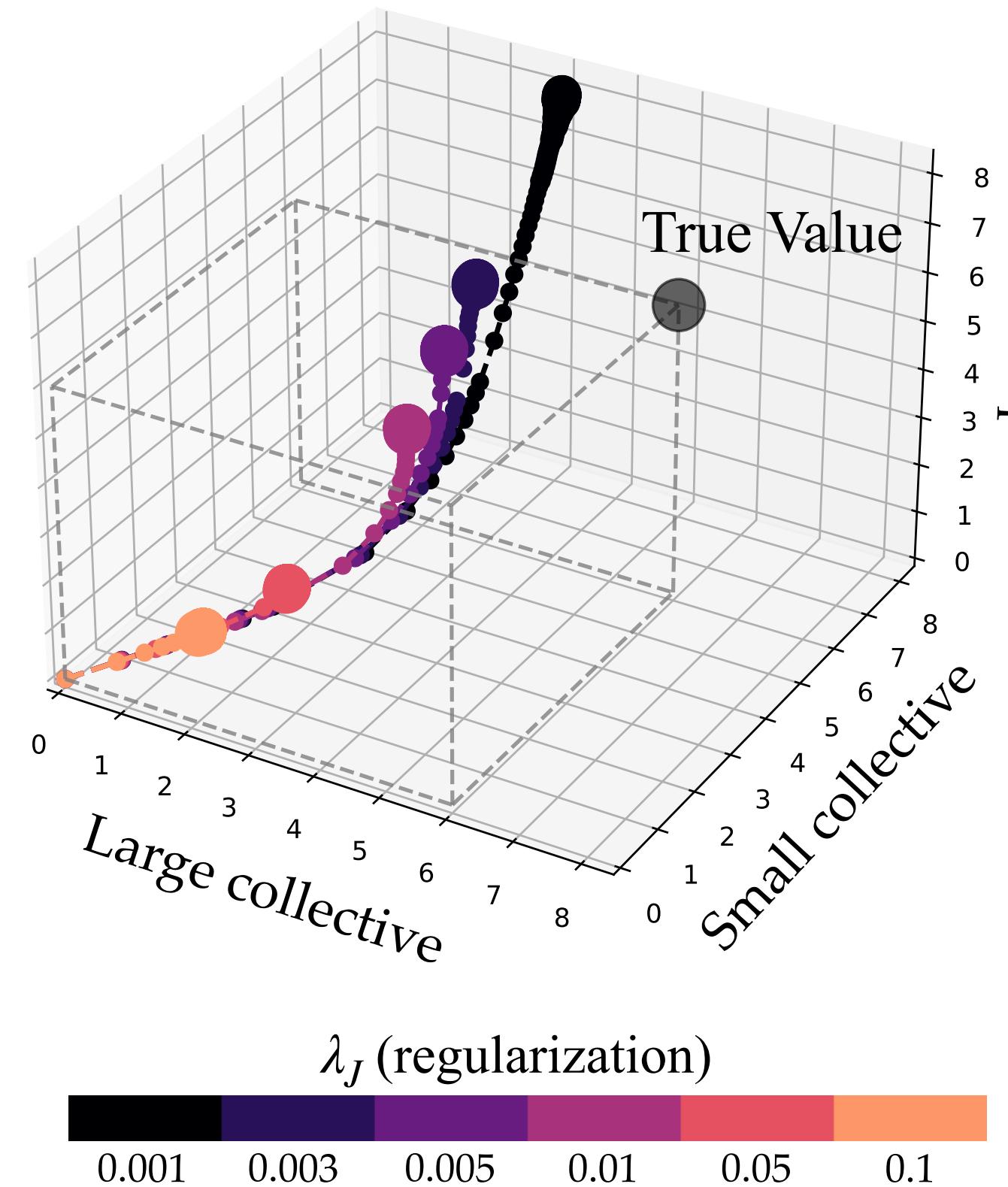
$L_2$  regularization  
 $2^{nd}$  order:  $\mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\theta_t)$



### III. Stochastic Boltzmann Machine

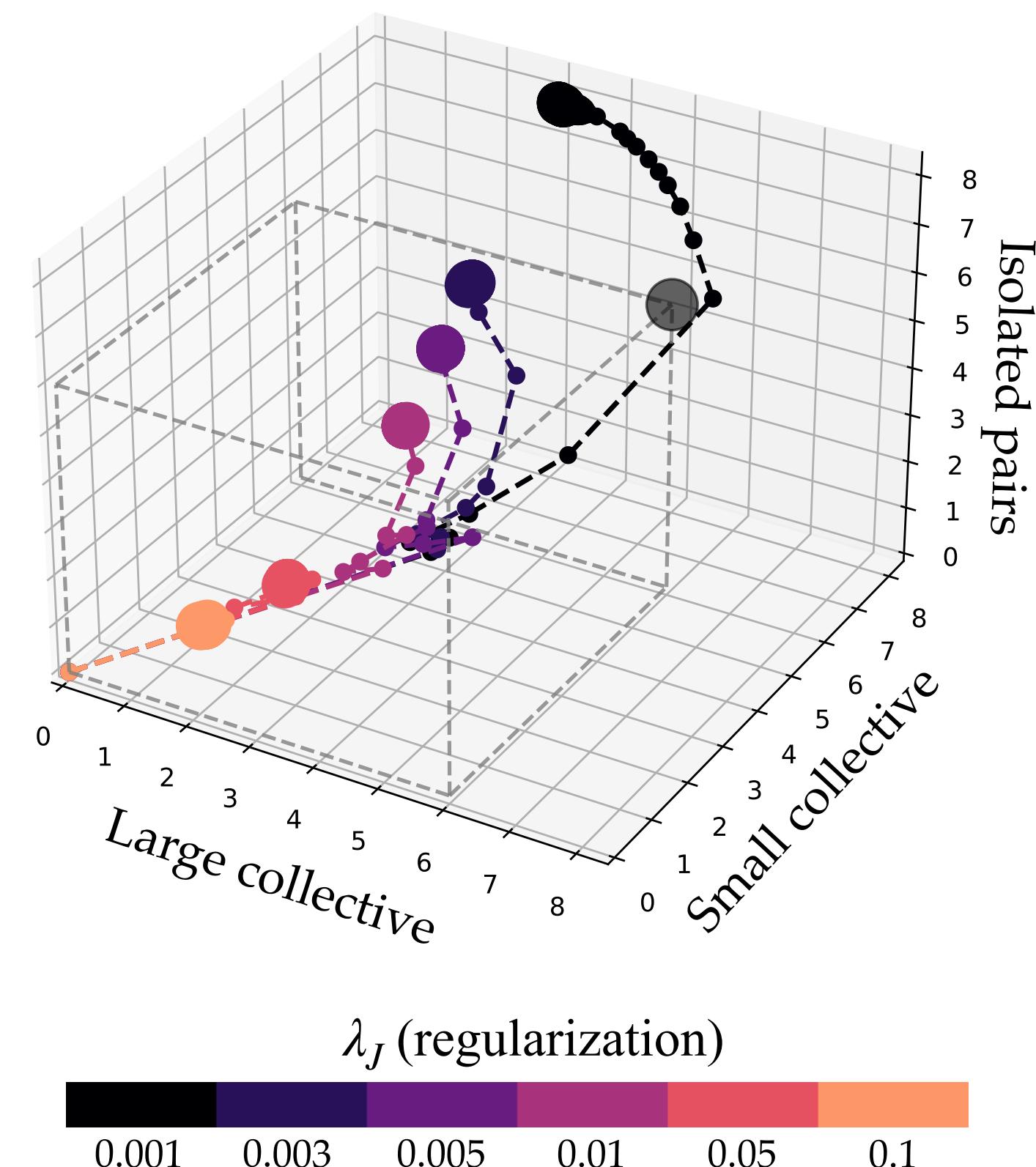
**BM**

$L_2$  regularization  
 $1^{st}$  order:  $\mathbf{p}_t = -\nabla f(\theta_t)$



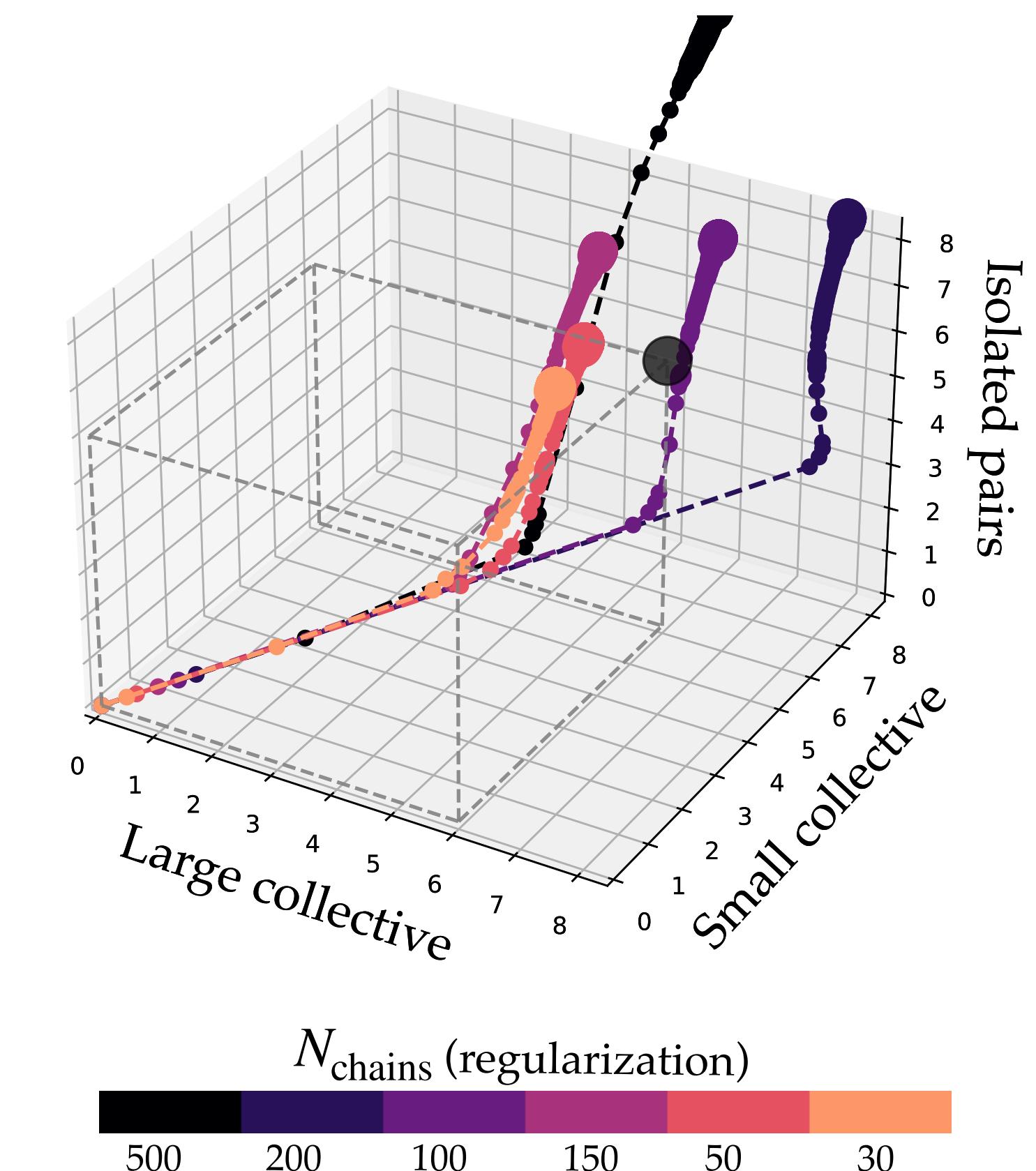
**BM**

$L_2$  regularization  
 $2^{nd}$  order:  $\mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\theta_t)$



**Stochastic Boltzmann Machine**

$L_2$  regularization  
 $2^{nd}$  order:  $\mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\theta_t)$

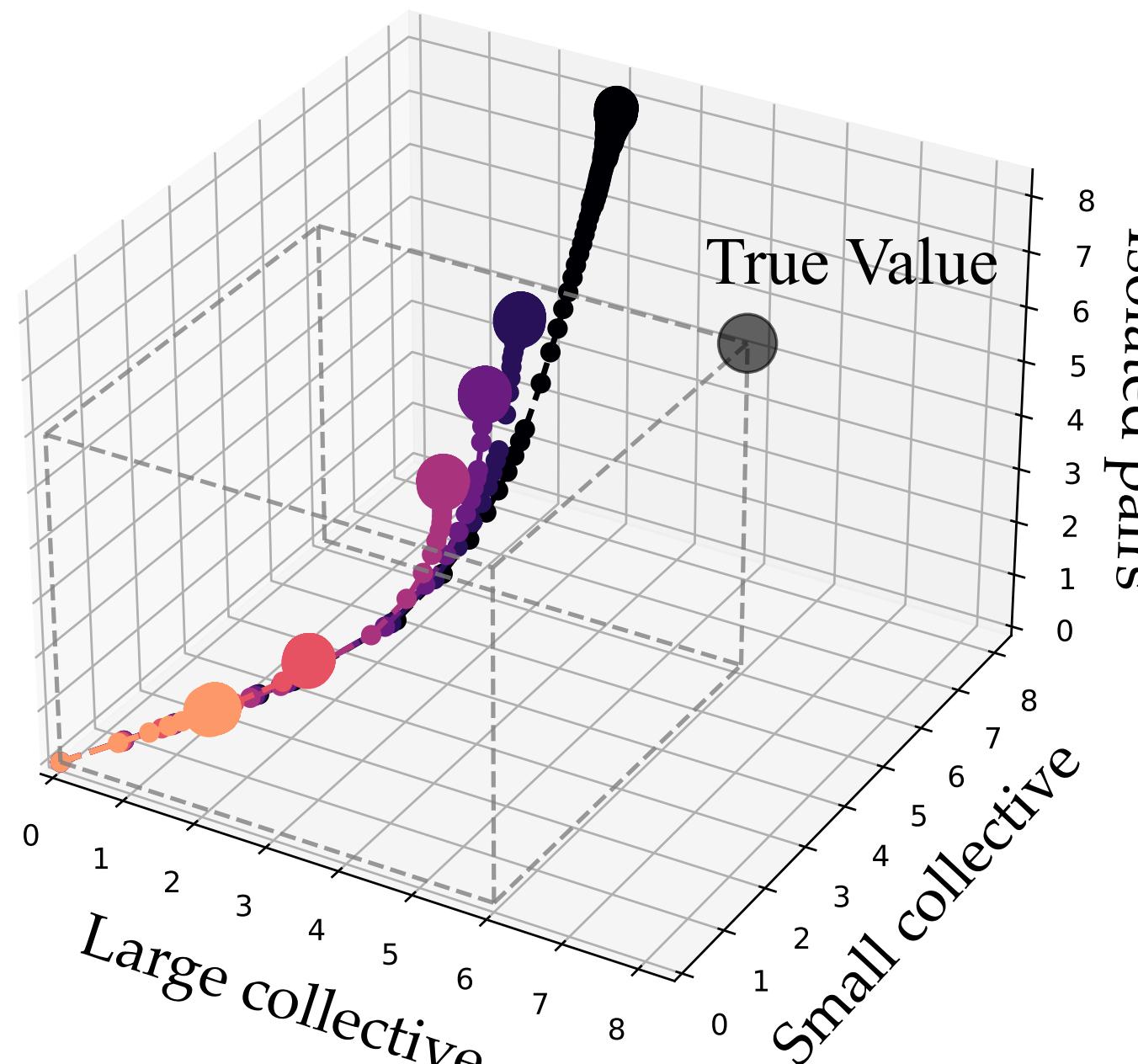


### III. Stochastic Boltzmann Machine

**BM**

$L_2$  regularization

$$1^{st} \text{ order: } \mathbf{p}_t = -\nabla f(\boldsymbol{\theta}_t)$$



$\lambda_J$  (regularization)



All models are inferred with  $k_{mc} = 10^5$ ,  $N_{iter} = 5000$

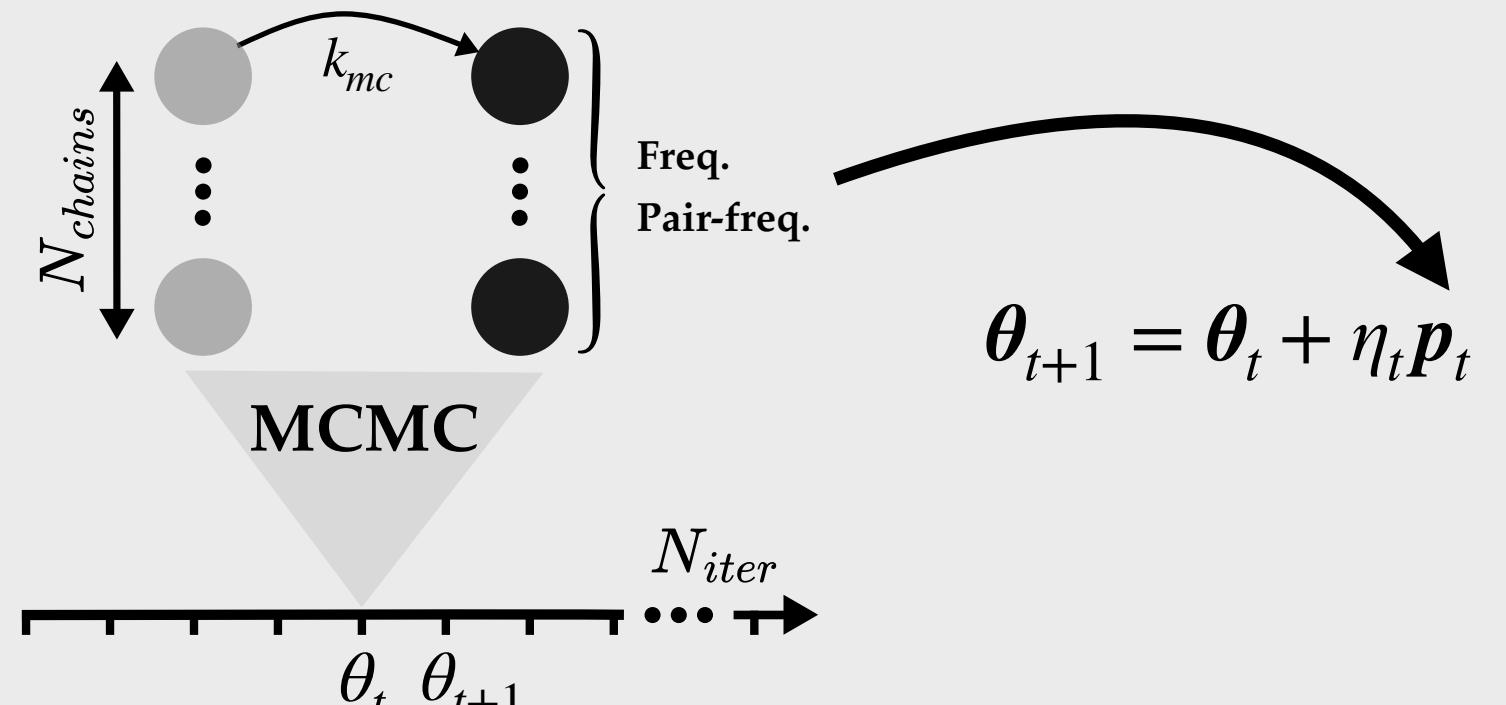
**Lowering  $N_{chains}$**

#### Observation

- ▶ Training BM: capture statistics
- ▶ But statistics lack reliability because of undersampling

#### Strategy

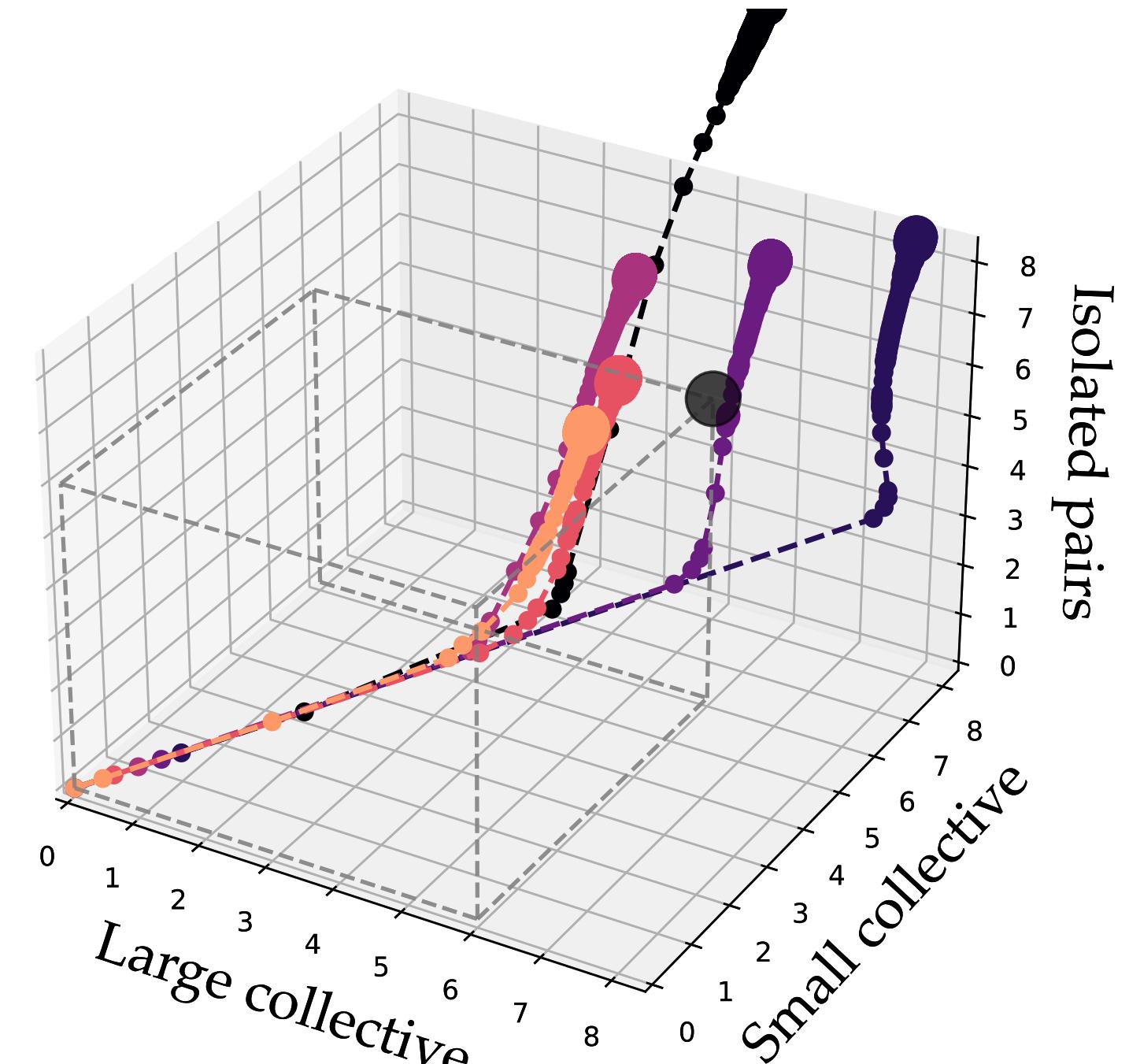
- ▶ Intentionally undersample the model to mirror data undersampling



**Stochastic Boltzmann Machine**

$L_2$  regularization

$$2^{nd} \text{ order: } \mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\boldsymbol{\theta}_t)$$



$N_{chains}$  (regularization)

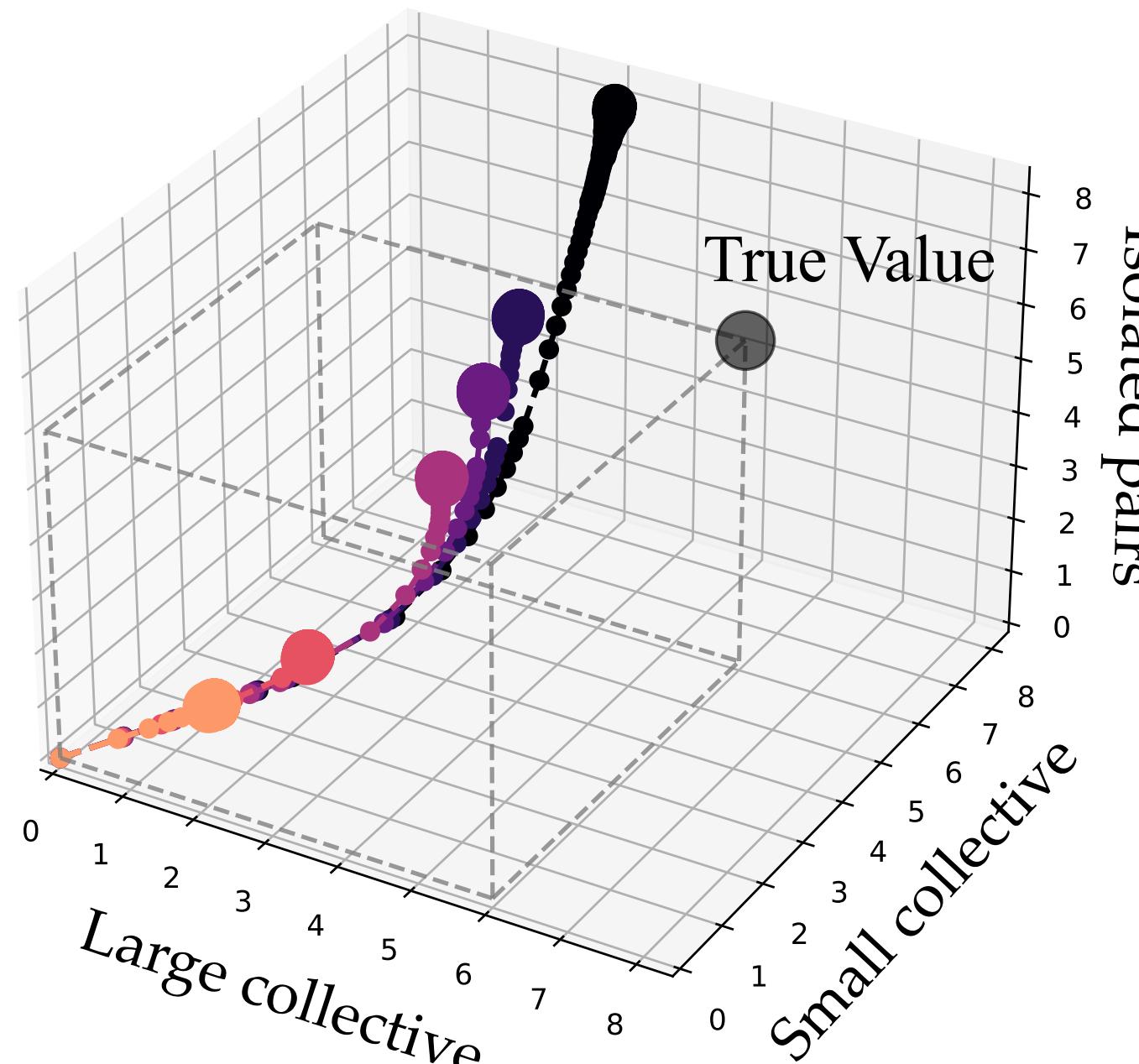


### III. Stochastic Boltzmann Machine

**BM**

$L_2$  regularization

$$1^{st} \text{ order: } \mathbf{p}_t = -\nabla f(\boldsymbol{\theta}_t)$$



$\lambda_J$  (regularization)



0.001 0.003 0.005 0.01 0.05 0.1

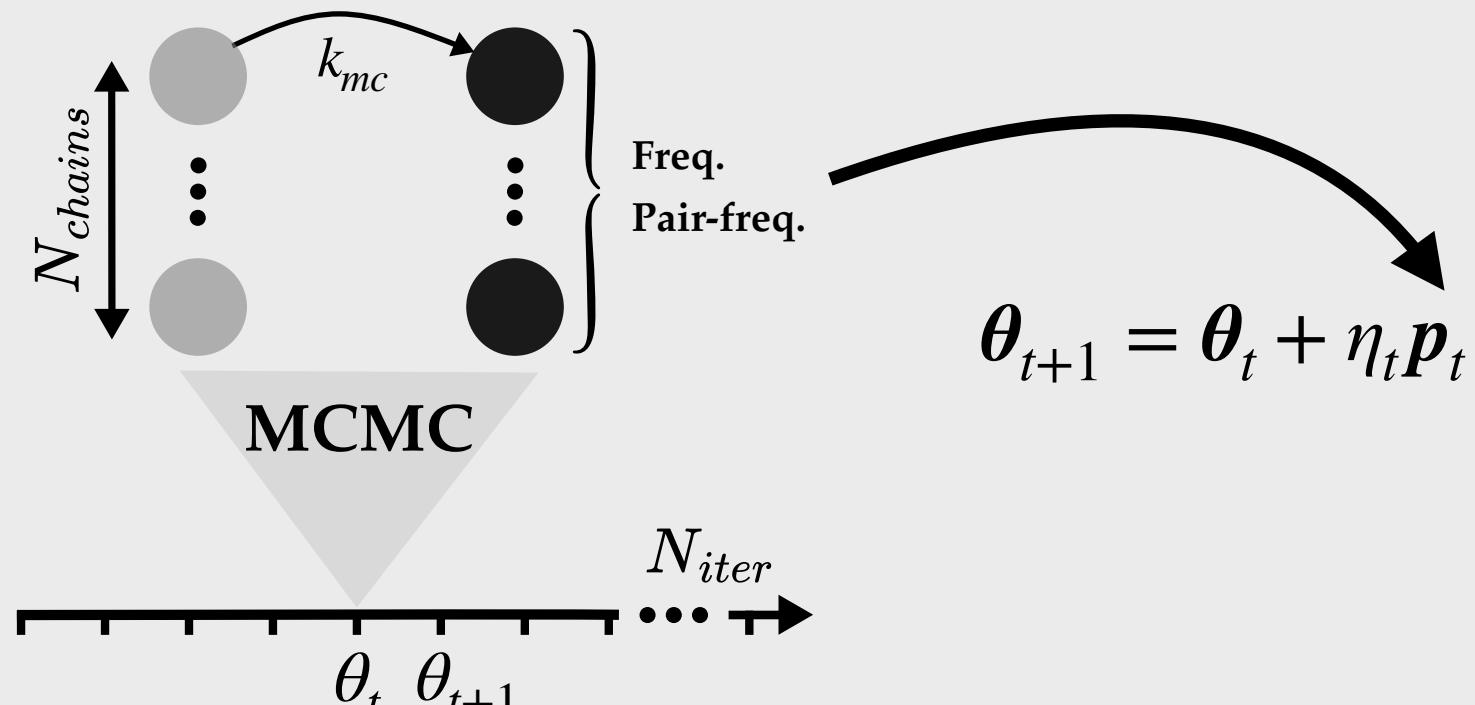
**Lowering  $N_{\text{chains}}$**

#### Observation

- ▶ Training BM: capture statistics
- ▶ But statistics lack reliability because of undersampling

#### Strategy

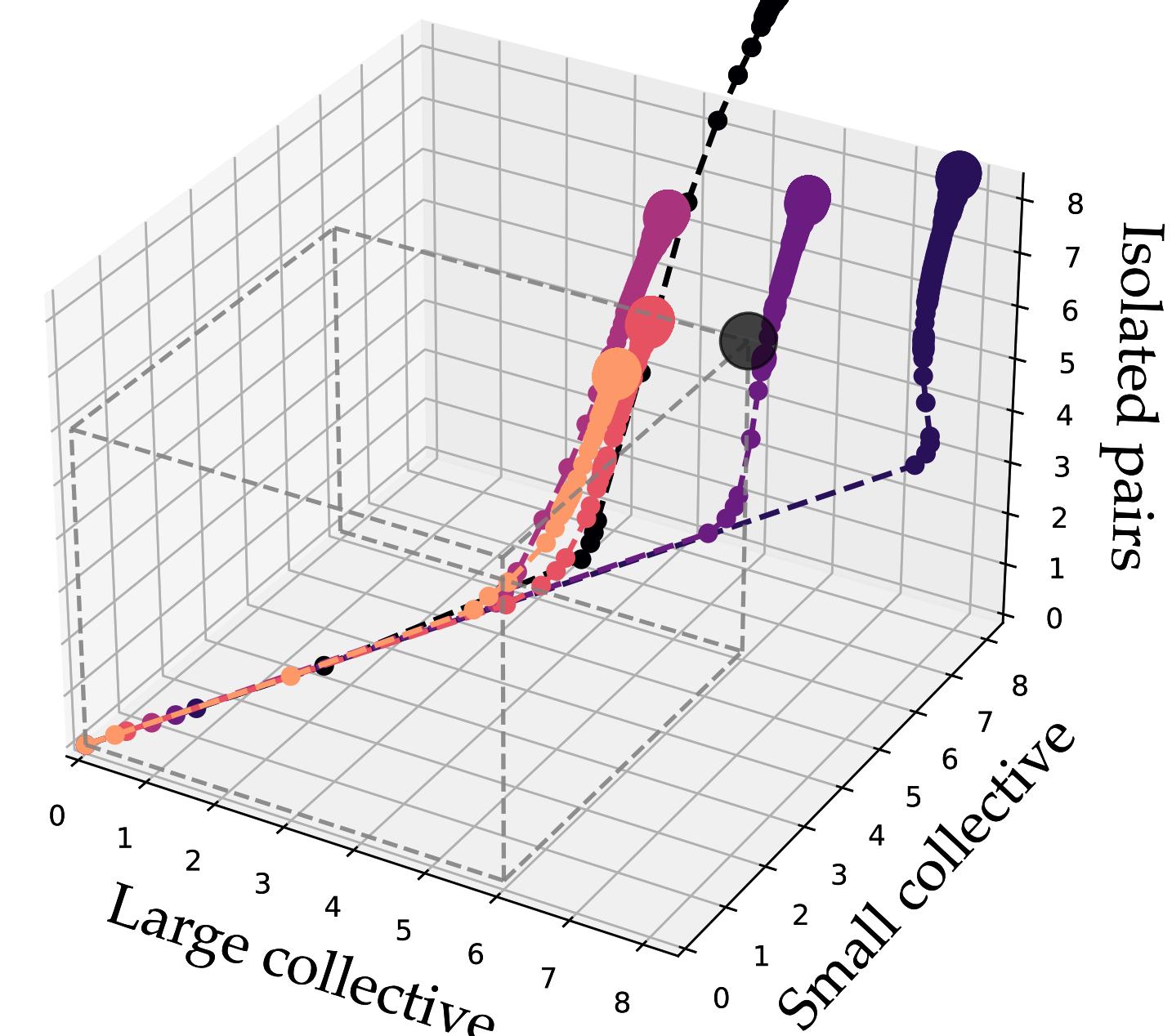
- ▶ Intentionally undersample the model to mirror data undersampling



**Stochastic Boltzmann Machine**

$L_2$  regularization

$$2^{nd} \text{ order: } \mathbf{p}_t = -\mathbf{B}_k^{-1} \nabla f(\boldsymbol{\theta}_t)$$



$N_{\text{chains}}$  (regularization)



500 200 100 150 50 30

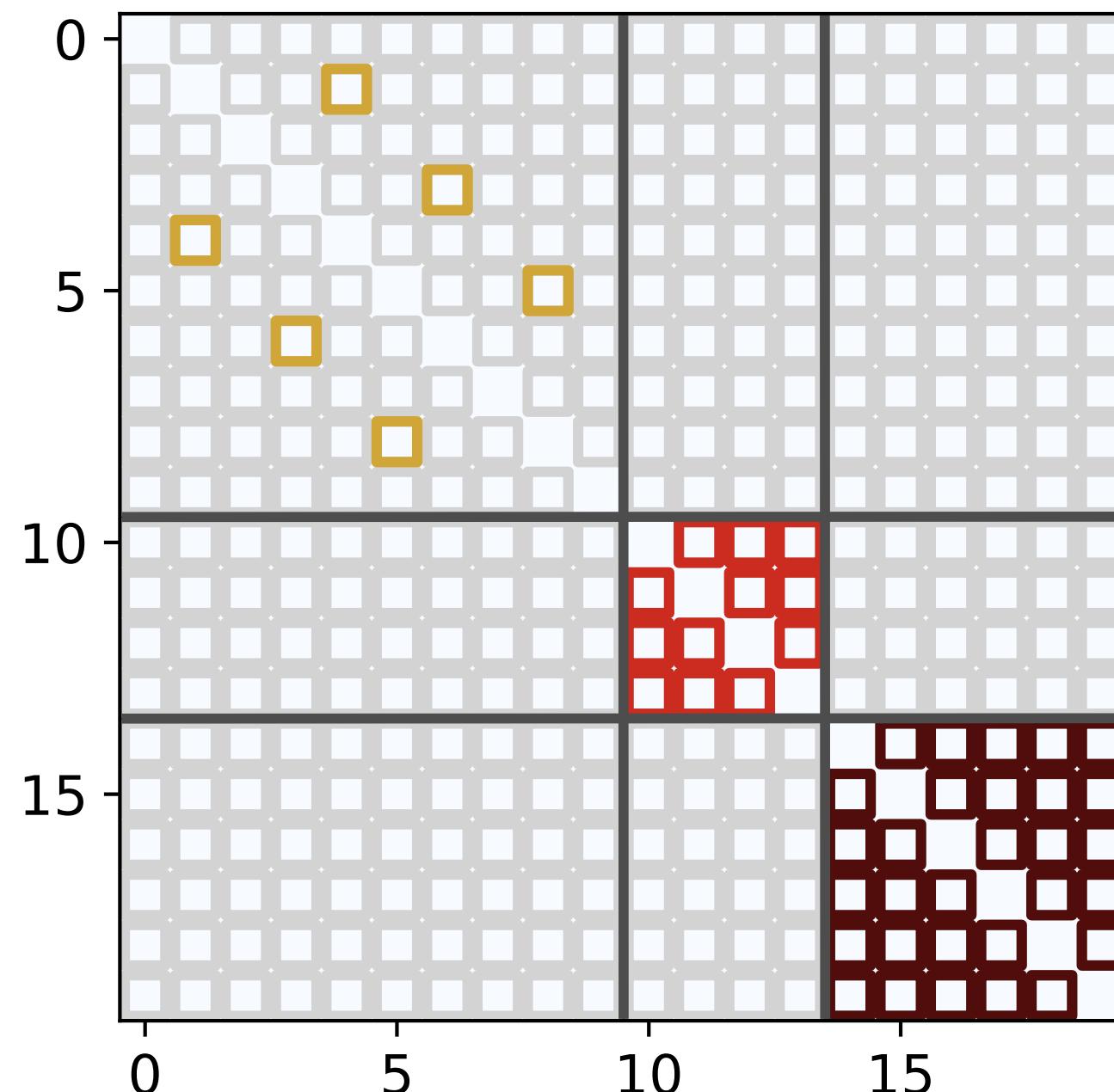
- ▶ Model undersampling:  $N_{\text{chains}}$
- ▶ Slow down convergence:  $m \sim 1$
- ▶ Early stopping:  $N_{\text{iter}}$

All models are inferred with  $k_{mc} = 10^5$ ,  $N_{\text{iter}} = 5000$

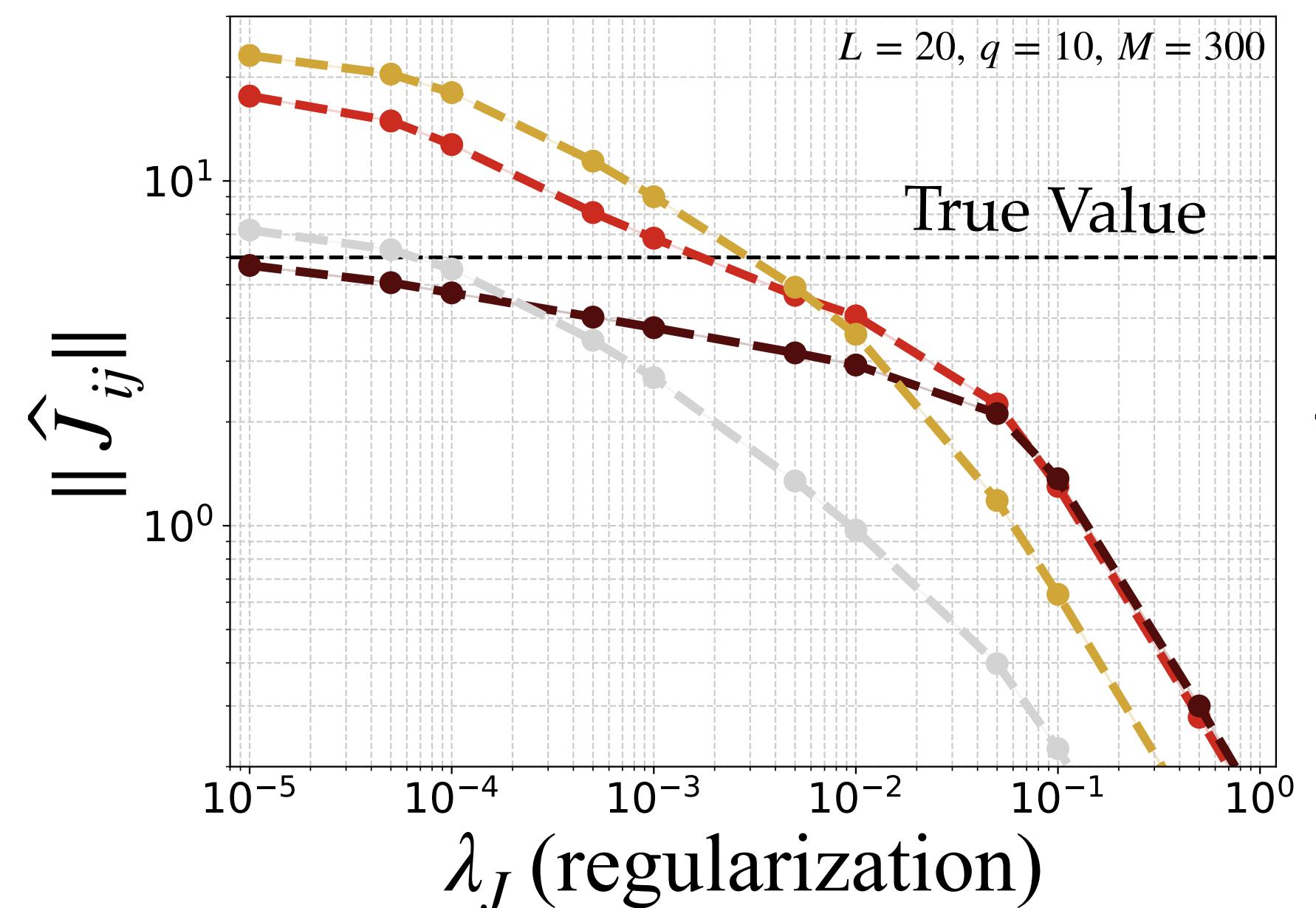
### III. Stochastic Boltzmann Machine

*Assess model performance with a toy model*

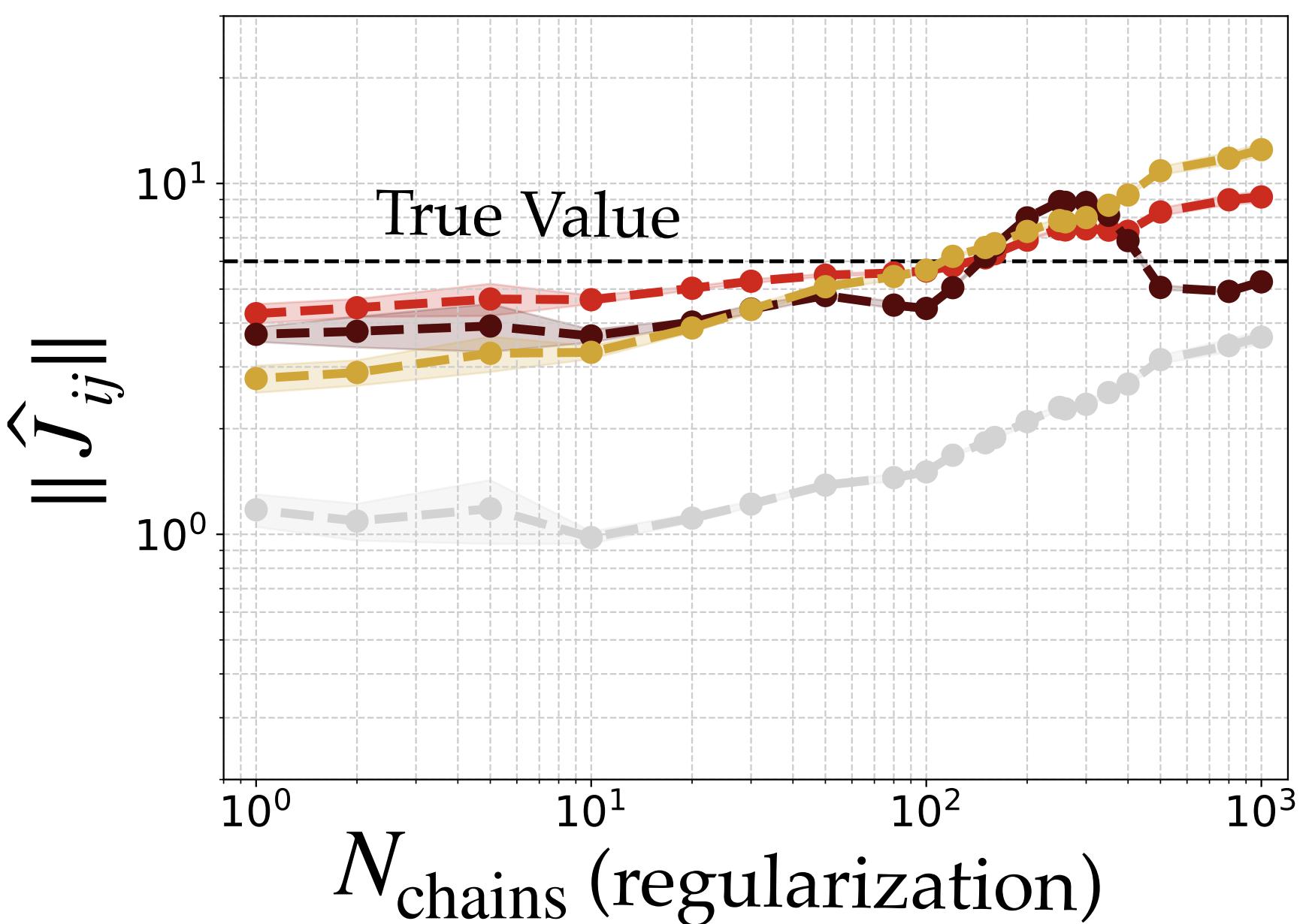
- Isolated pairs
- Small collective
- Large collective
- Non interacting



BM



SBM



#### Stochastic Boltzmann Machine

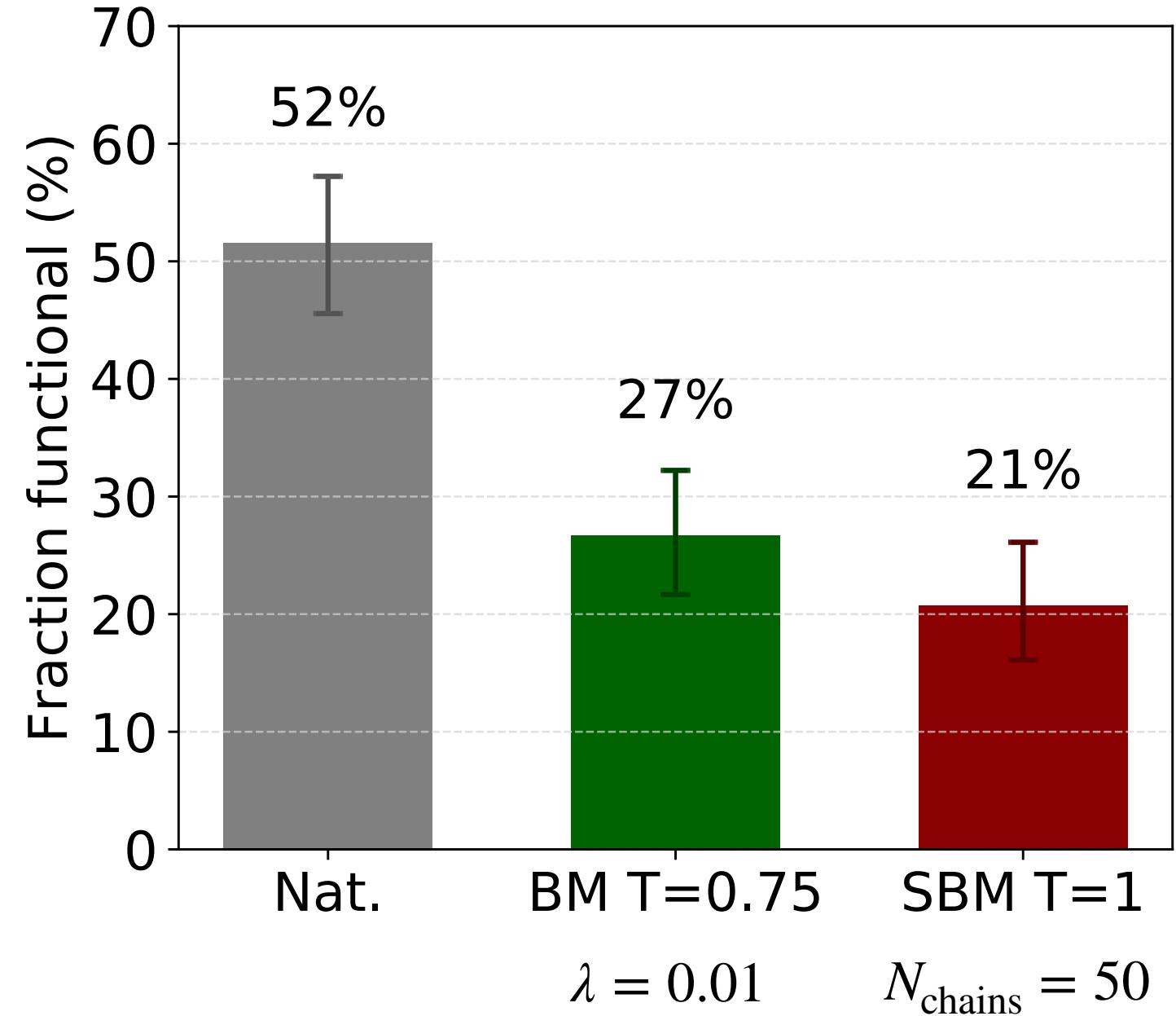
- ▶ Toy model: correct undersampling-induced biases
- ▶ Real data?

### III. Stochastic Boltzmann Machine

*Relevance to real proteins: application to the chorismate mutase family*

#### Fidelity

*Fraction of functional sequences*



#### Novelty

*Distribution of identity to the nearest natural sequence*

#### Diversity

*% of taxonomic families represented*

The experiments were performed by Emily Hinds

The models are inferred with  $k_{\text{mc}} = 10^5$ ,  $N_{\text{iter}} = 400$ ,  $\theta = 0.3$

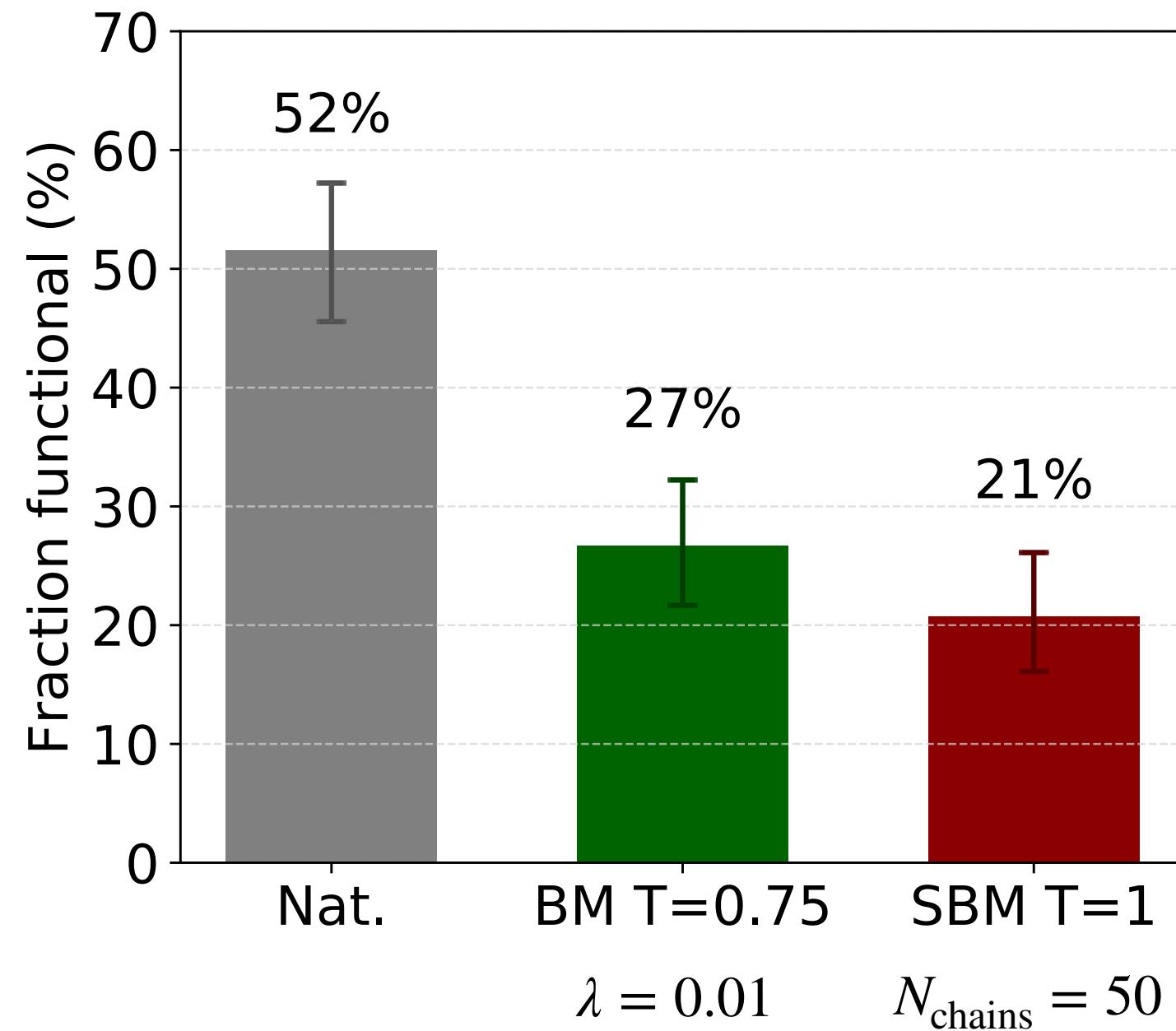
890, 193 and 180 sequences were tested for the natural dataset, the SBM and the BM models respectively

### III. Stochastic Boltzmann Machine

*Relevance to real proteins: application to the chorismate mutase family*

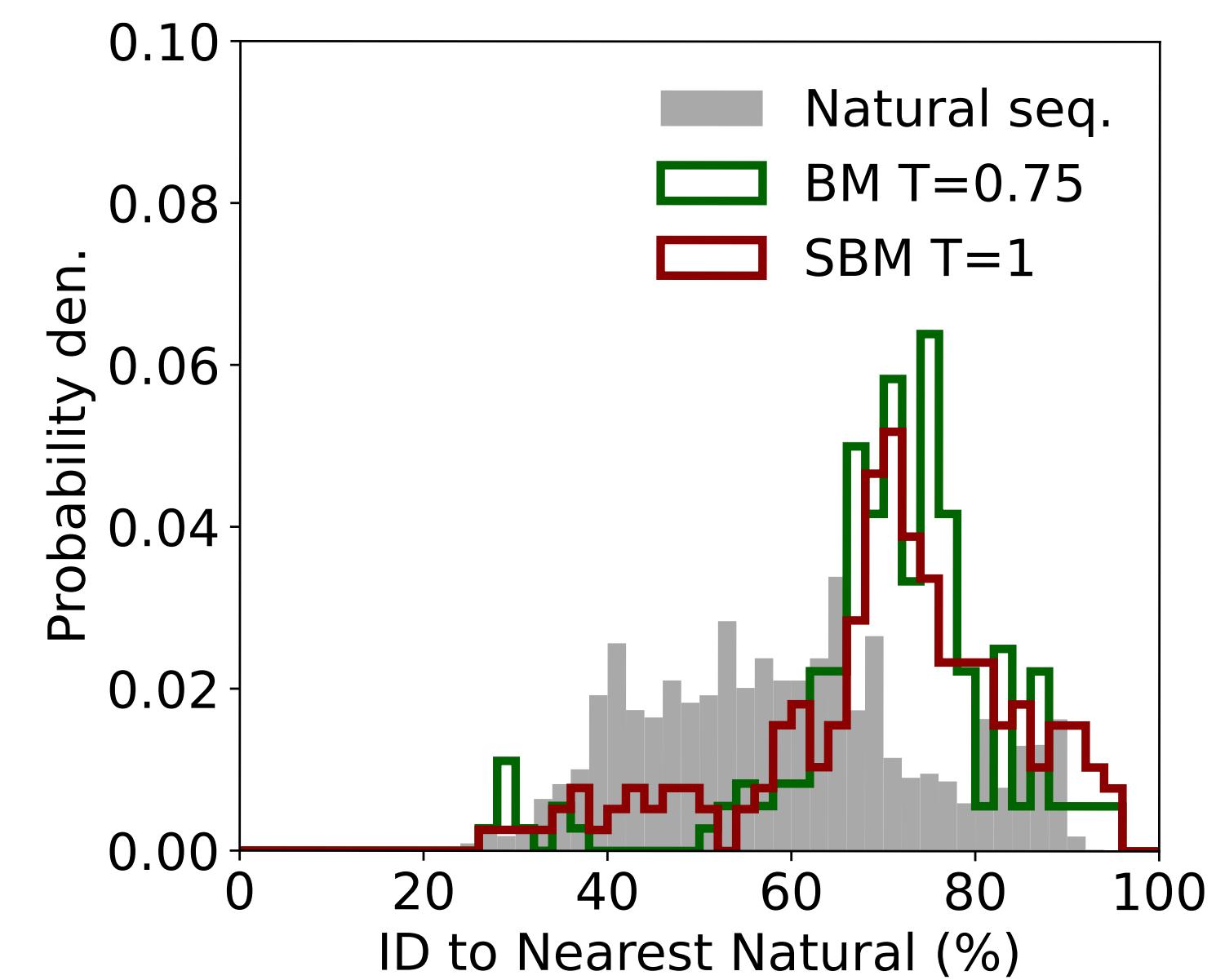
#### Fidelity

*Fraction of functional sequences*



#### Novelty

*Distribution of identity to the nearest natural sequence*



#### Diversity

*% of taxonomic families represented*

The experiments were performed by Emily Hinds

The models are inferred with  $k_{\text{mc}} = 10^5$ ,  $N_{\text{iter}} = 400$ ,  $\theta = 0.3$

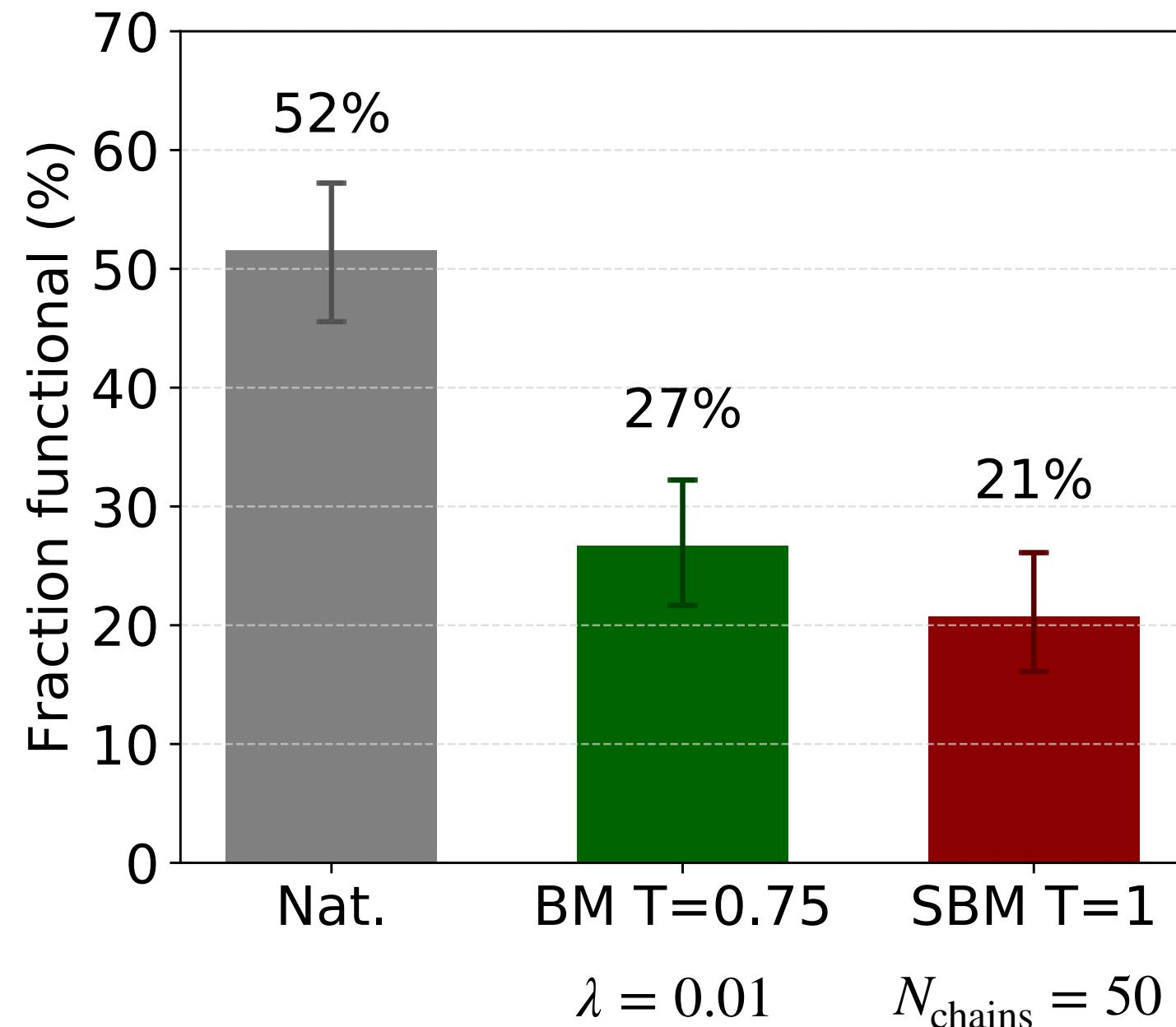
890, 193 and 180 sequences were tested for the natural dataset, the SBM and the BM models respectively

### III. Stochastic Boltzmann Machine

*Relevance to real proteins: application to the chorismate mutase family*

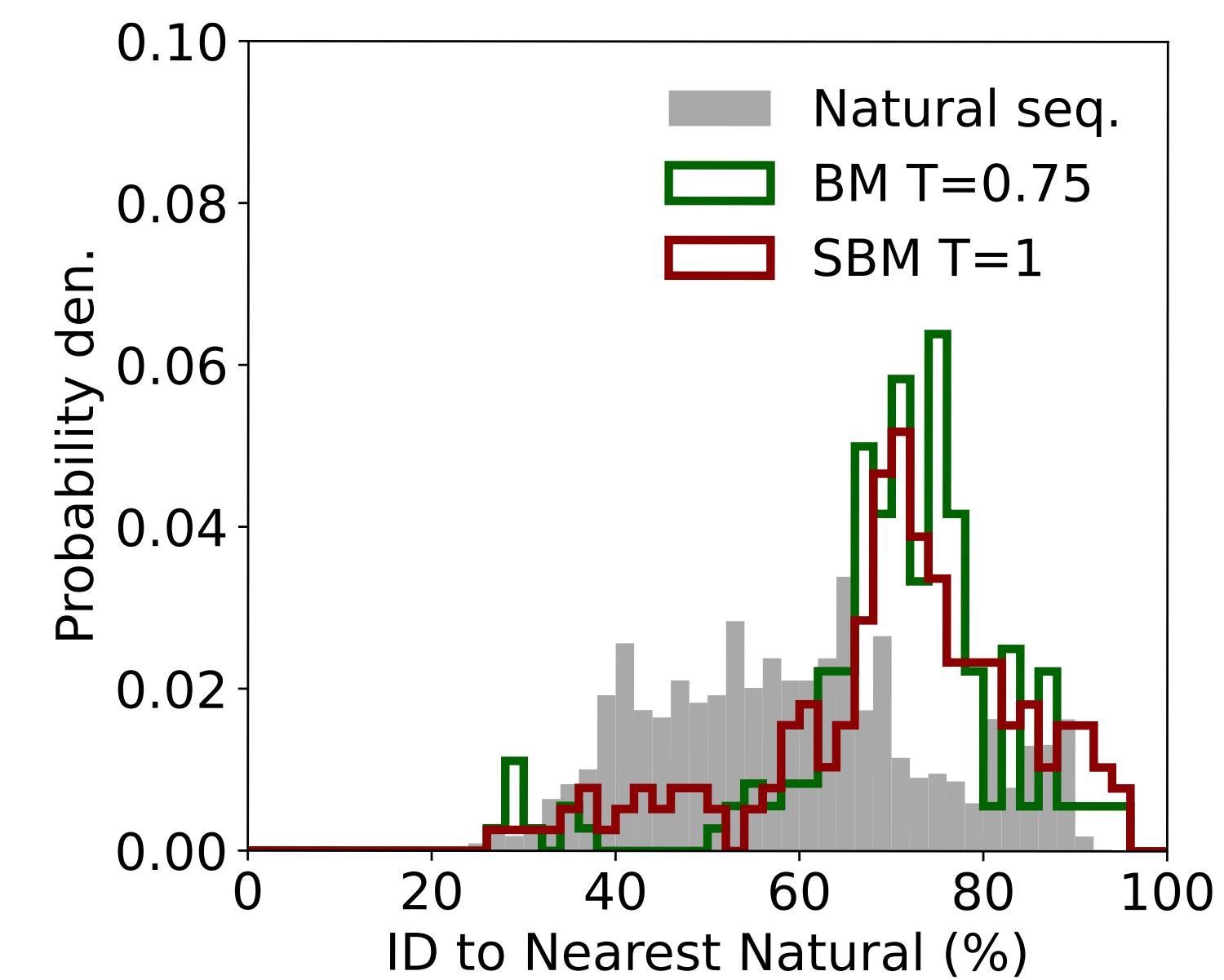
#### Fidelity

*Fraction of functional sequences*



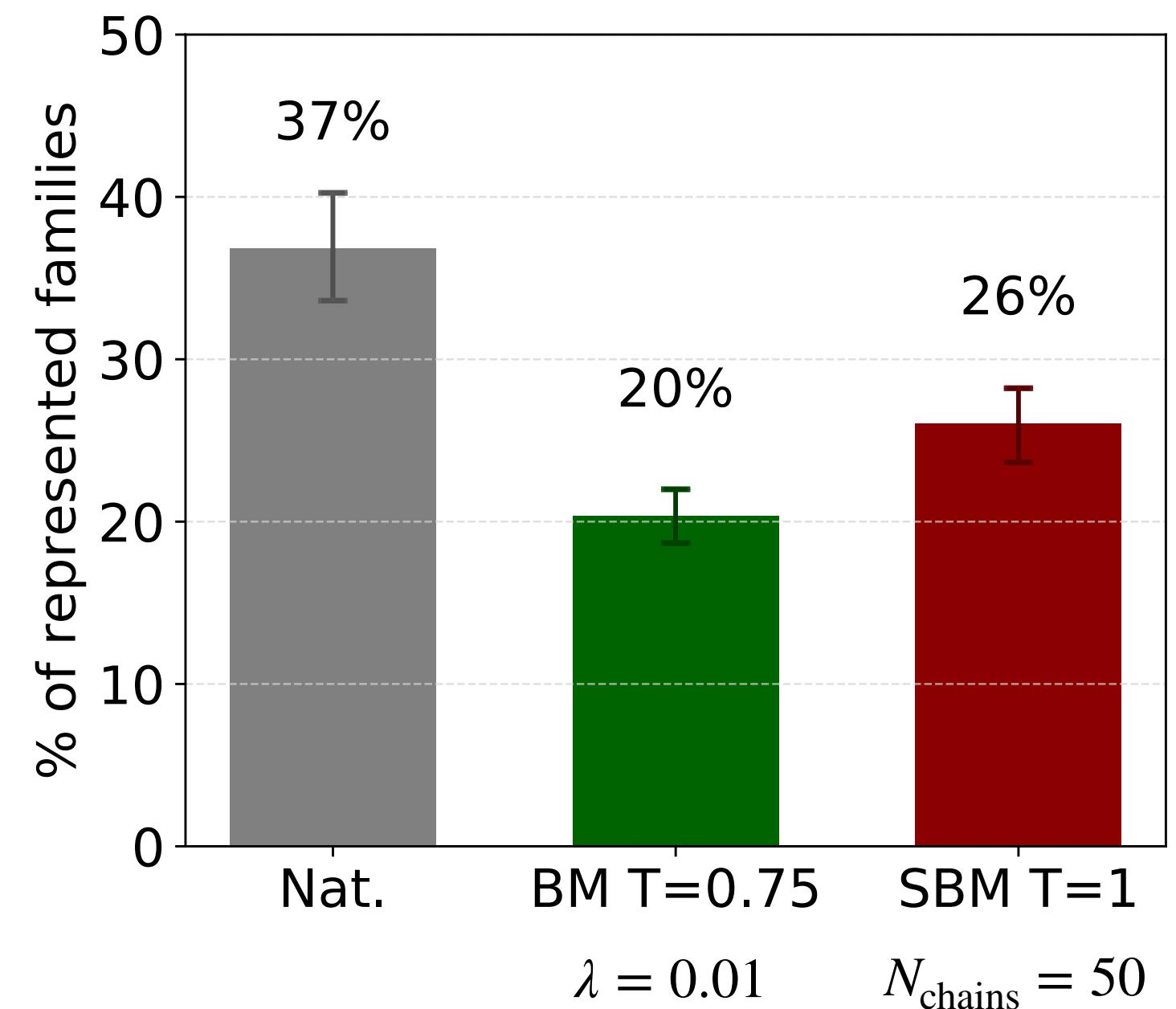
#### Novelty

*Distribution of identity to the nearest natural sequence*



#### Diversity

*% of taxonomic families represented*



#### Stochastic Boltzmann Machine

- ▶ Toy model: correct undersampling-induced biases
- ▶ Real data: generative without low-temperature sampling

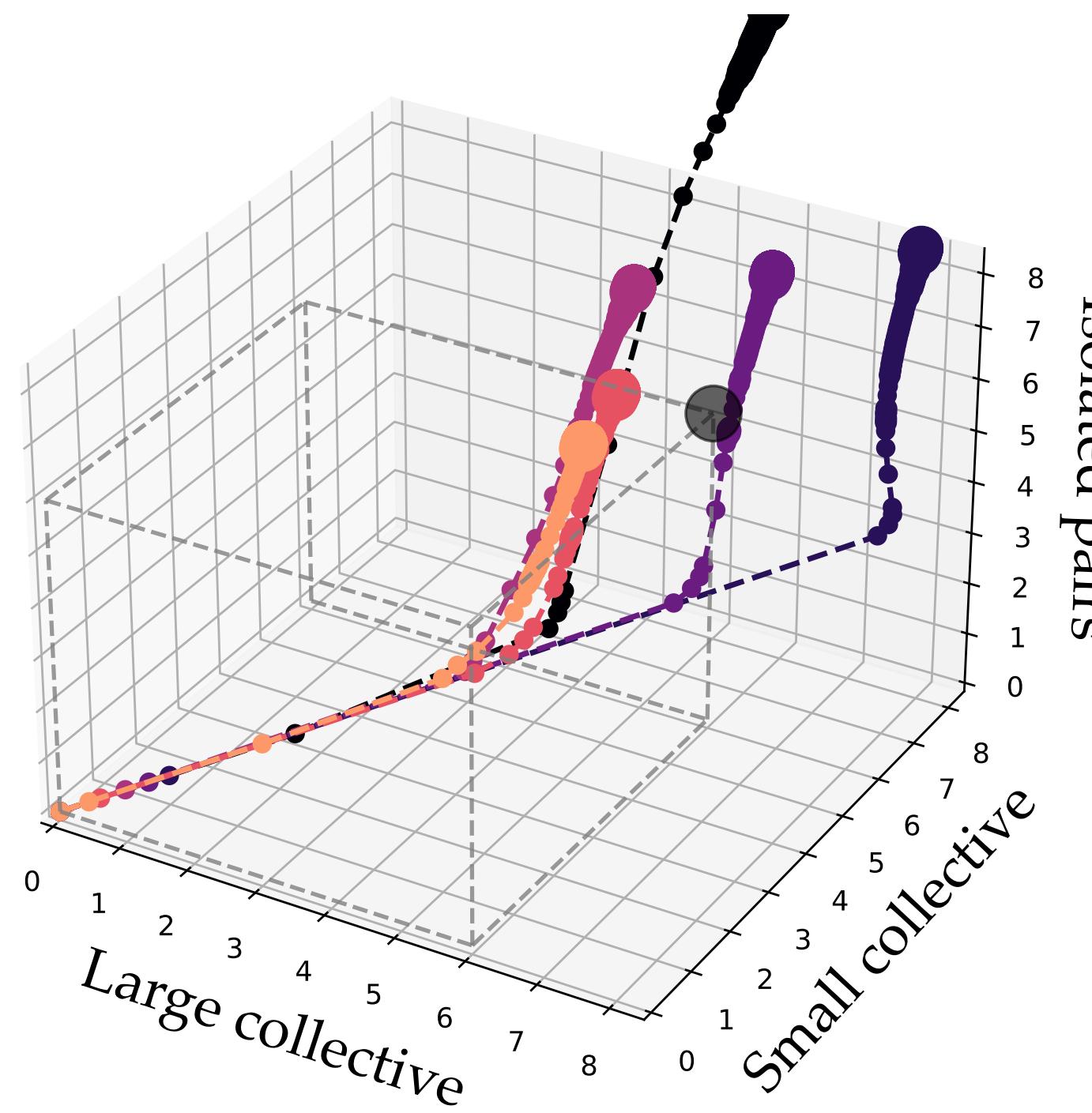
The experiments were performed by Emily Hinds

The models are inferred with  $k_{\text{mc}} = 10^5$ ,  $N_{\text{iter}} = 400$ ,  $\theta = 0.3$

890, 193 and 180 sequences were tested for the natural dataset, the SBM and the BM models respectively

# The undersampling problem

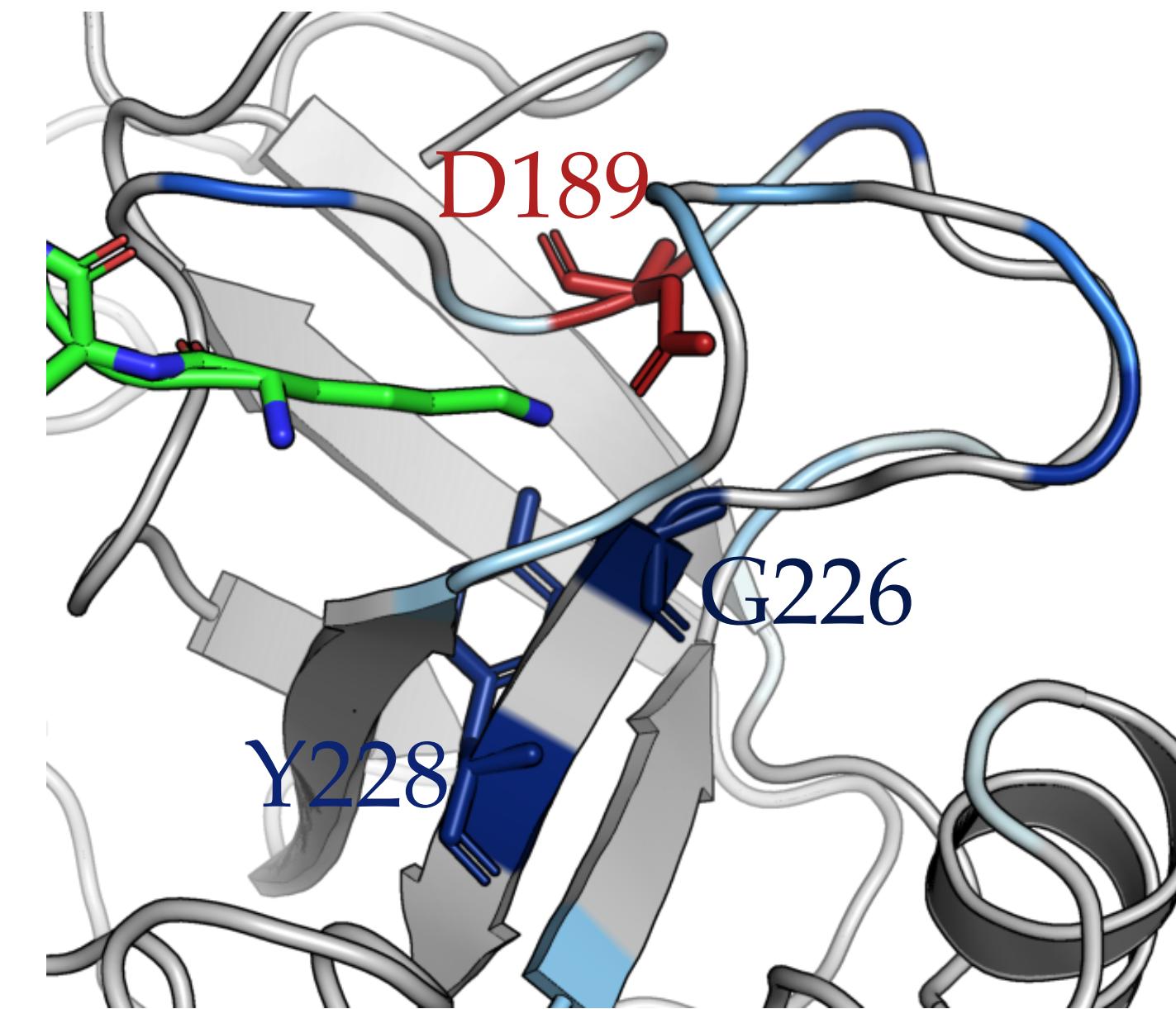
Overcoming undersampling induced biases



*In collaboration with Emily Hinds, Yaakov Kleeorin, Rama Ranganathan (University of Chicago, USA)*

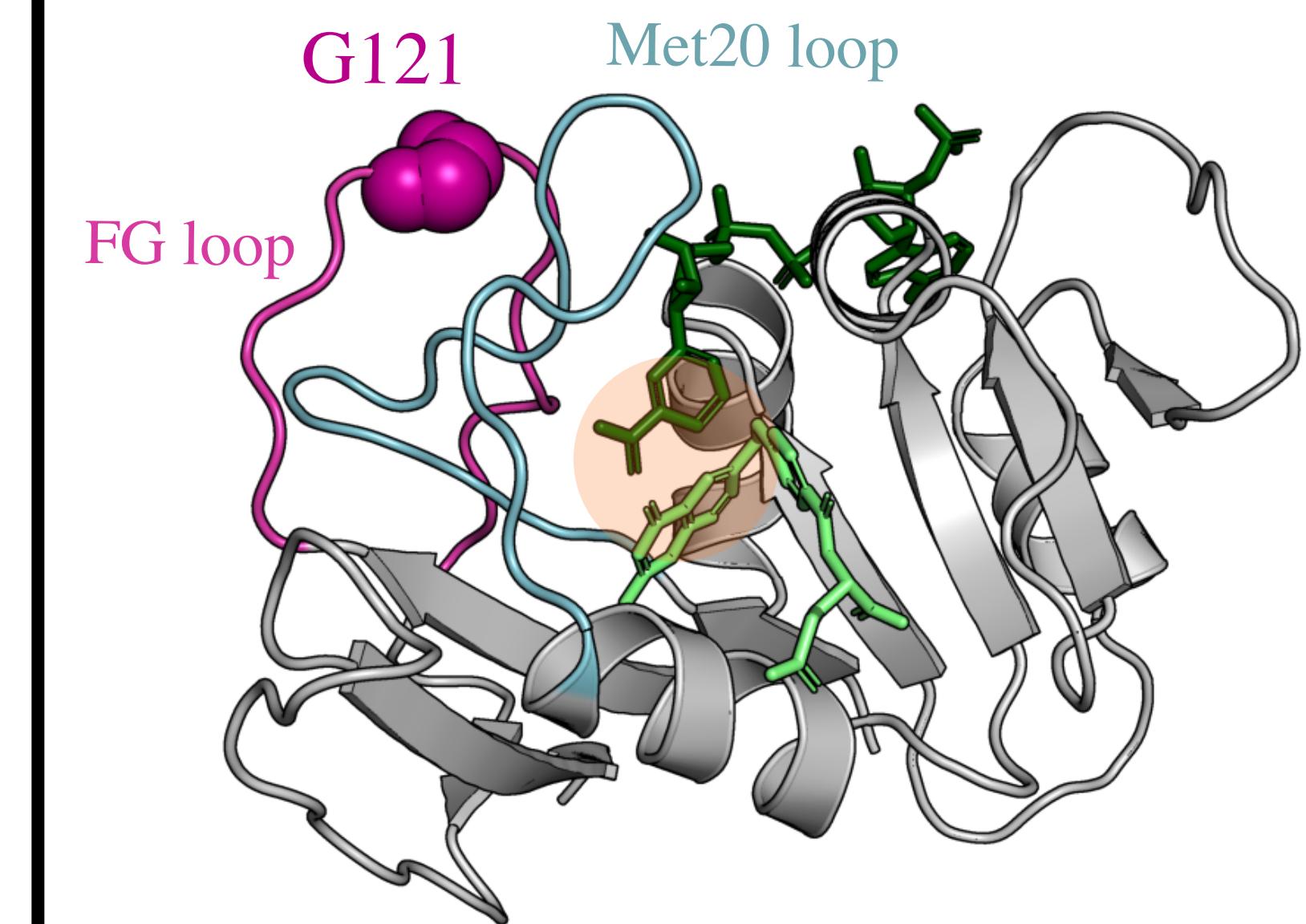
# Investigate protein properties with statistical learning

Specificity mechanism in S1A family



*In collaboration with Amaury Paveyranne, Timothé Lucas, Shoichi Yip, Clément Nizak (LJP, Sorbonne University, France)*

Allosteric network of *E. Coli* DHFR



*In collaboration with Paul Guenon, Damien Laage, Guillaume Stirnemann (ENS, France), Clément Nizak (LJP, France), Karolina Filipowska, Kim Reynolds (University of Texas, USA)*

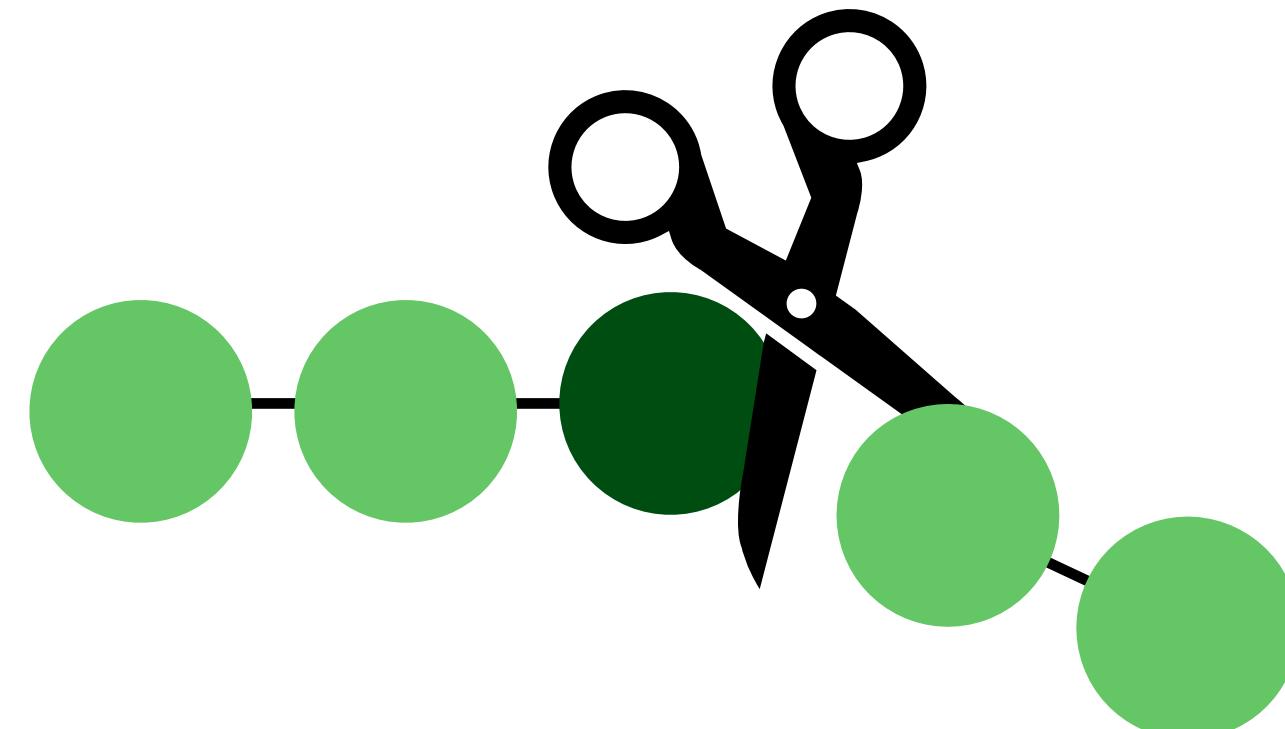
# Investigate the determinant of specificity within S1A family

*In collaboration with Amaury Paveyranne, Timothé Lucas,  
Shoichi Yip, Clément Nizak (LJP, Sorbonne University, France)*

# S1A serine proteases protein family

## The basics

- ▶ Catalyzes peptide bound hydrolysis
- ▶ Many substrate specificities



Trypsin

Positively charged

Chymotrypsin

Aromatic

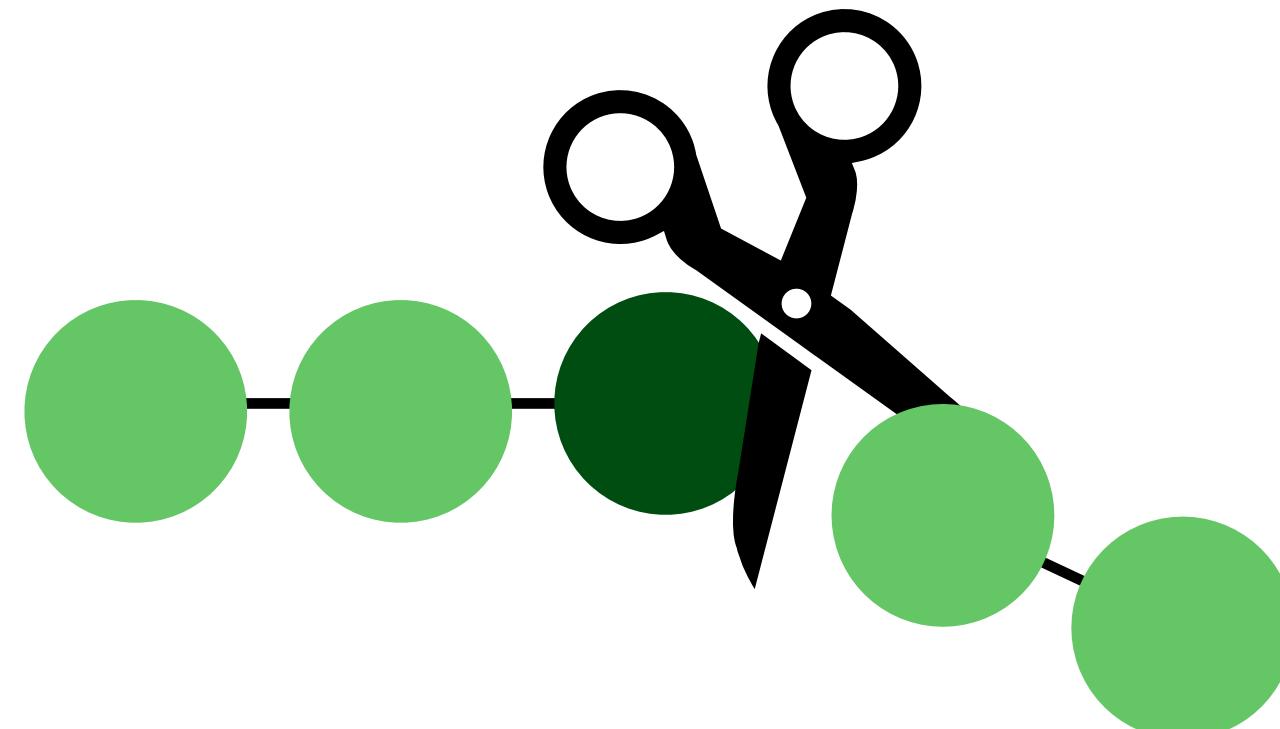
# S1A serine proteases protein family

## The basics

- ▶ Catalyzes peptide bound hydrolysis
- ▶ Many substrate specificities

## Specificity

- ▶ Cleavage after a specific amino-acid
- ▶ Efficiency can vary up to  $10^5$ -fold depending on the amino acid
- ▶ Mechanism not fully understood



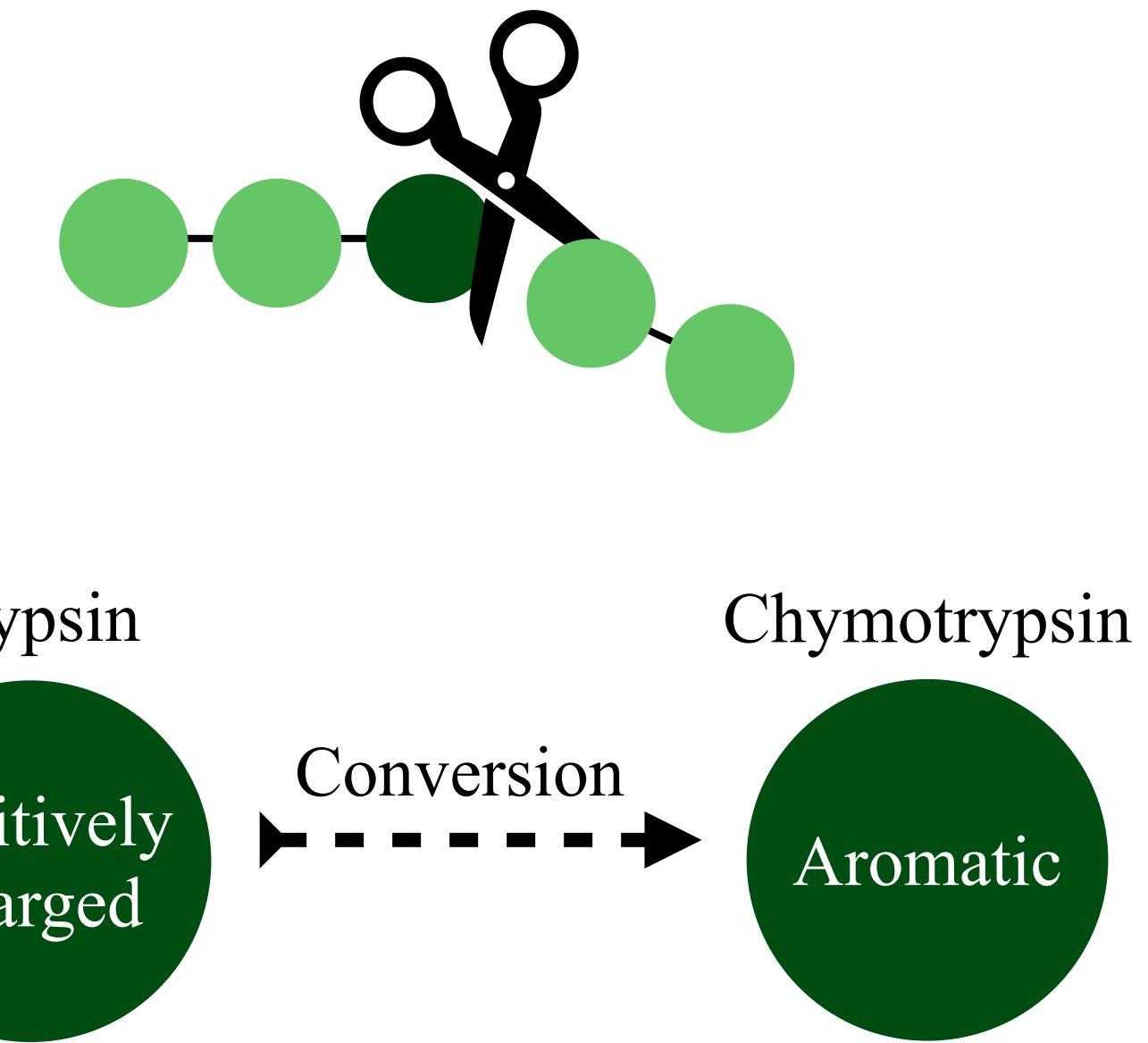
Trypsin



Chymotrypsin



# A rare success in specificity conversion: Trypsin → Chymotrypsin



Gráf et al., PNAS 1988

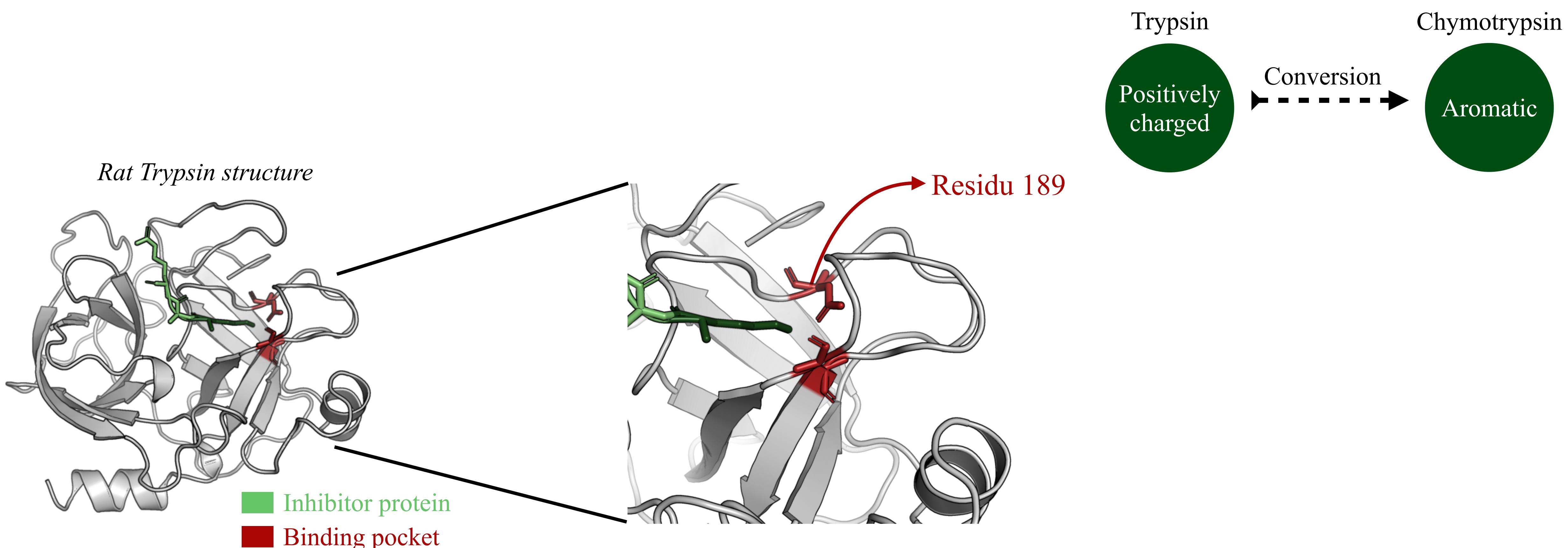
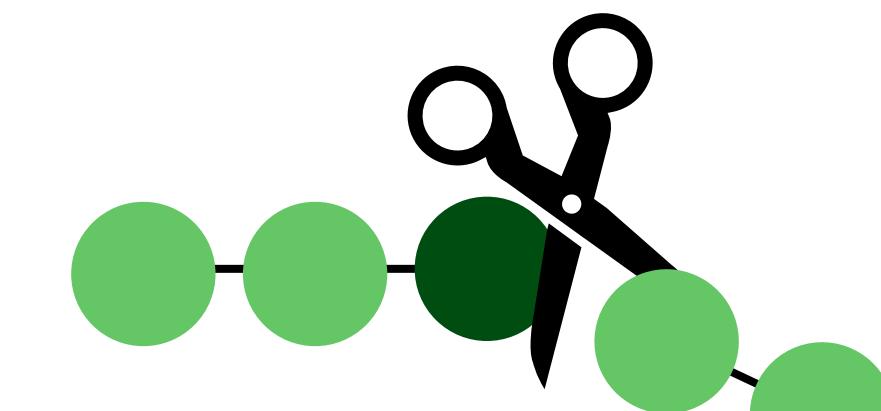
Hedstrom et al., Science 1992

Hedstrom et al., Biochemistry 1994

# A rare success in specificity conversion: Trypsin → Chymotrypsin

## Structural approach

- ▶ Focus on binding pocket
- ▶ Residu 189 sufficient to convert specificity? No (Gráf *et al.* 1988)

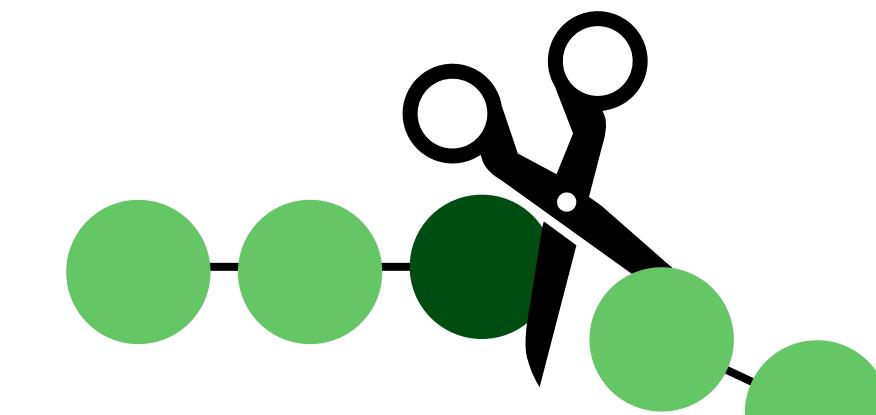


Gráf et al., PNAS 1988  
 Hedstrom et al., Science 1992  
 Hedstrom et al., Biochemistry 1994

# A rare success in specificity conversion: Trypsin → Chymotrypsin

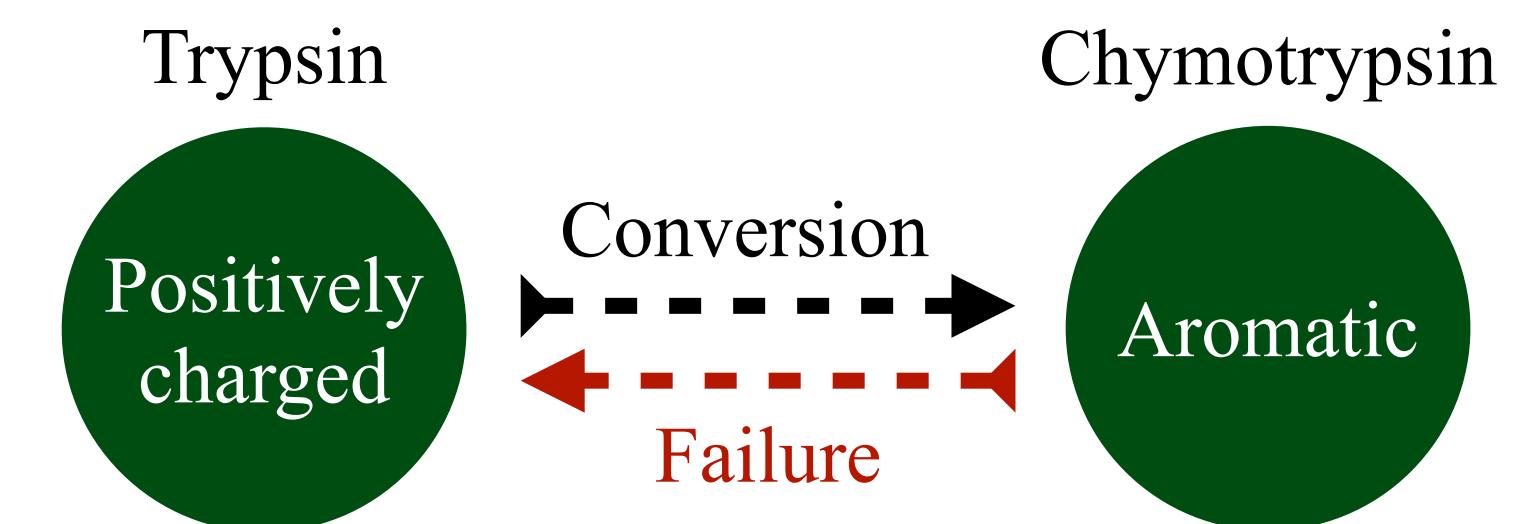
## Structural approach

- ▶ Focus on binding pocket
- ▶ Residu 189 sufficient to convert specificity? No (Gráf *et al.* 1988)

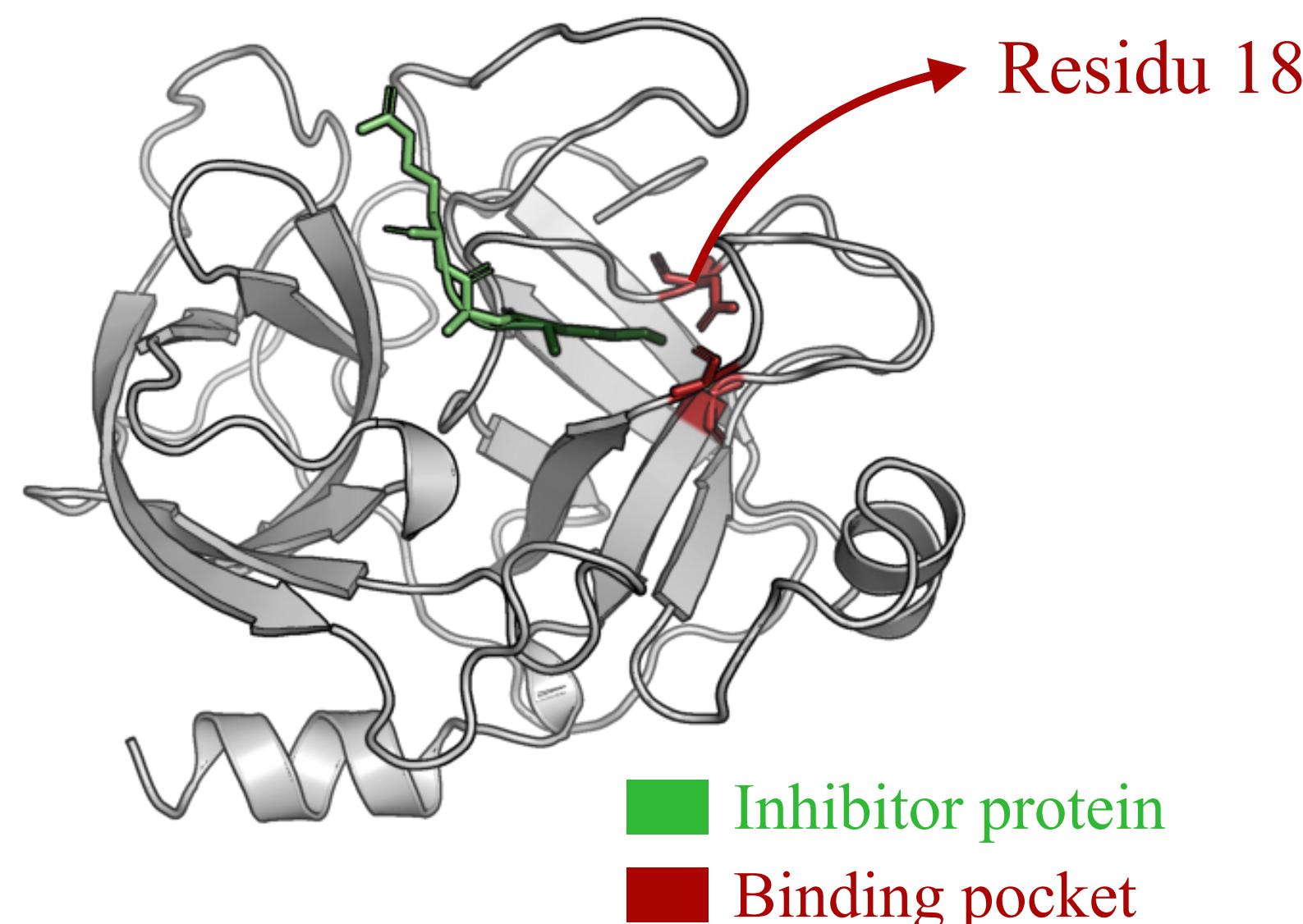


## Evolutionary approach

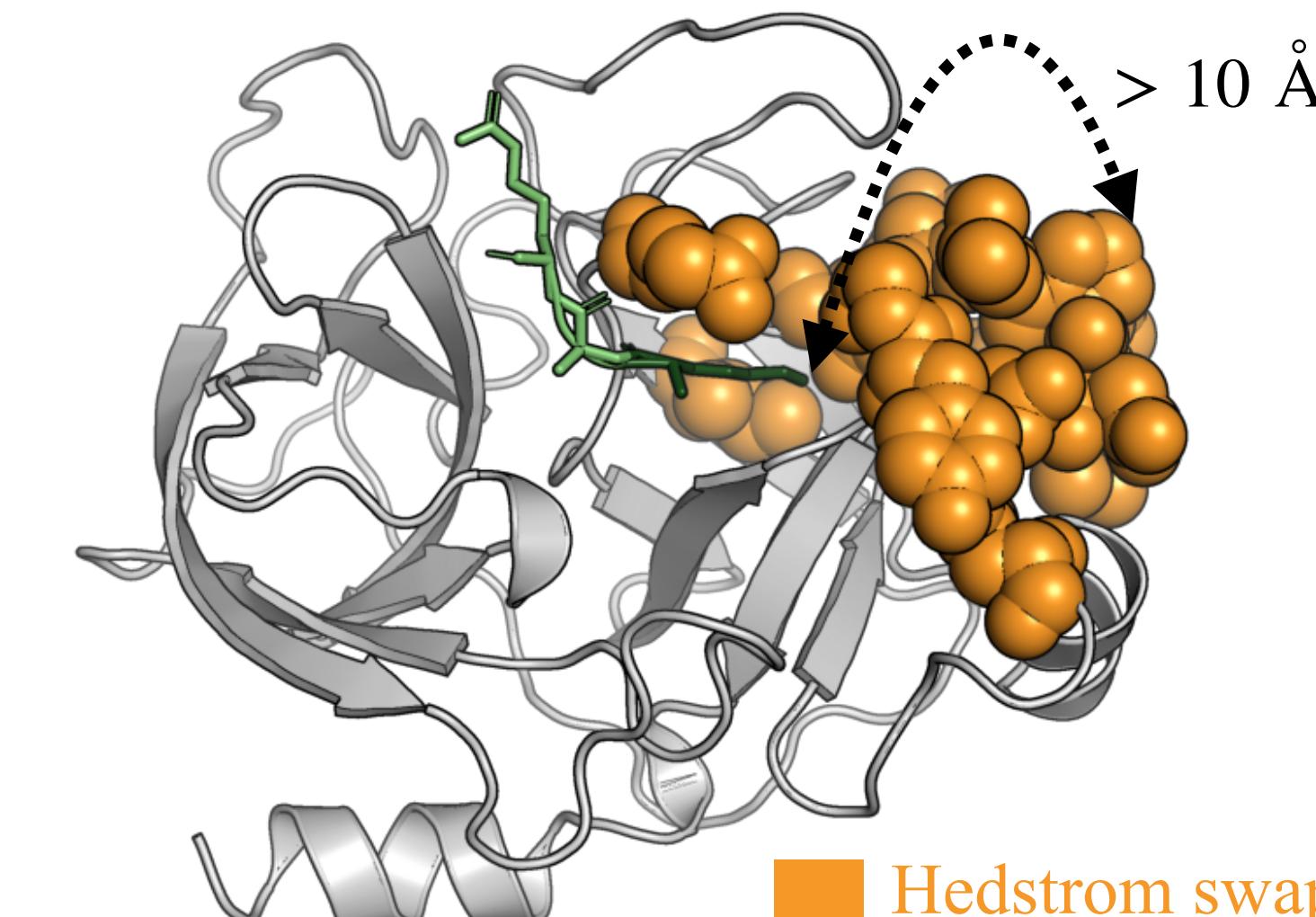
- ▶ Analysis of trypsin and chymotrypsin sequences
- ▶ Conversion with 16 mutations (Hedstrom *et al.* 1994)



*Rat Trypsin structure*



*Rat Trypsin structure*

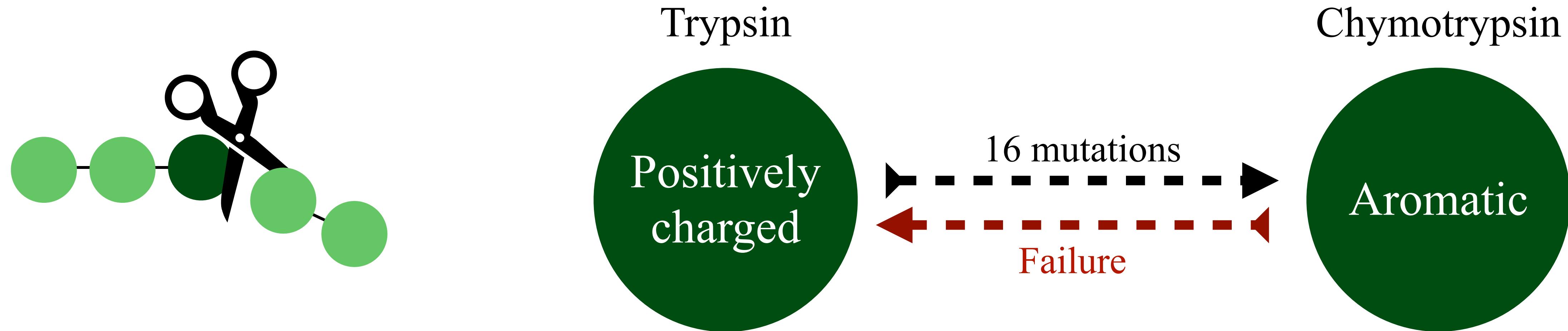


Gráf *et al.*, PNAS 1988

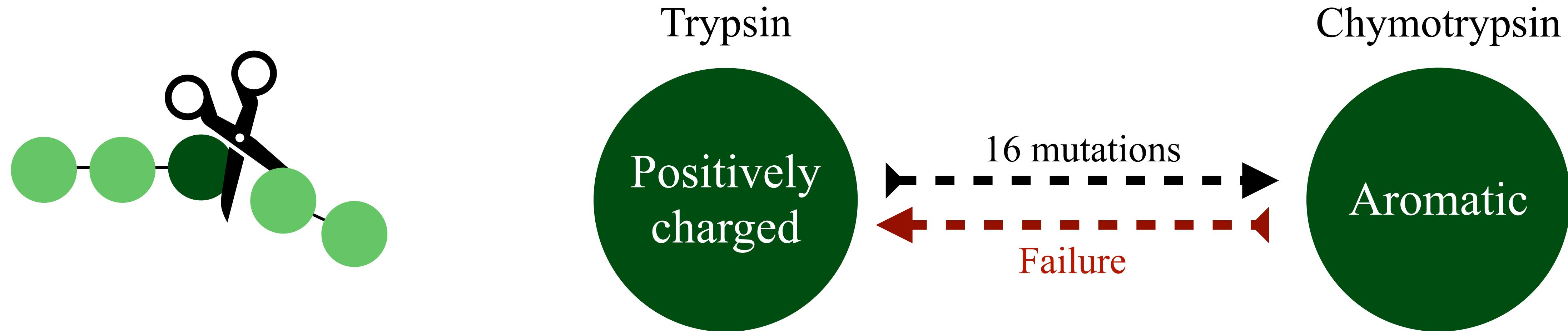
Hedstrom *et al.*, Science 1992

Hedstrom *et al.*, Biochemistry 1994

# Boltzmann Machine-guided specificity conversion

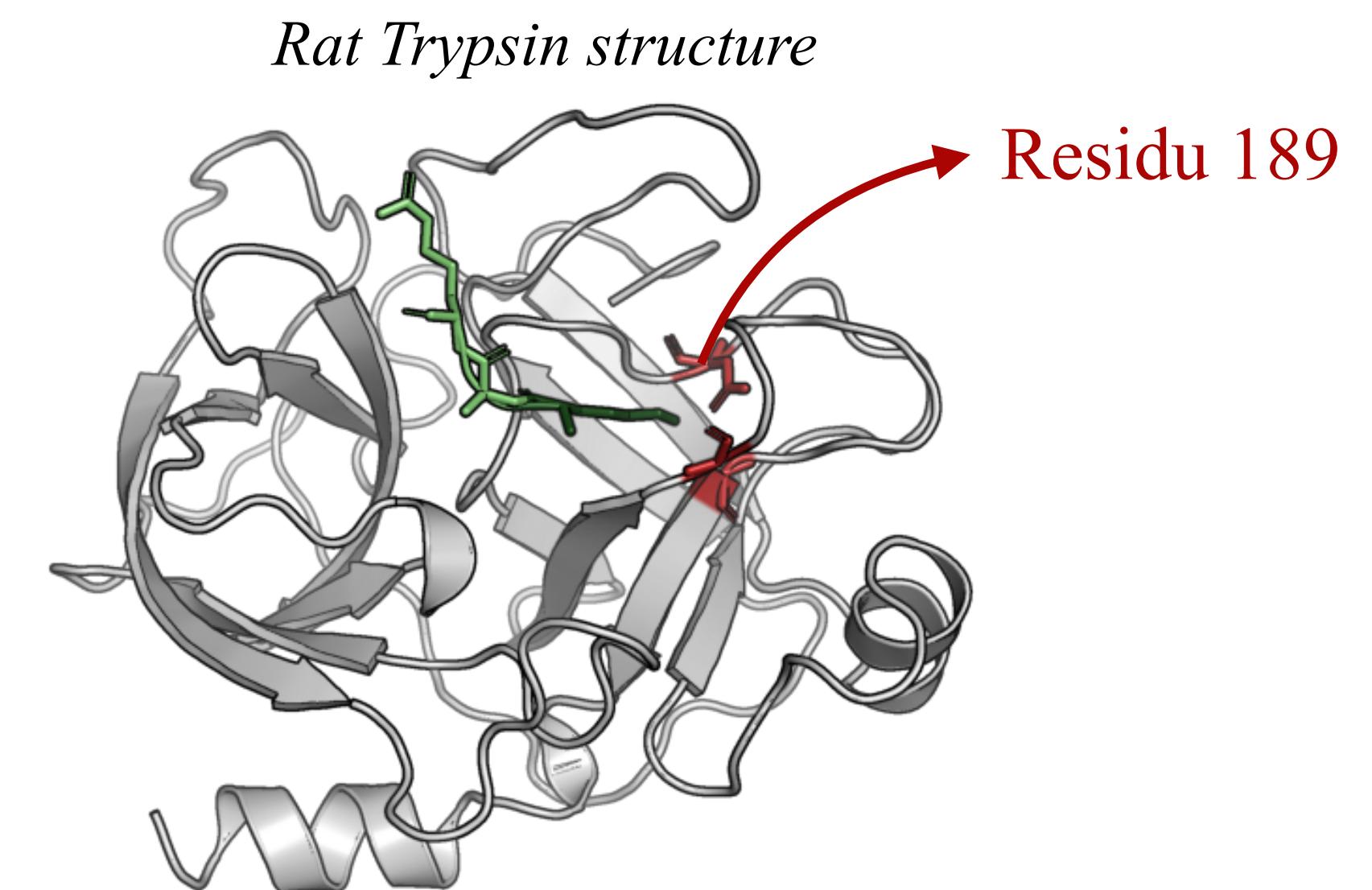


# Boltzmann Machine-guided specificity conversion

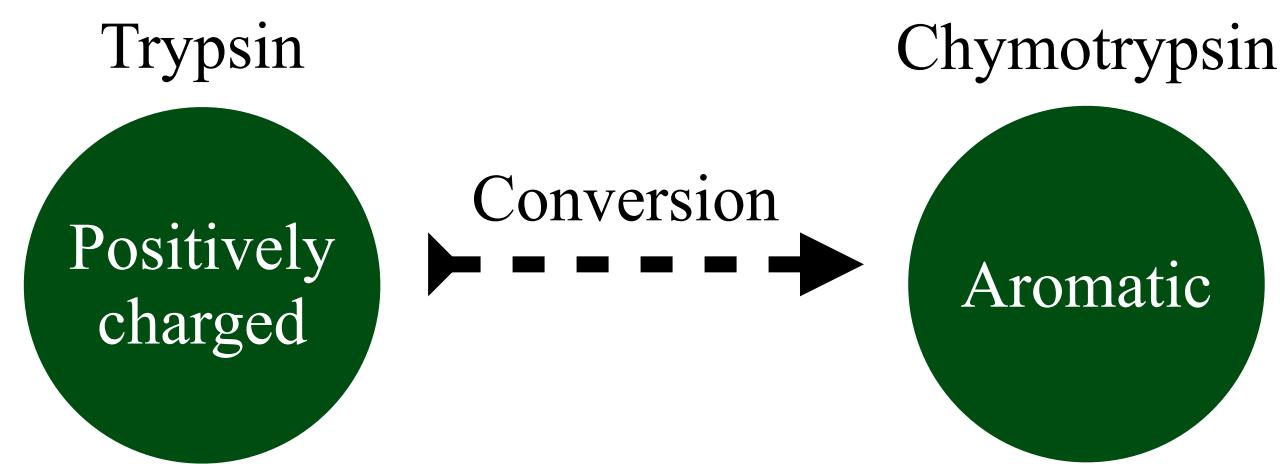


## Our approach

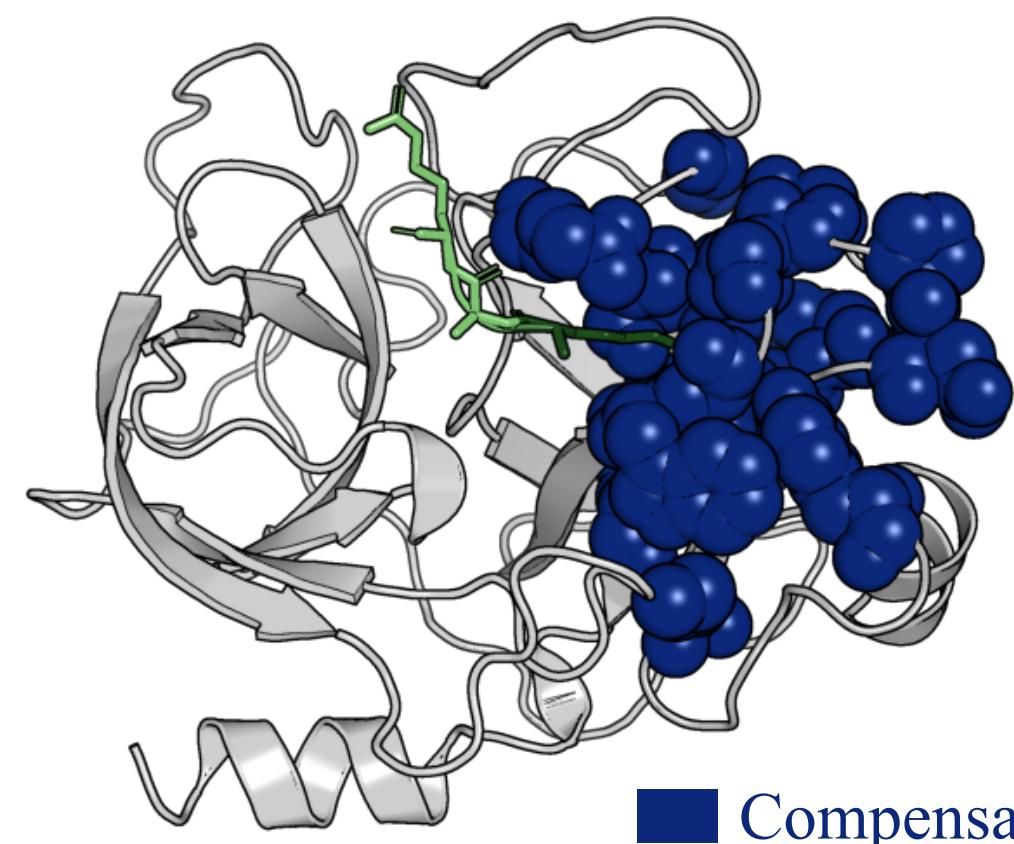
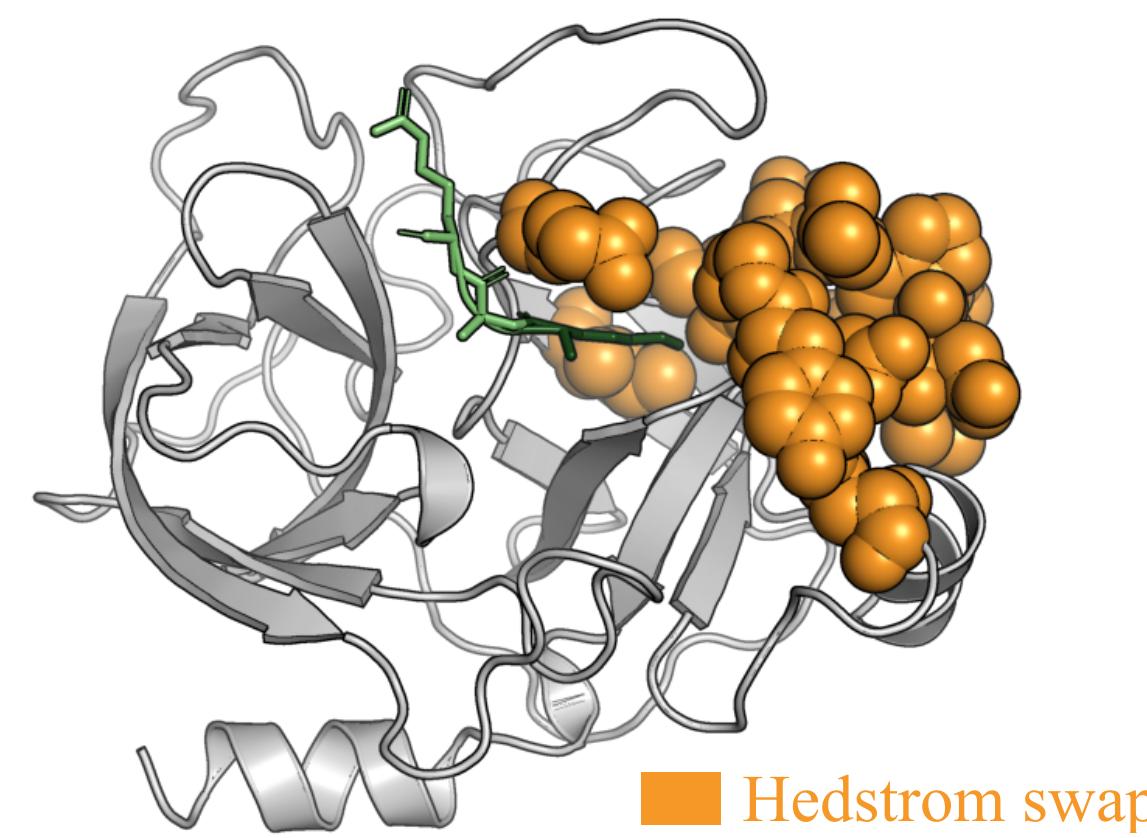
- ▶ Impose a mutation at site 189
- ▶ Predict compensatory mutations using Boltzmann Machine model
- ▶ Identify positive epistatic interactions using model couplings



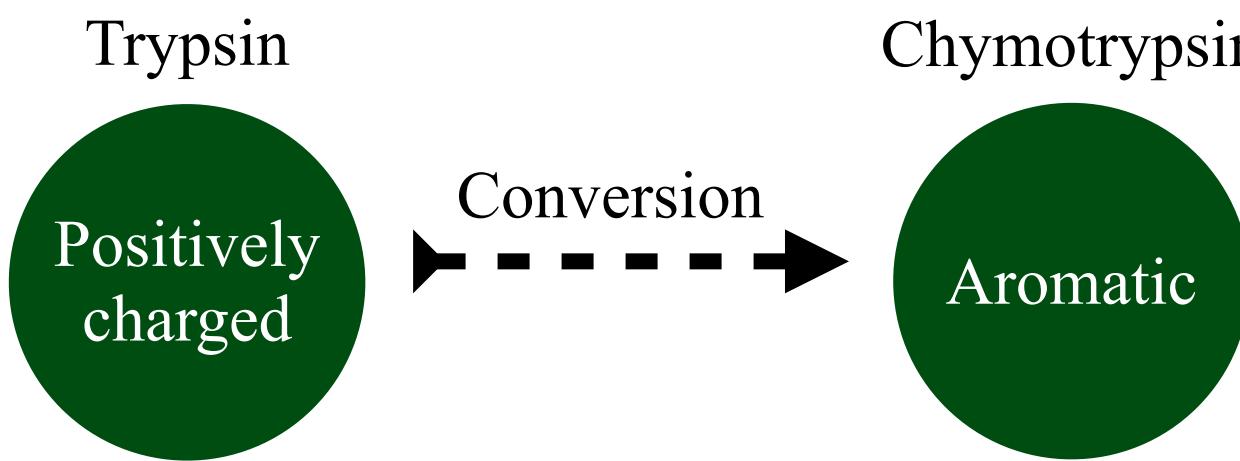
# Boltzmann Machine-guided specificity conversion



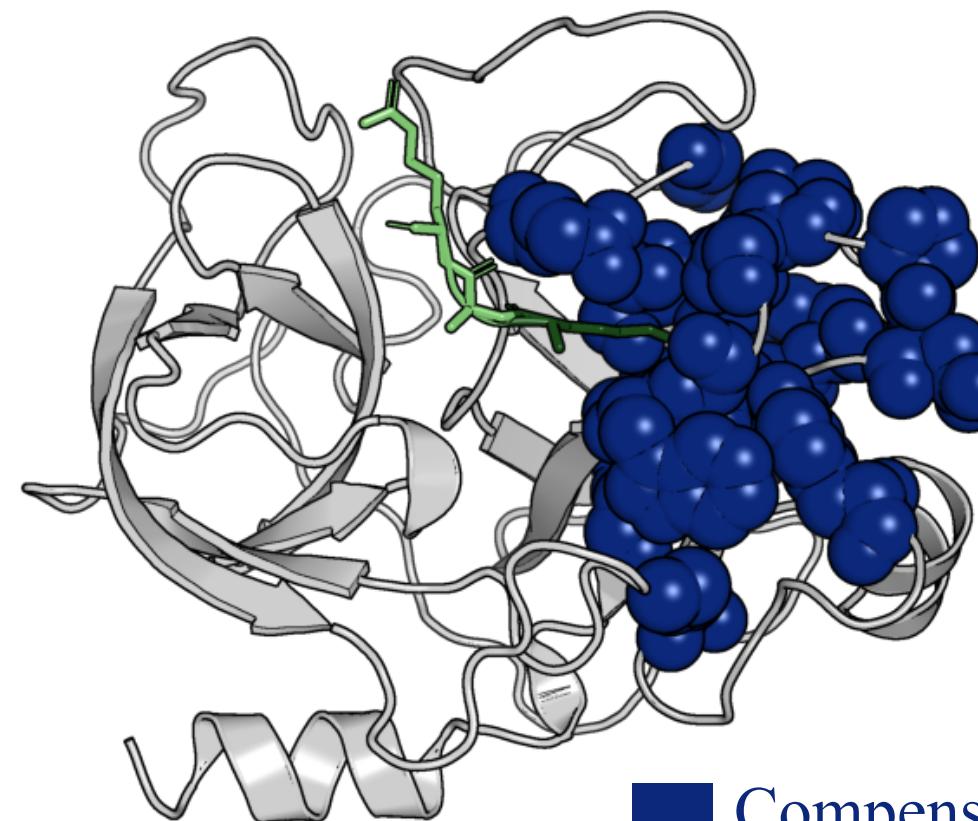
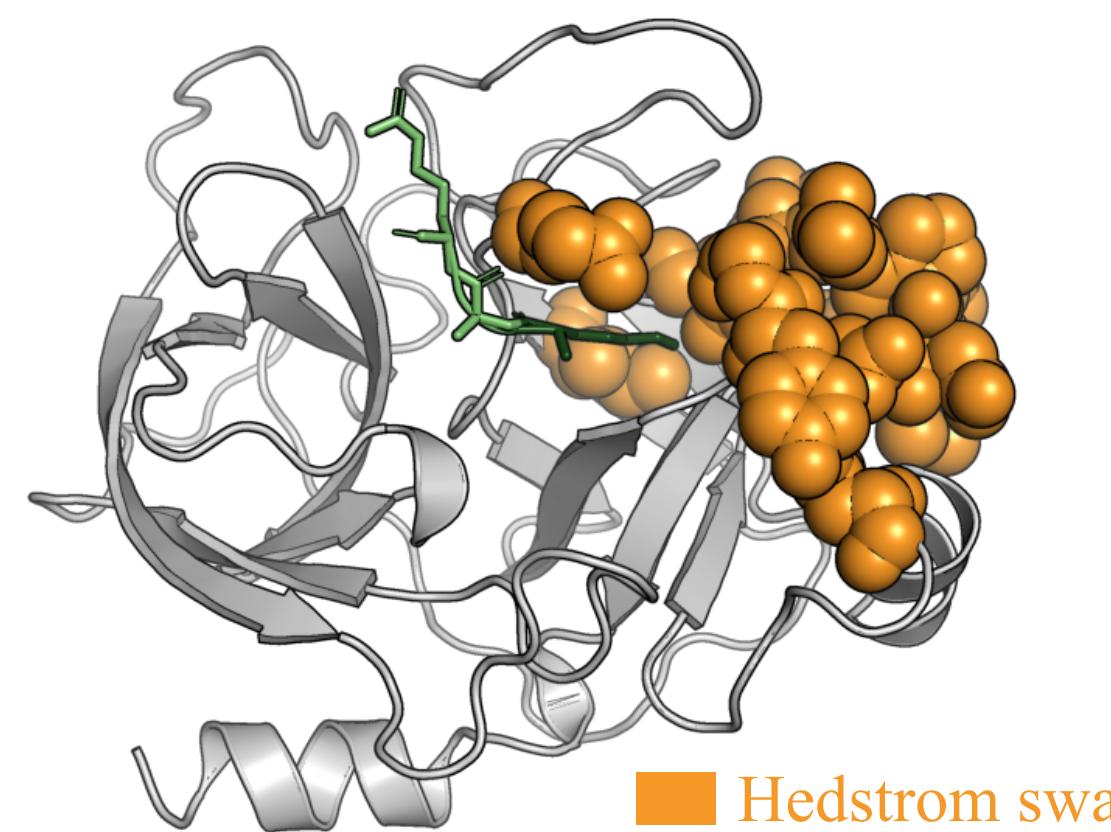
## In silico comparison with Hedstrom



# Boltzmann Machine-guided specificity conversion

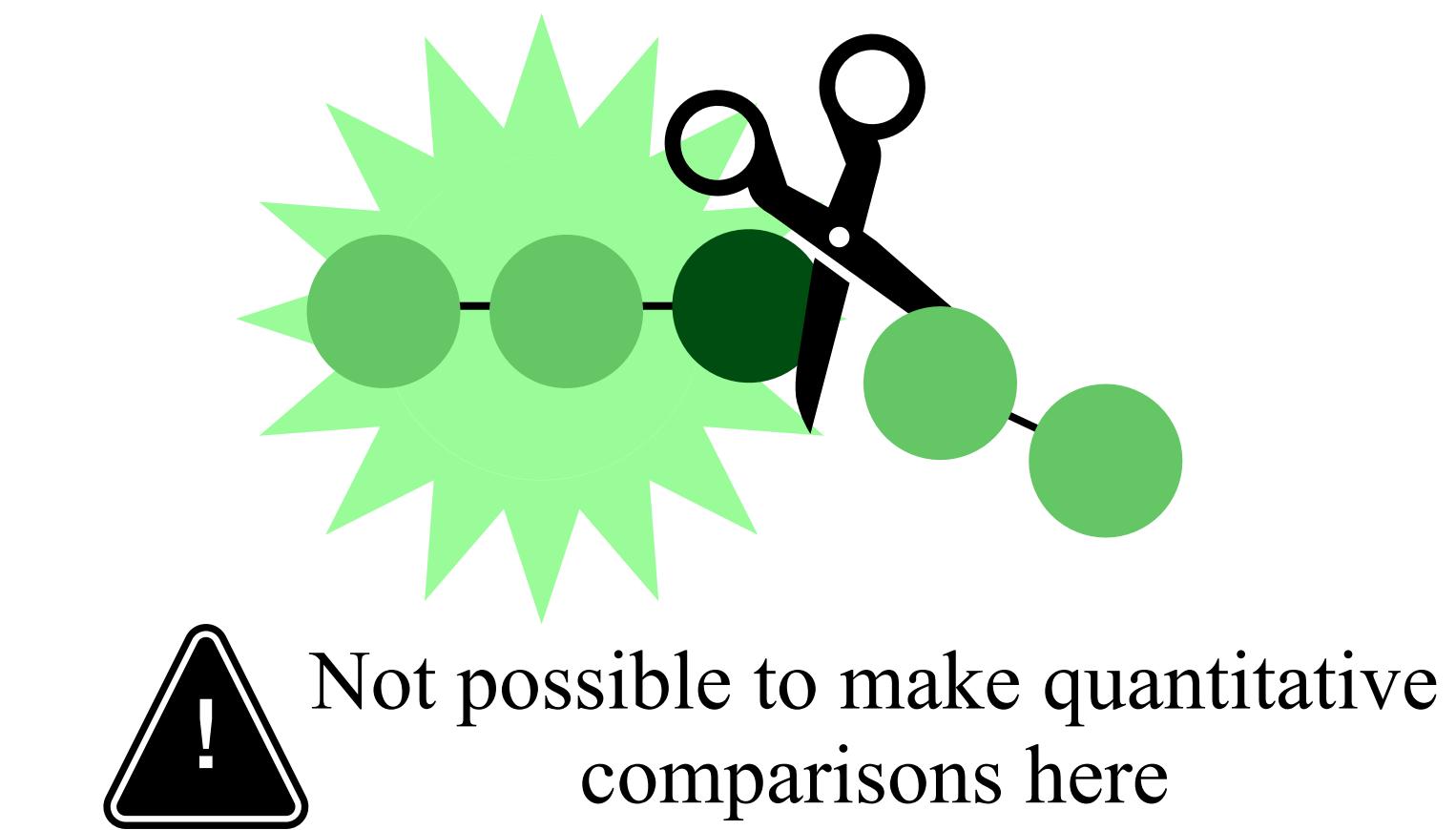
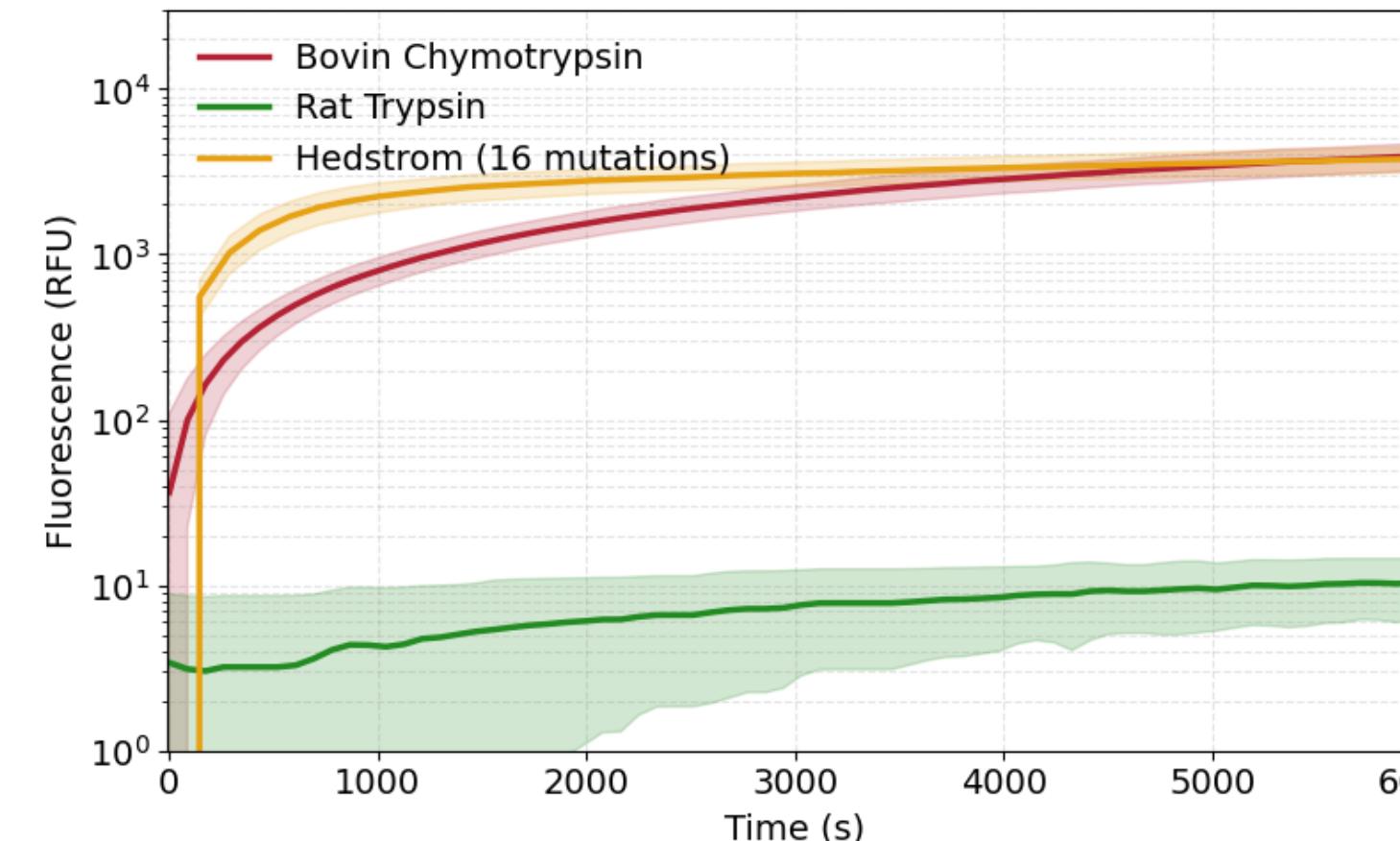


## In silico comparison with Hedstrom

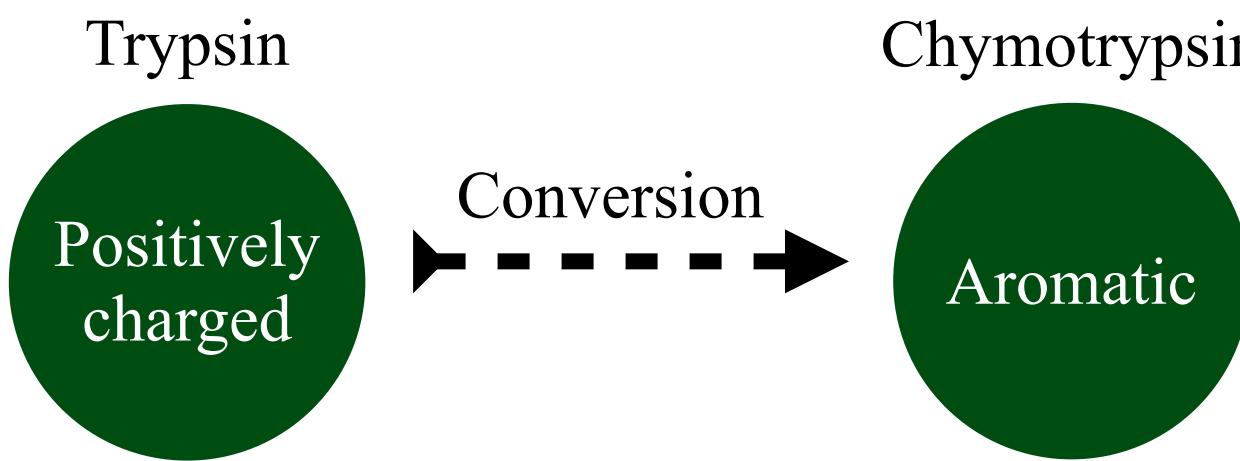


## Preliminary experimental results

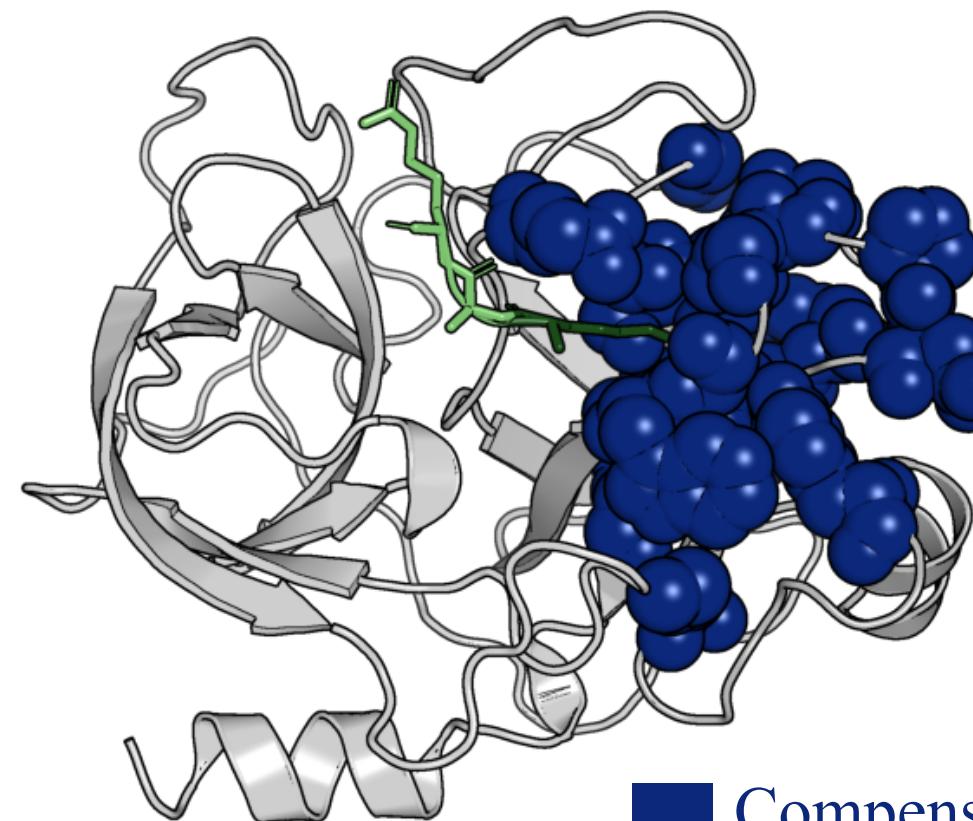
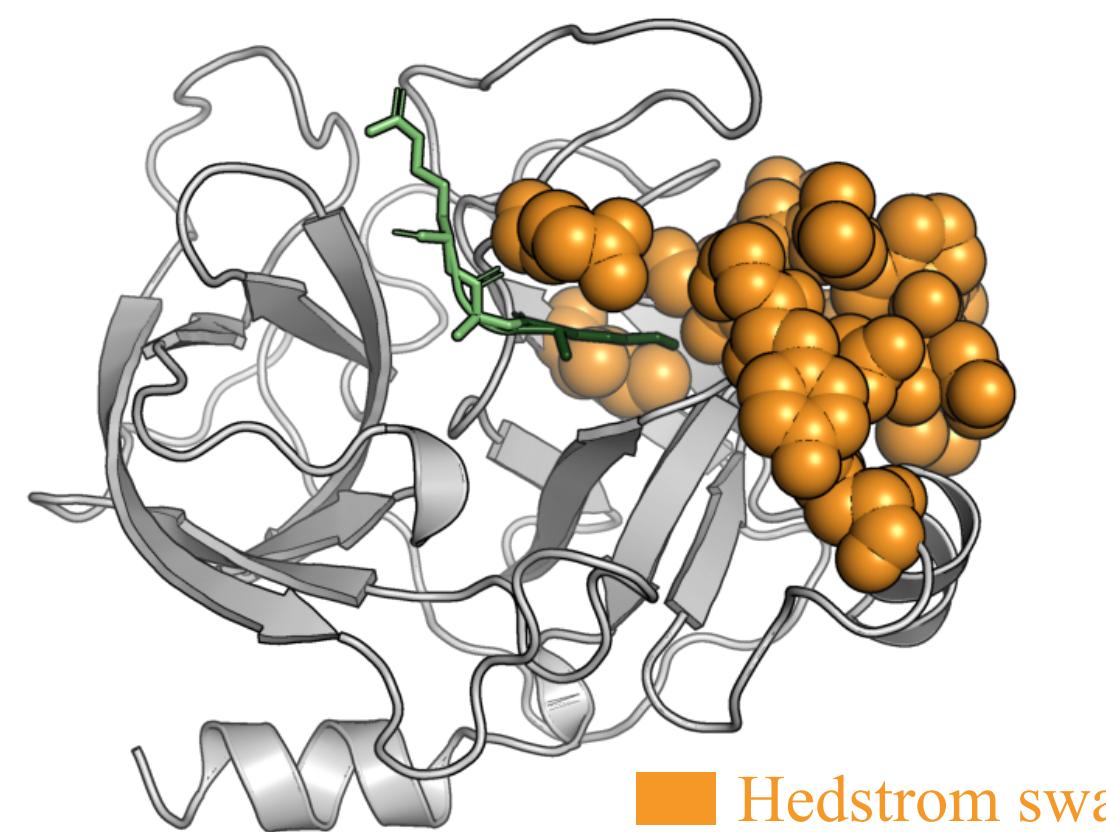
### Chymotrypsin activity



# Boltzmann Machine-guided specificity conversion

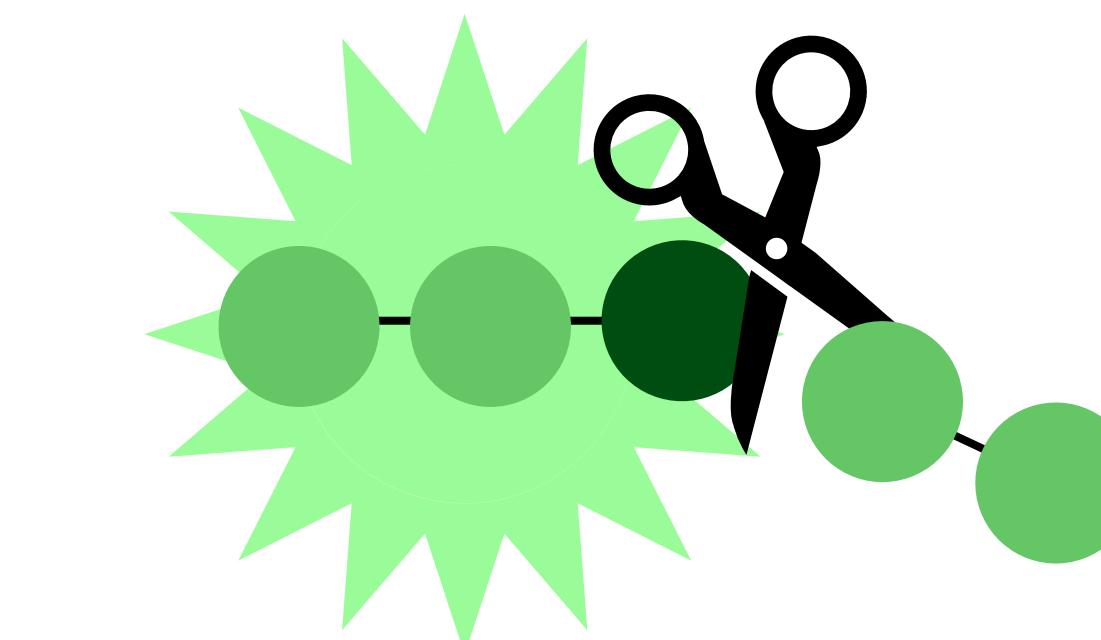
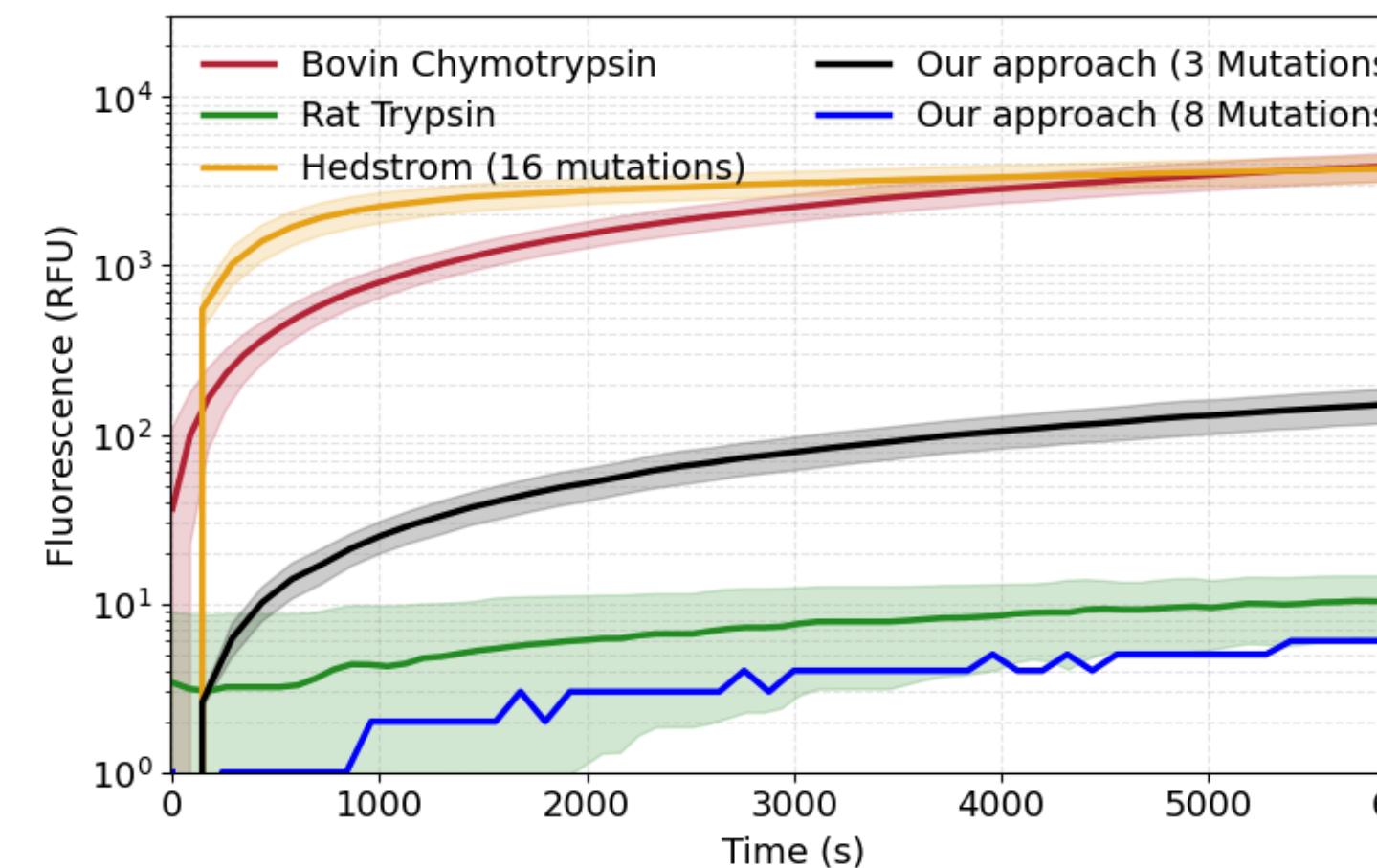


## In silico comparison with Hedstrom



## Preliminary experimental results

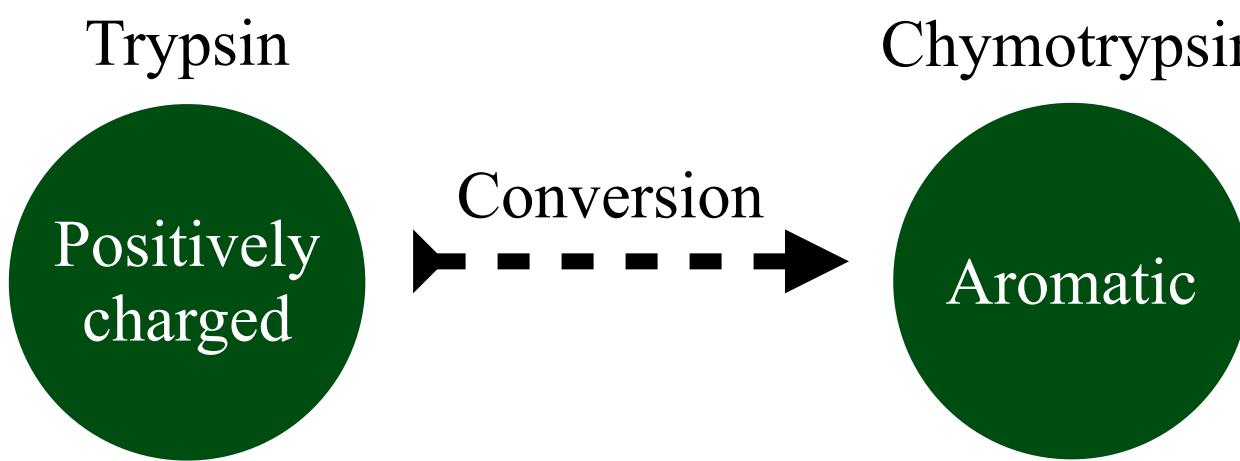
### Chymotrypsin activity



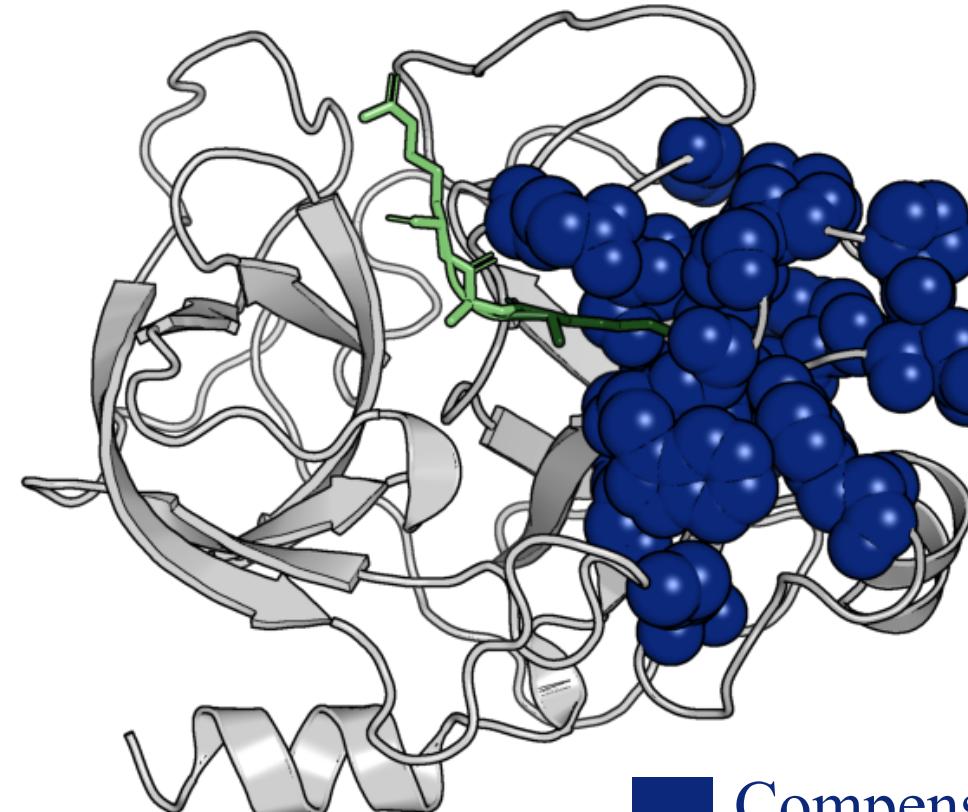
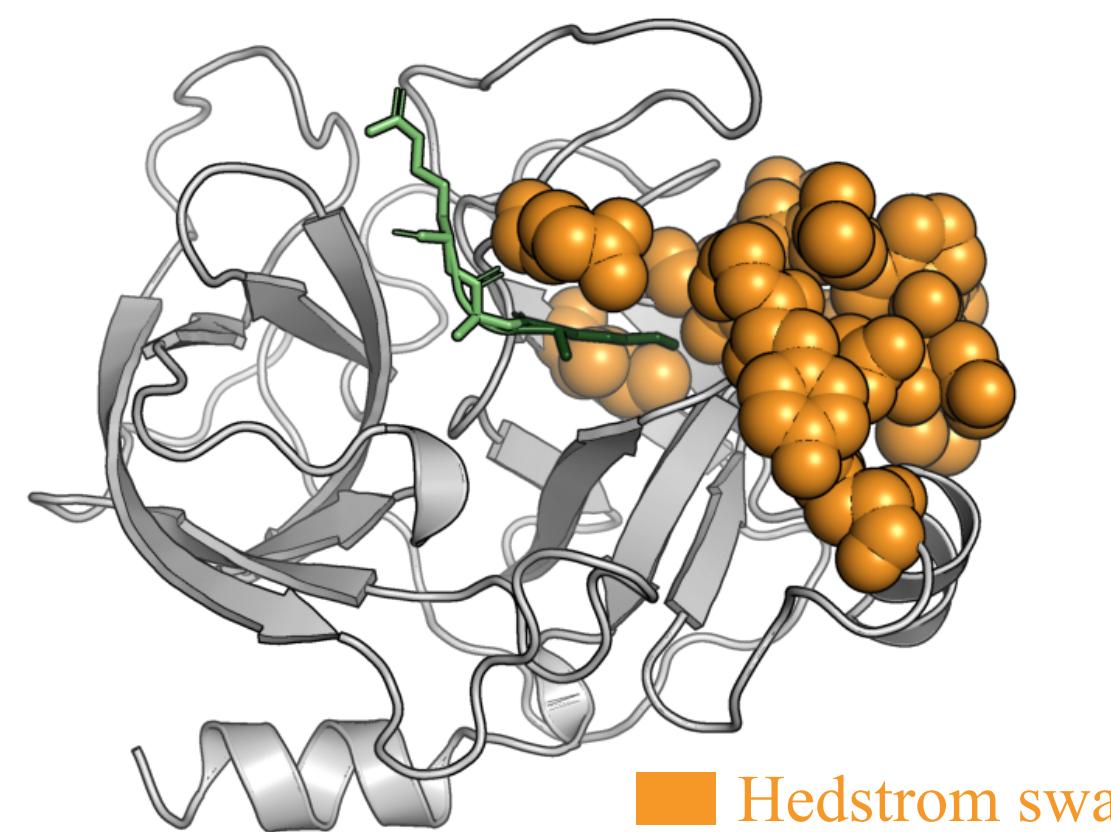
⚠ Not possible to make quantitative comparisons here

- Evidence of chymotrypsin activity 3 mutations away from rat trypsin

# Boltzmann Machine-guided specificity conversion

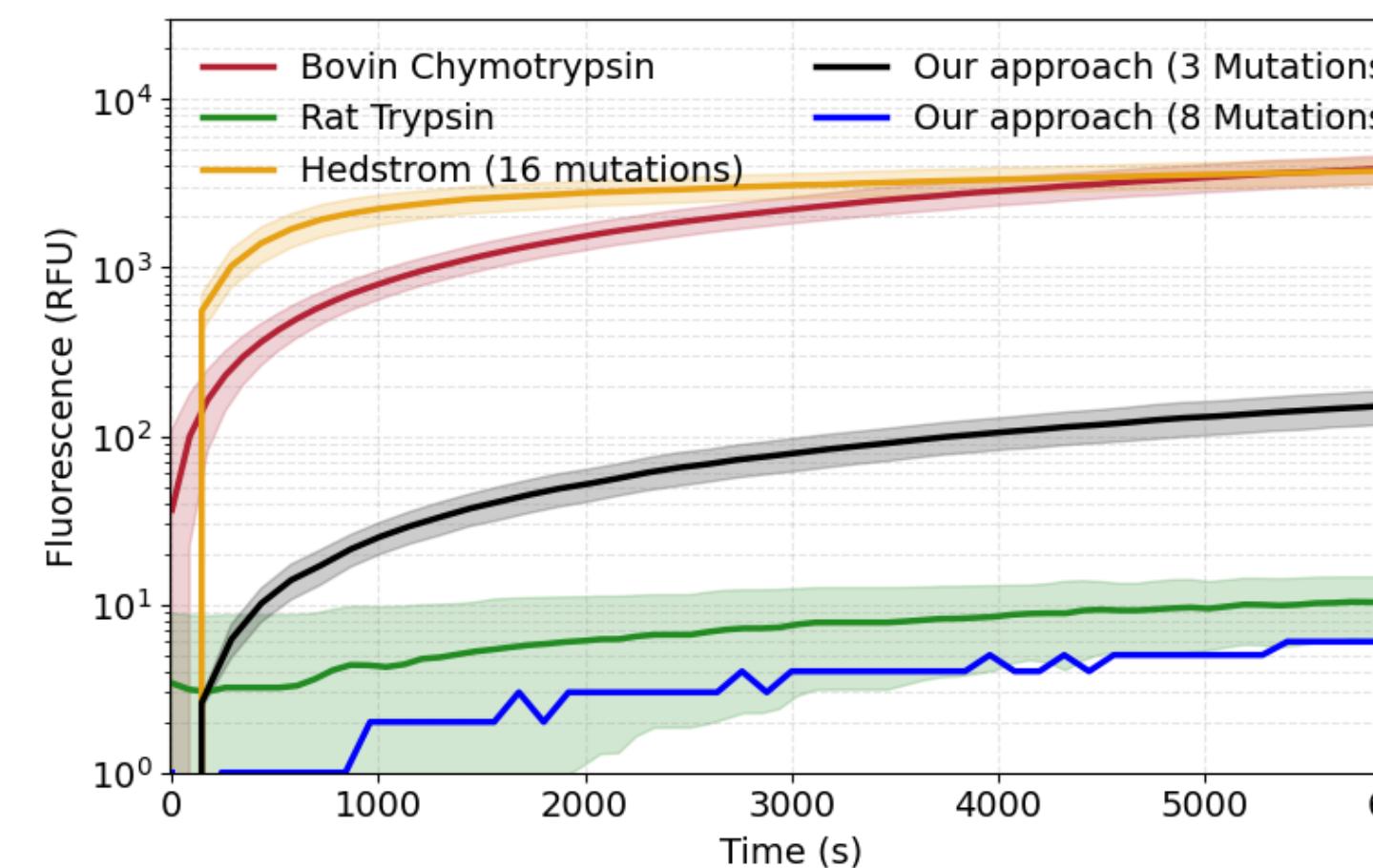


## In silico comparison with Hedstrom

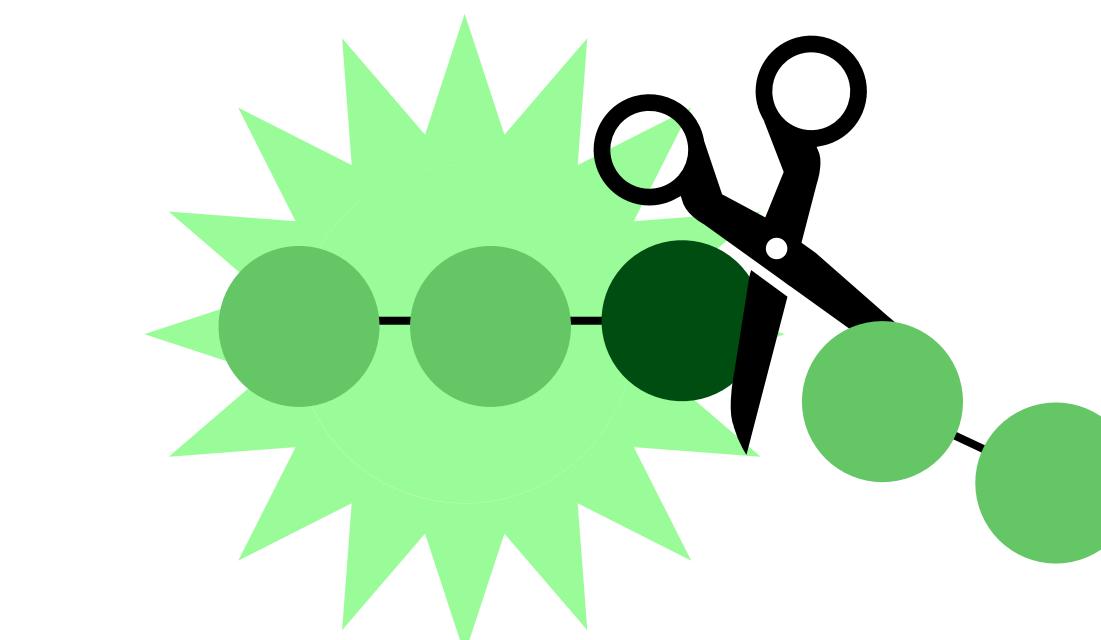
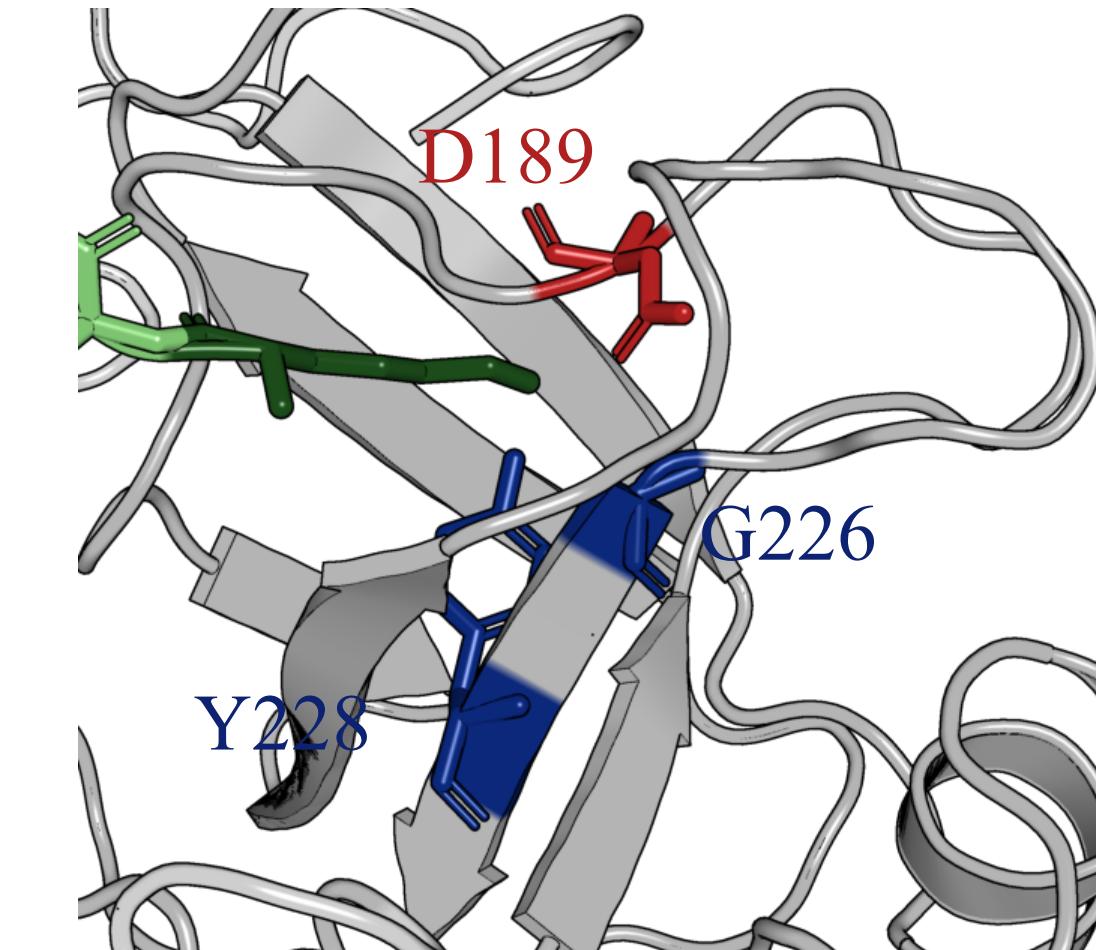


## Preliminary experimental results

### Chymotrypsin activity



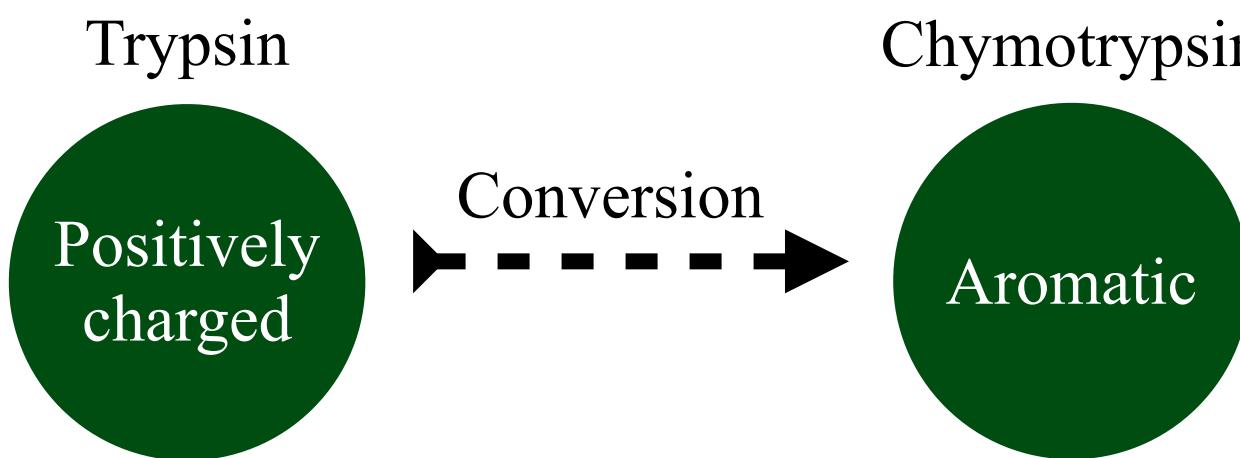
### Zoom on the binding pocket



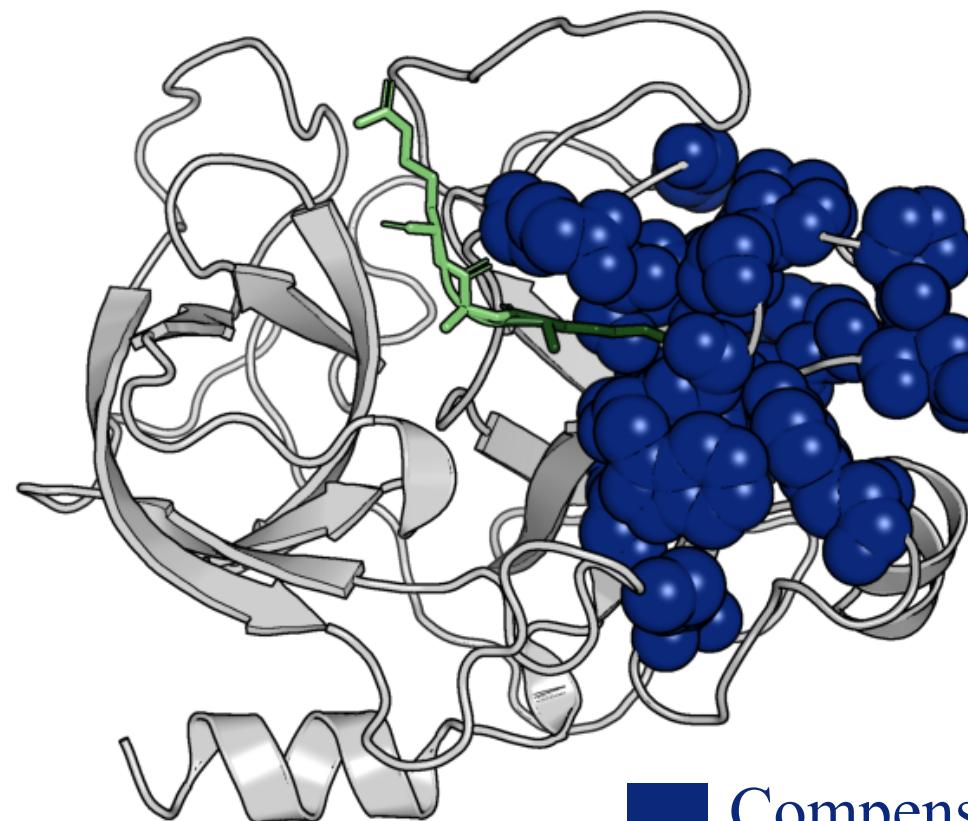
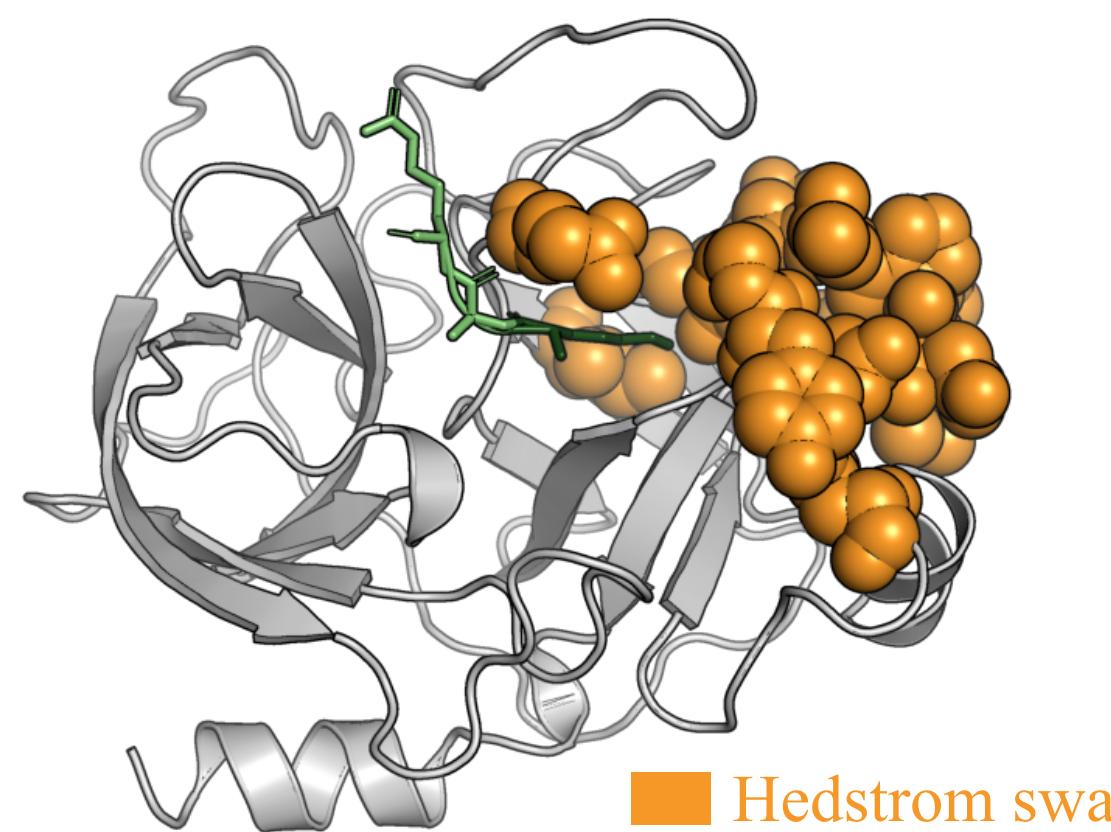
Not possible to make quantitative comparisons here

- ▶ Evidence of chymotrypsin activity 3 mutations away from rat trypsin
- ▶ Mutations located in the binding pocket

# Boltzmann Machine-guided specificity conversion

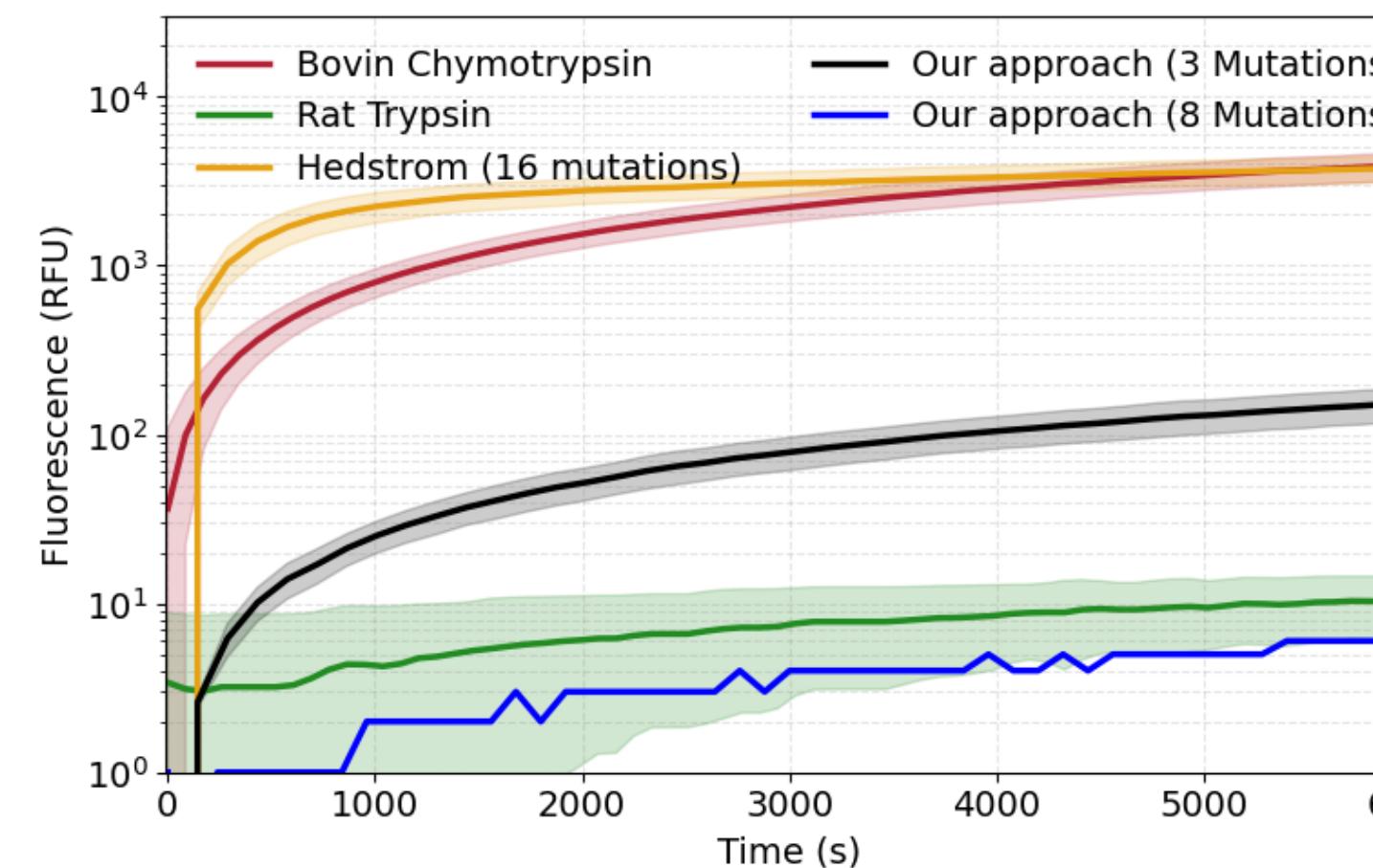


## In silico comparison with Hedstrom

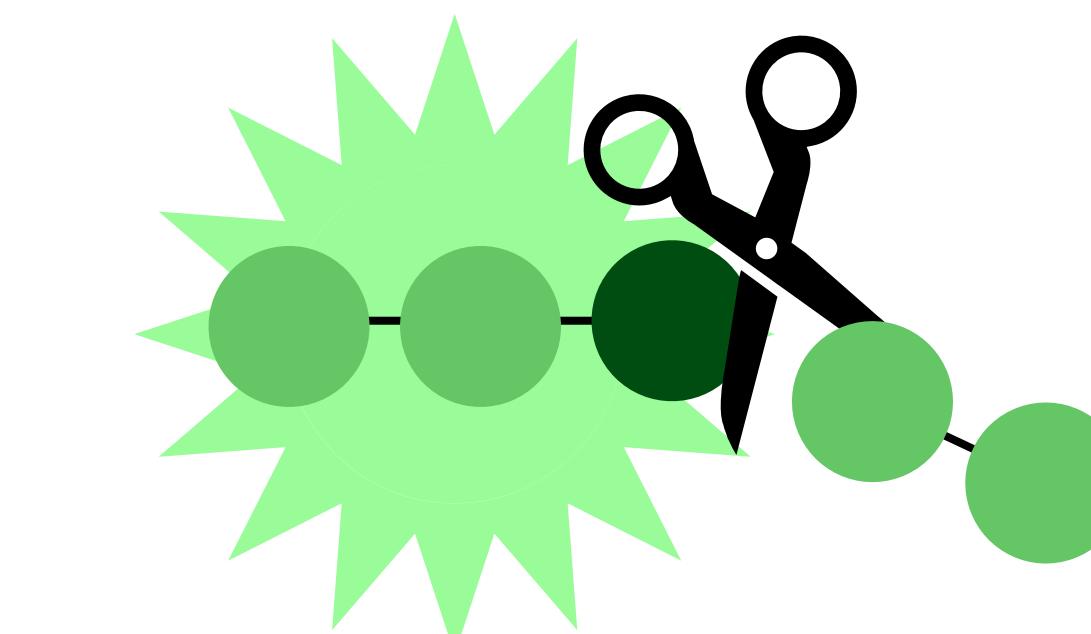
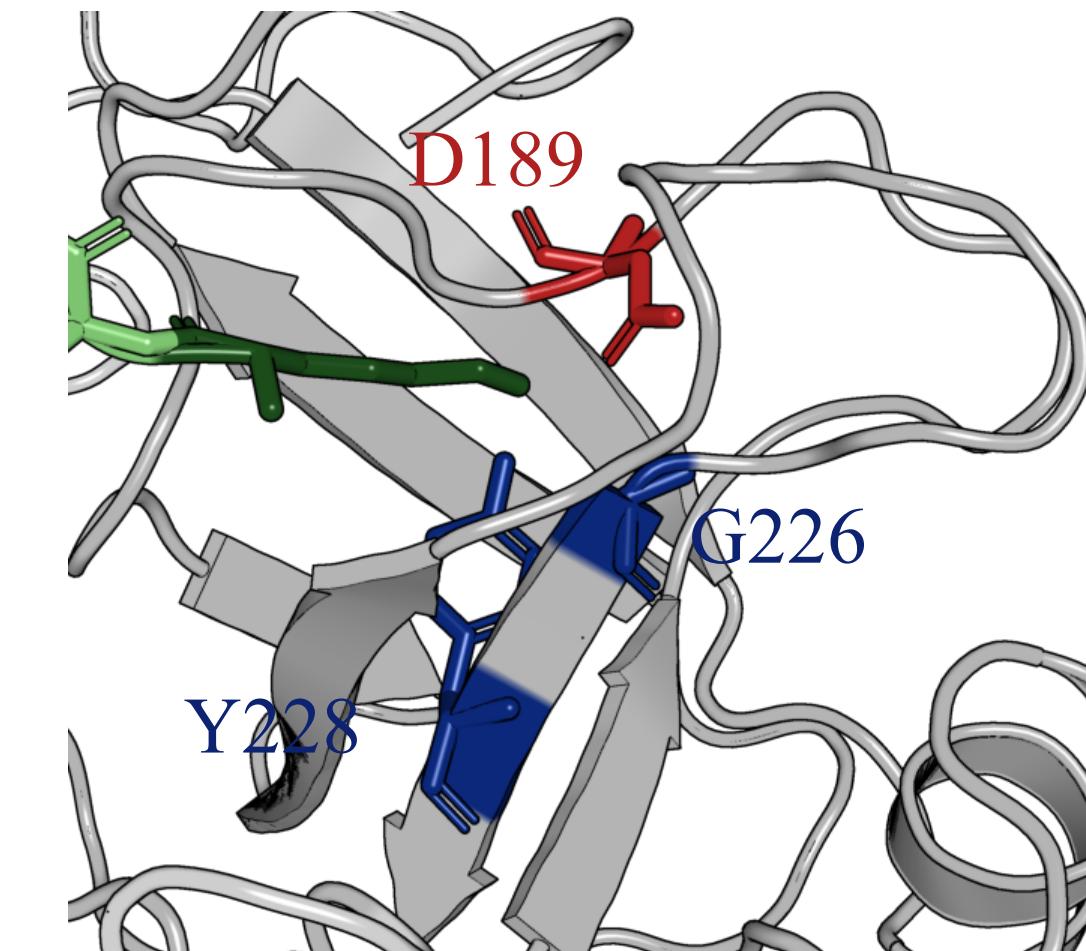


## Preliminary experimental results

### Chymotrypsin activity



### Zoom on the binding pocket



⚠ Not possible to make quantitative comparisons here

- ▶ Evidence of chymotrypsin activity 3 mutations away from rat trypsin
- ▶ Mutations located in the binding pocket

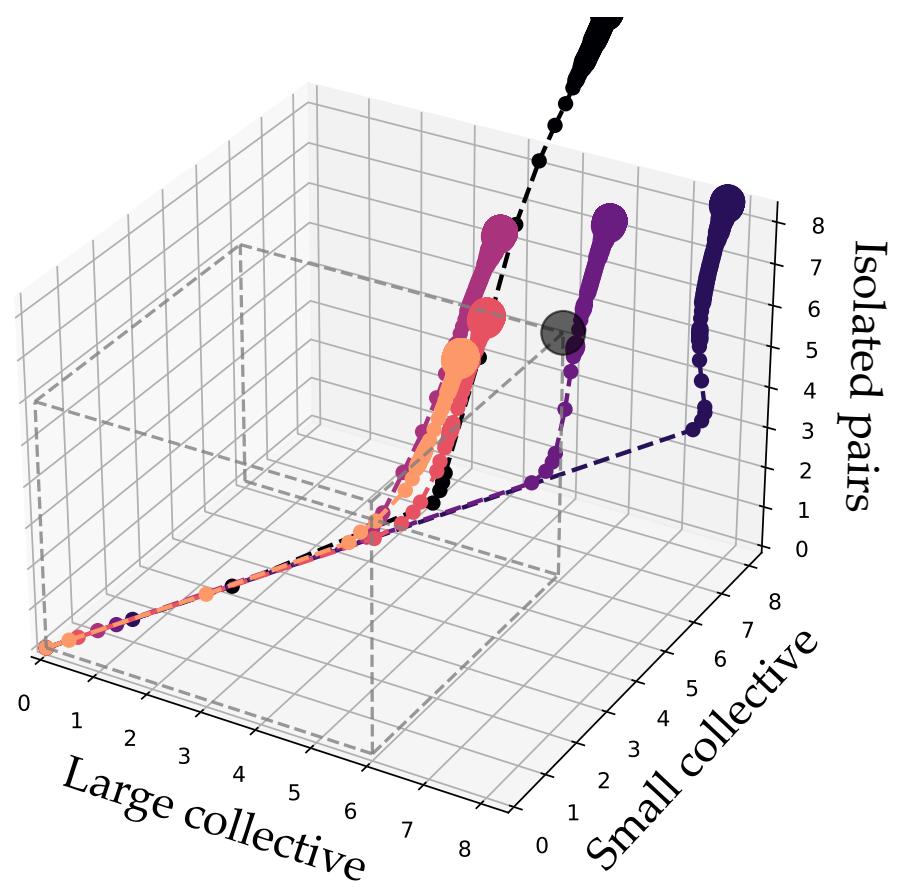
## Other results & next steps:

- ▶ More precise characterization of functional mutants
- ▶ Other conversions

# The undersampling problem

How to infer rich statistical  
structure from limited data?

Stochastic Boltzmann Machine



**Toy model:** correct  
undersampling-induced biases

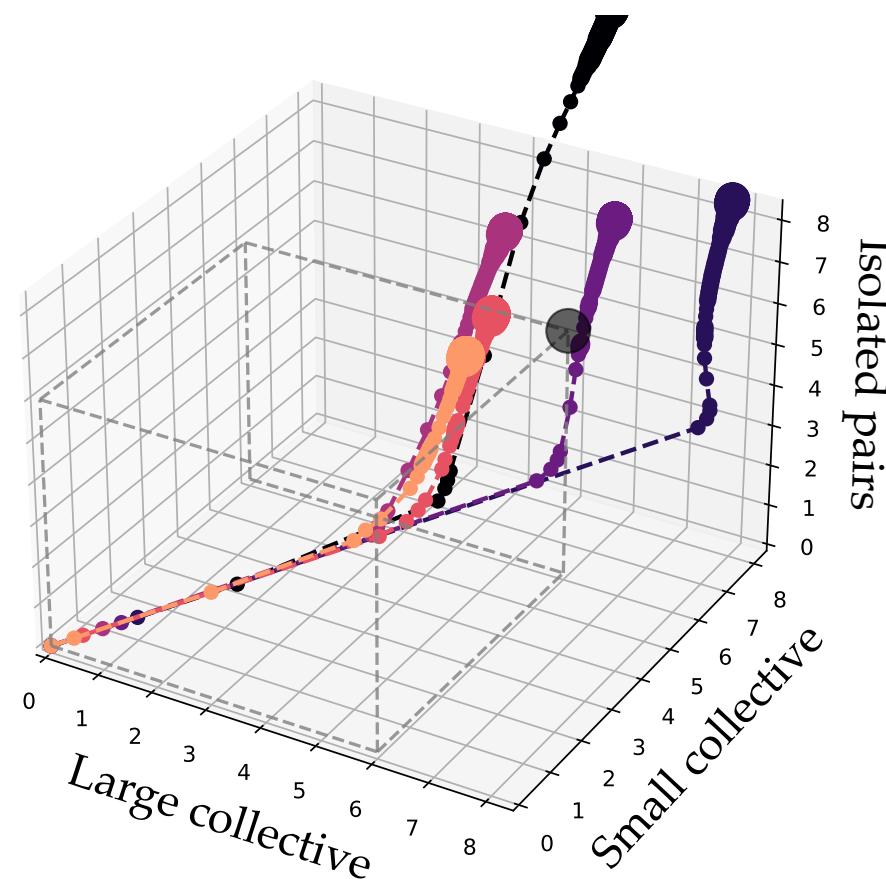
**Real data:** model generative  
without modification

Theoretical understanding?

# The undersampling problem

How to infer rich statistical structure from limited data?

Stochastic Boltzmann Machine



**Toy model:** correct undersampling-induced biases

**Real data:** model generative without modification

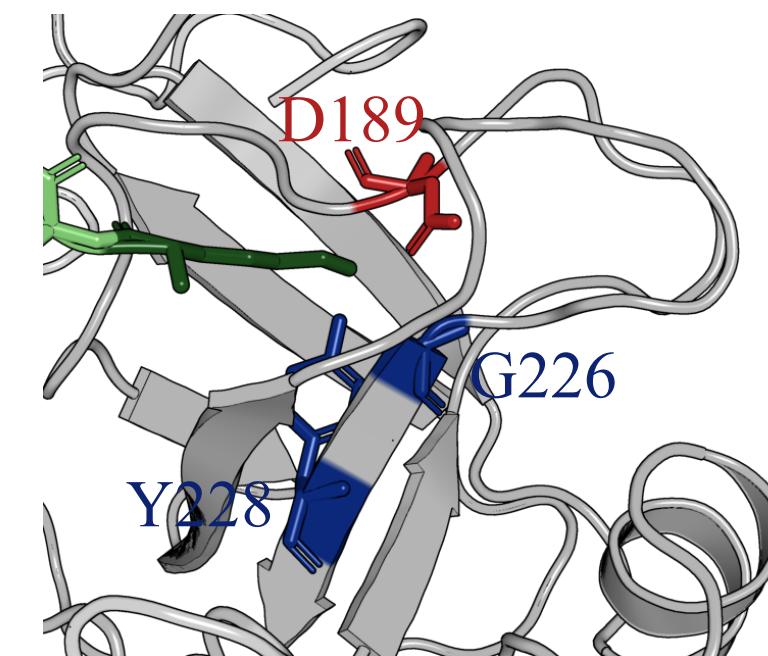
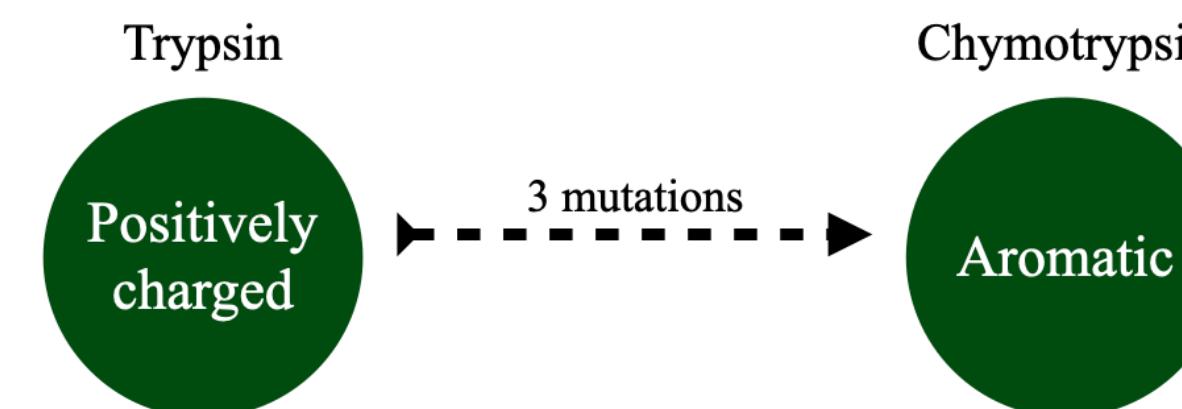
Theoretical understanding?

In collaboration with Emily Hinds, Yaakov Kleeorin, Rama Ranganathan (University of Chicago, USA)

# Investigate protein properties with statistical learning

Specificity conversion in S1A family?

Predict compensatory mutations with Boltzmann Machine



Evidence of chymotrypsin activity 3 mutations away from rat trypsin

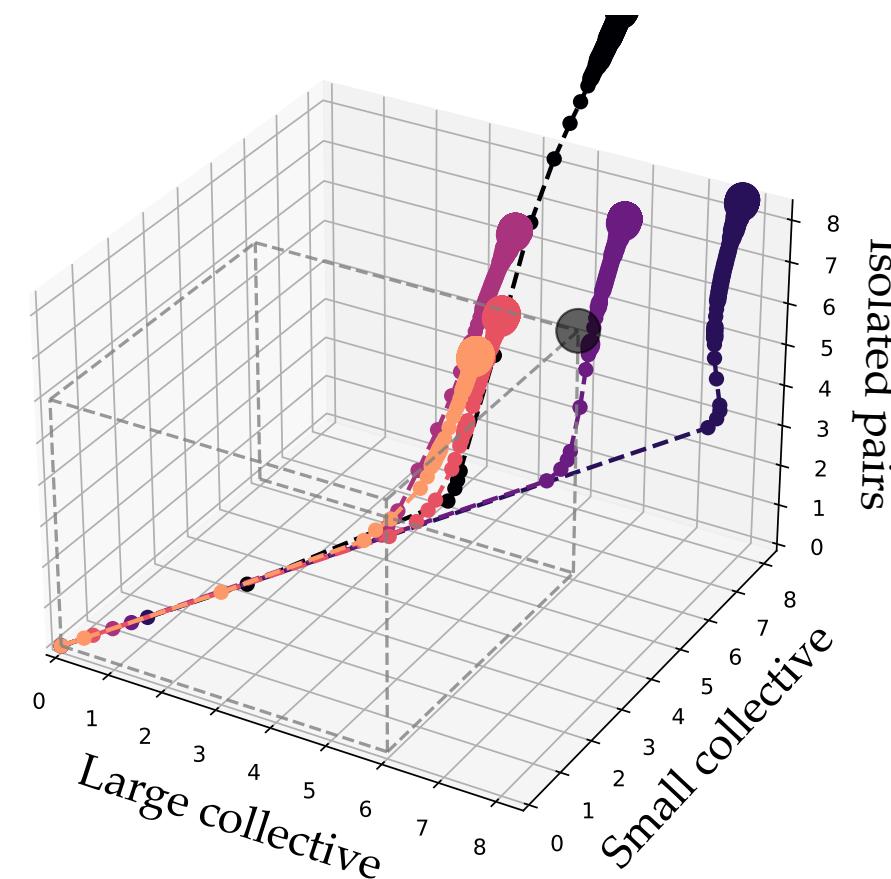
Better experimental characterization

In collaboration with Amaury Paveyranne, Timothé Lucas, Shoichi Yip, Clément Nizak (LJP, Sorbonne University, France)

# The undersampling problem

How to infer rich statistical structure from limited data?

Stochastic Boltzmann Machine



**Toy model:** correct undersampling-induced biases

**Real data:** model generative without modification

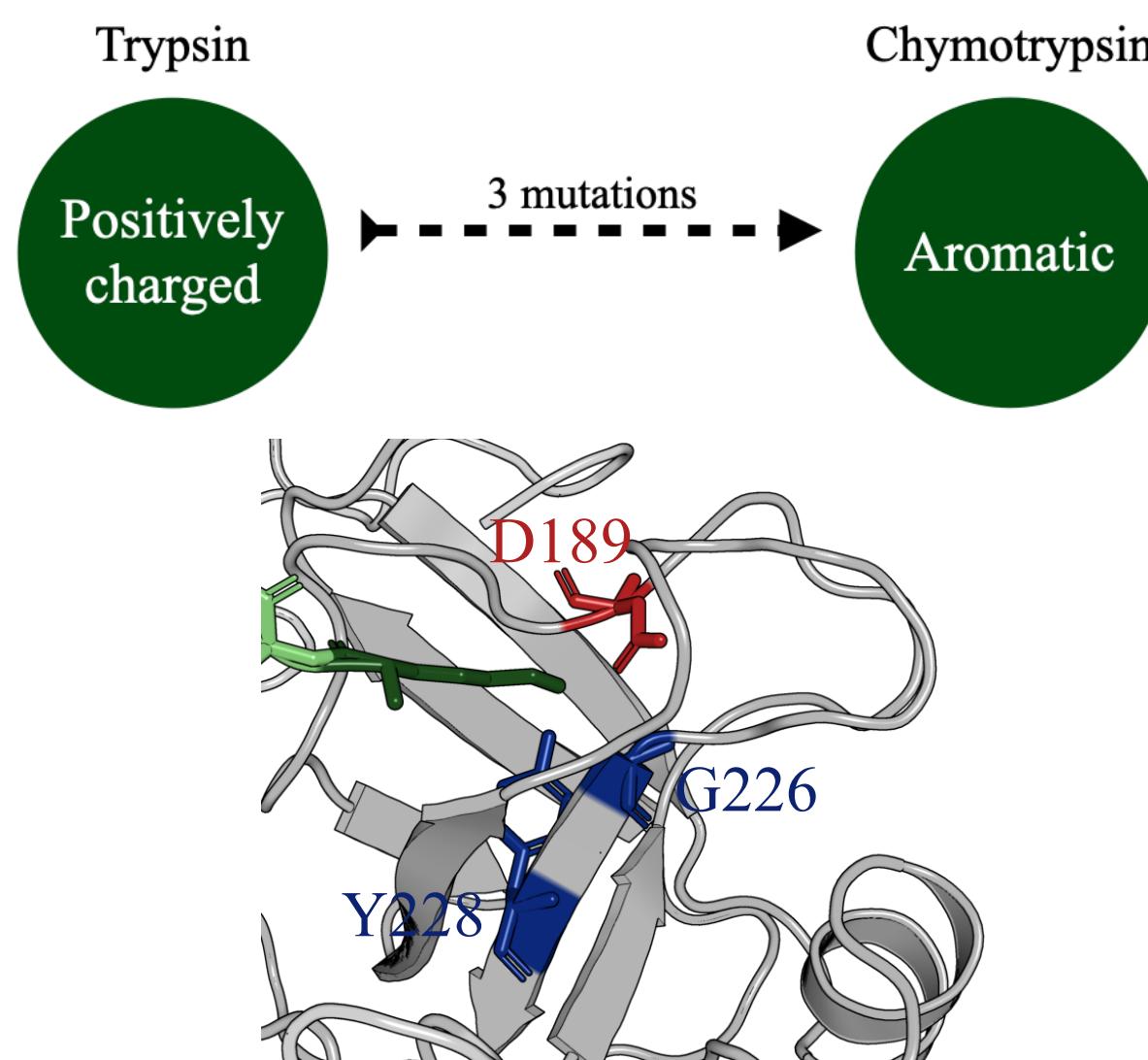
Theoretical understanding?

In collaboration with Emily Hinds, Yaakov Kleeorin, Rama Ranganathan (University of Chicago, USA)

# Investigate protein properties with statistical learning

Specificity conversion in S1A family?

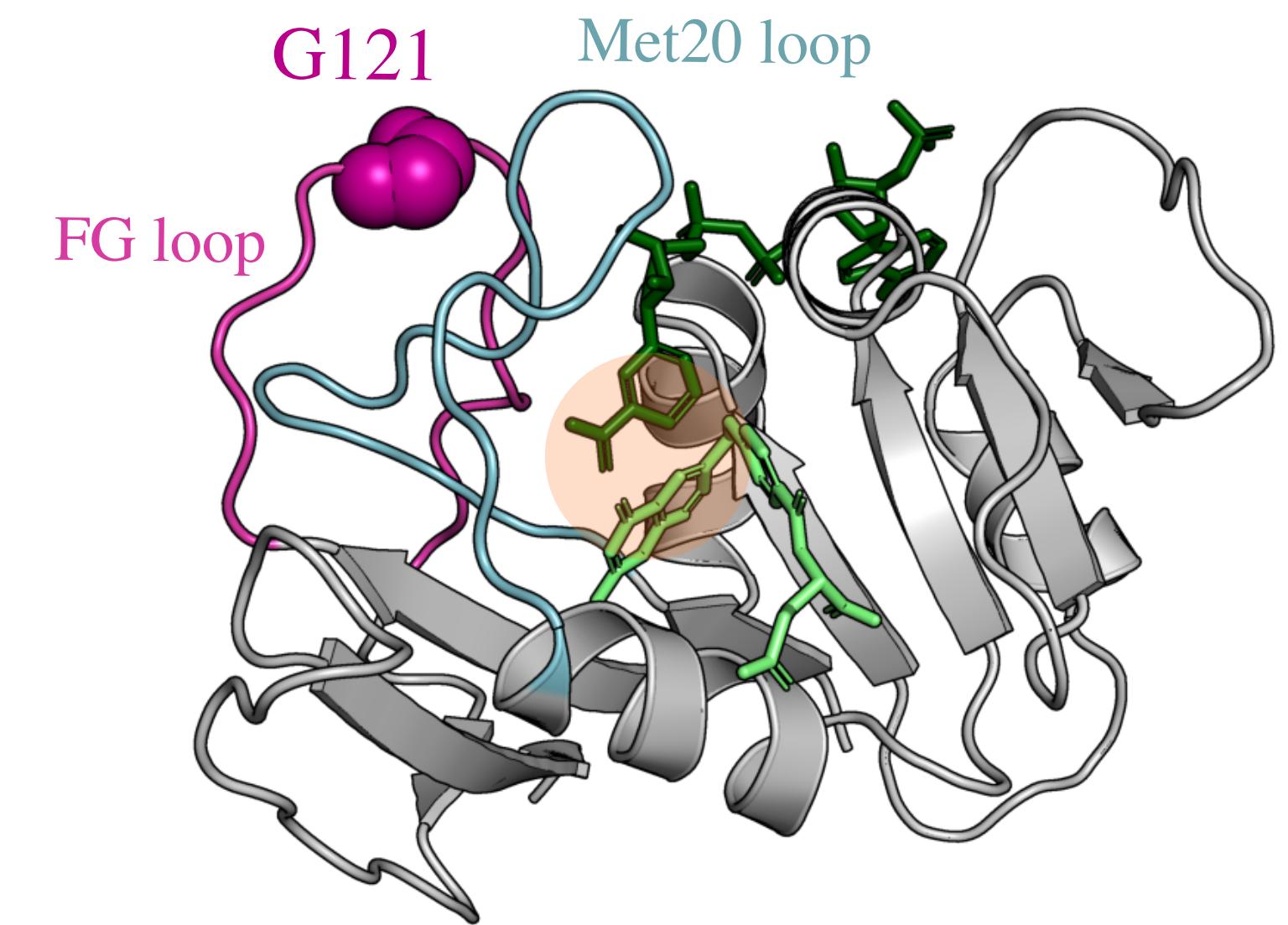
Predict compensatory mutations with Boltzmann Machine



Evidence of chymotrypsin activity  
3 mutations away from rat trypsin  
Better experimental characterization

In collaboration with Amaury Paveyranne, Timothé Lucas, Shoichi Yip, Clément Nizak (LJP, Sorbonne University, France)

Allosteric network of *E. Coli* DHFR



Support for an allosteric mechanism obtained through molecular dynamics simulations

In collaboration with Paul Guenon, Damien Laage, Guillaume Stirnemann (ENS, France), Clément Nizak (LJP, France), Karolina Filipowska, Kim Reynolds (University of Texas, USA)

# Collaborators

Rama Ranganathan

Emily Hinds

Nadav Benhamou Goldfajn

Miranda Moore

Yaakov Kleeorin

Amaury Paveyranne

Timothé Lucas

Shoichi Yip

Clément Nizak

Paul Guenon

Damien Laage

Guillaume Stirnemann

Karolina Filipowska

Kim Reynolds

# Nantes & co

Ma famille, mes amis

Les p'tits lapins

Les colocs

Snoopy, Sido, Pinpin

# Directeurs de thèse

Ivan Junier

Olivier Rivoire

# Jury

Anne-Florence Bitbol

Cyril Furtlehner

Thierry Mora

Marco Ribezzi

# Paris

Le laboratoire Gulliver

Le bureau 414 du LJP

# Grenoble

Le laboratoire TIMC

L'équipe TrEE

# Chicago

Ranganathan Lab'

# Merci

