

# Overcoming Undersampling in Generative Models for Biological Sequences

Marion Chauveau<sup>1,3</sup>, Yaakov Kleeorin<sup>2</sup>, Ivan Junier<sup>1</sup>, Olivier Rivoire<sup>3</sup>

<sup>1</sup>TIMC, Univ. Grenoble Alpes, <sup>2</sup>University of Chicago, <sup>3</sup>Gulliver Laboratory, ESPCI Paris

## Generative Models

### Fundamental problems

Undersampling: Inference relies on limited datasets [2]

Model Evaluation: How to score generative power ?

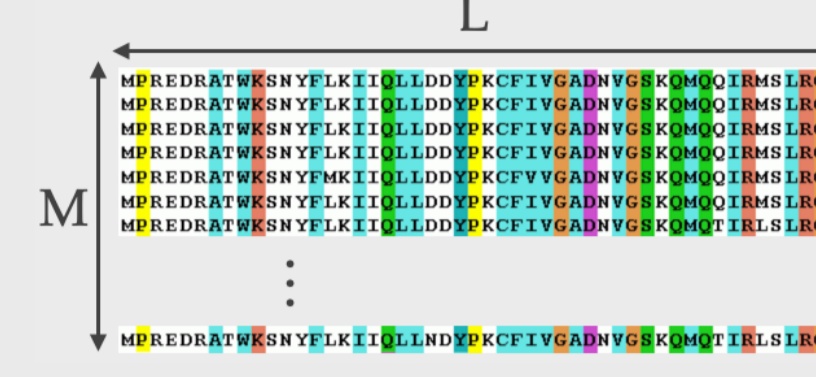
### Current approach

L2-Regularization

Comparison of empirical statistics with model statistics

## Potts Model

Maximum entropy model trained to match the empirical one and two-body frequencies of amino acids [1]

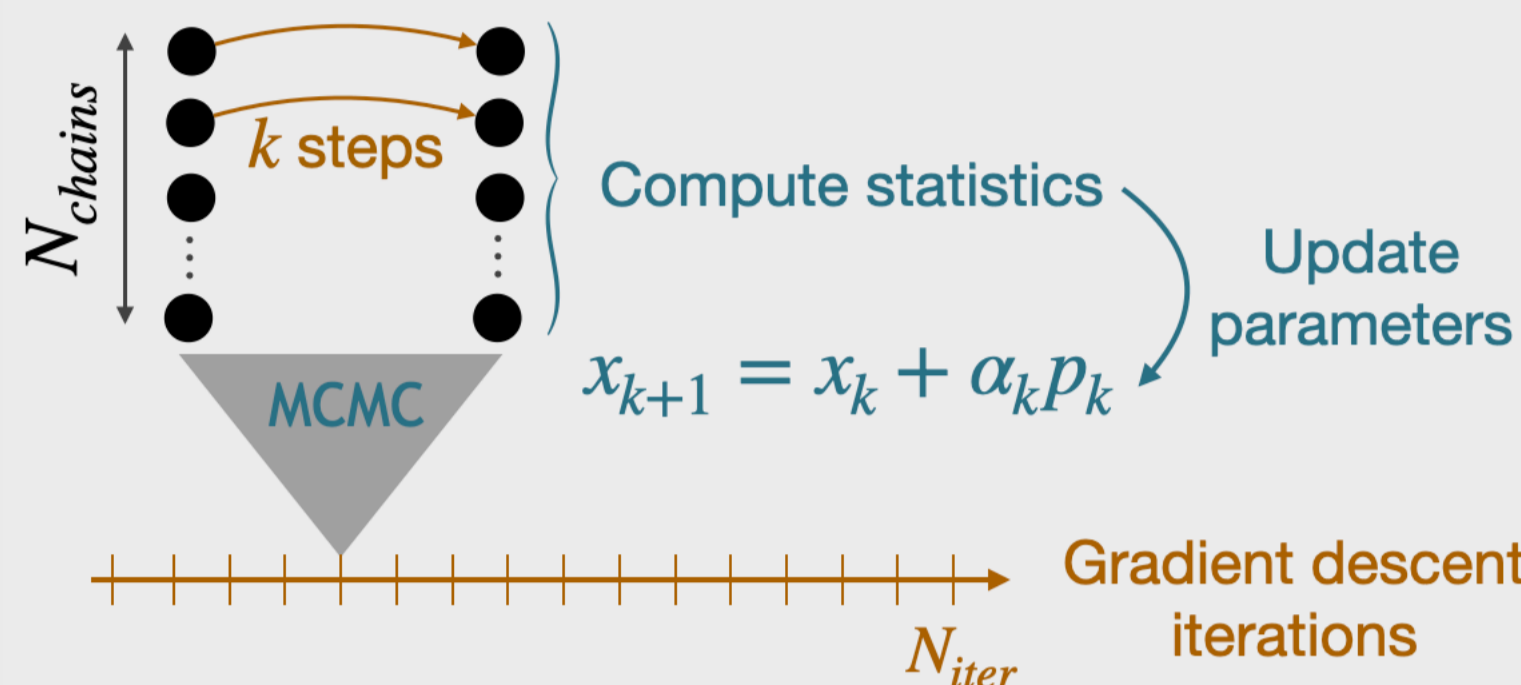


Protein:  $\sigma = \{\sigma_1, \dots, \sigma_L\}$

$\sigma_i$  = a.a. at site  $i$  ( $q$  possibilities)

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

## Methods



### Problems with BM & L2-regul. :

- Energy discrepancy between natural & artificial sequences
- Bias in the inference of heterogeneous couplings
- Computationally slow

### Optimization Method

### Hyperparameters

### Regularization

### Current: Boltzmann Machine (BM) [1]

Steepest Gradient Descent

$$p_k = -\nabla f_k$$

Log-Likelihood

### New: Stochastic Boltzmann Machine (SBM)

L-BFGS Gradient Descent [4]

$$p_k = -B_k^{-1} \nabla f_k$$

Low rank approximation of the Hessian

Nb of Gradient Descent iterations:  $N_{iter}$   
Nb of MCMC steps and chains:  $k, N_{chains}$

L2-regularization:  $\lambda_h, \lambda_J$

Rank of the approximated Hessian:  $m$

$(\lambda_h, \lambda_J)$

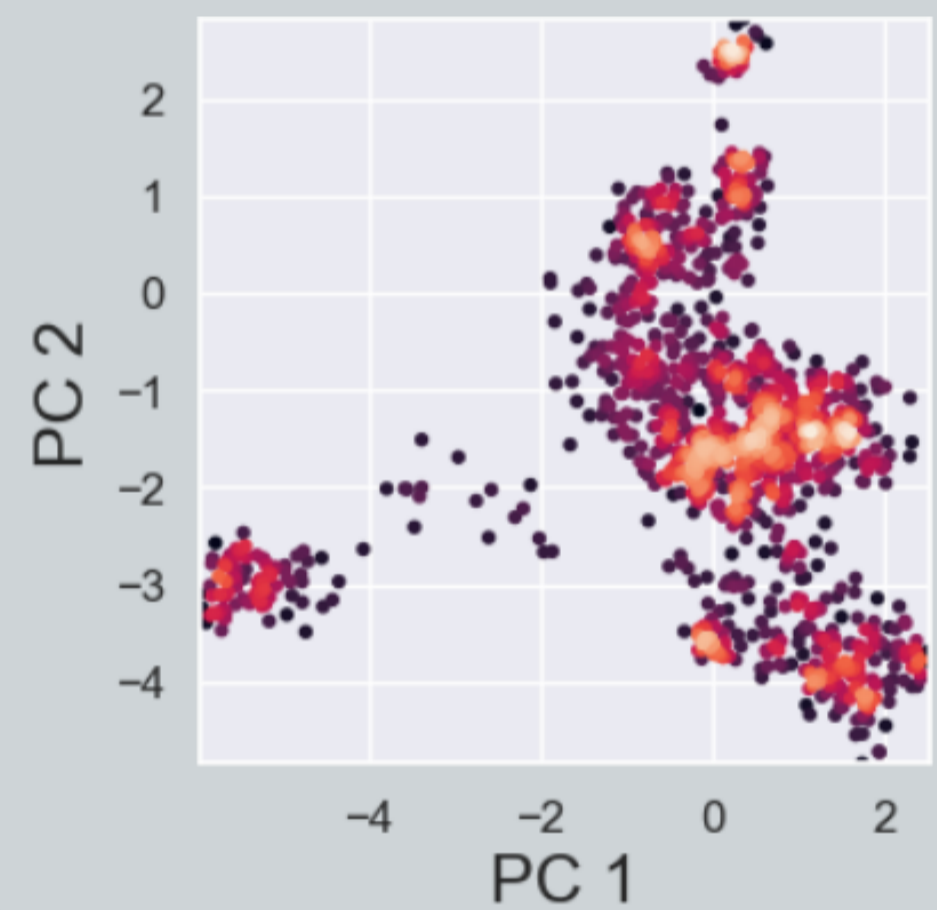
$(m, N_{chains}, N_{iter})$

## Energy discrepancy between natural & artificial sequences

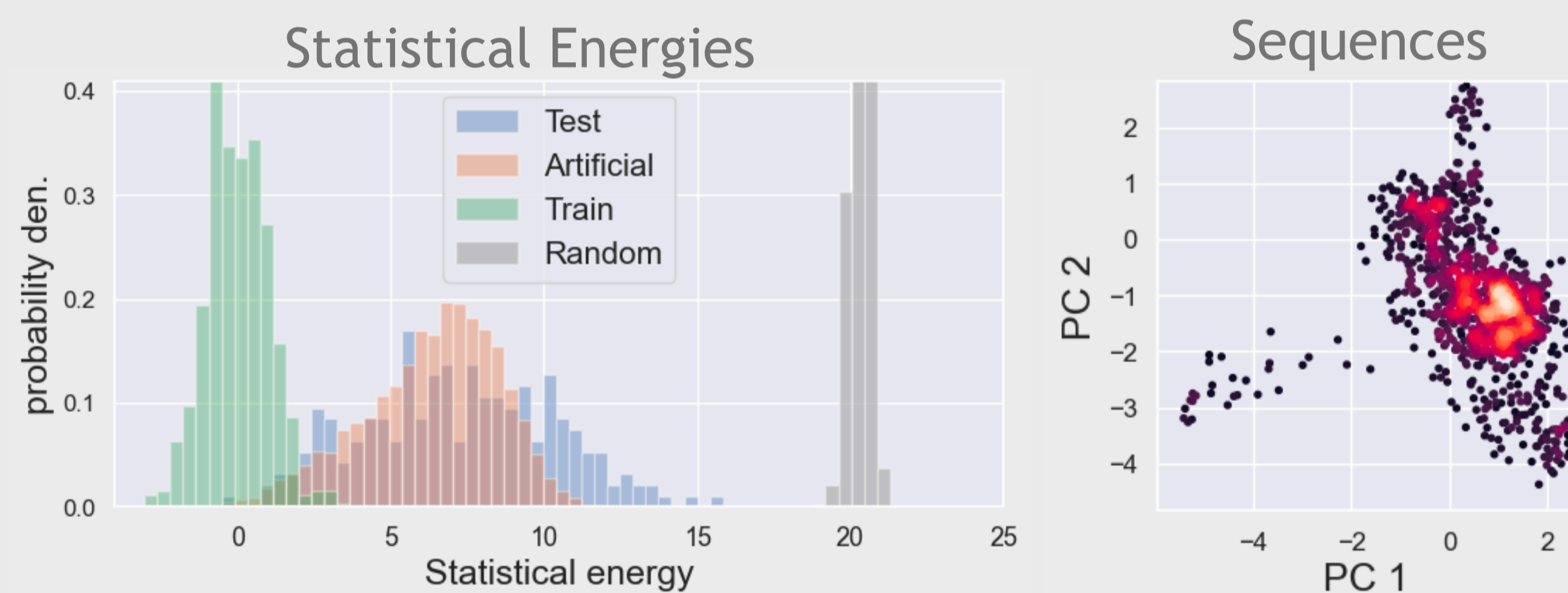
### Chorismate Mutase family [3]

$M_{train} = 904, M_{test} = 226$

$L = 96$

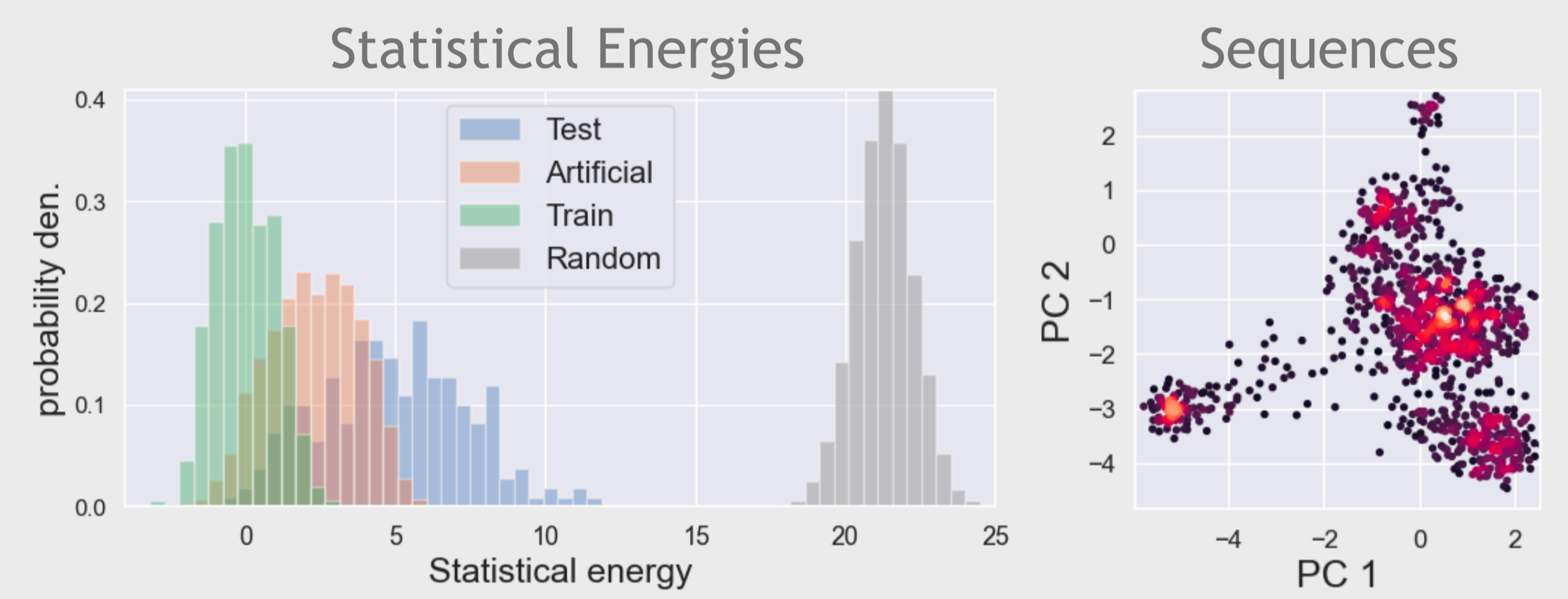


### BM ( $\lambda = 0.01, T = 1$ )



- ➔ L2-regularization of BM leads to an energy discrepancy between natural & artificial sequences
- ➔ PCA shows a loss of diversity in the artificial sequences

### SBM ( $N_{chains} = 40, m = 1, N_{iter} = 500, T = 1$ )



- ➔ Energy discrepancy between natural and artificial sequences is smaller than for the BM method
- ➔ PCA shows a greater diversity than for the BM method

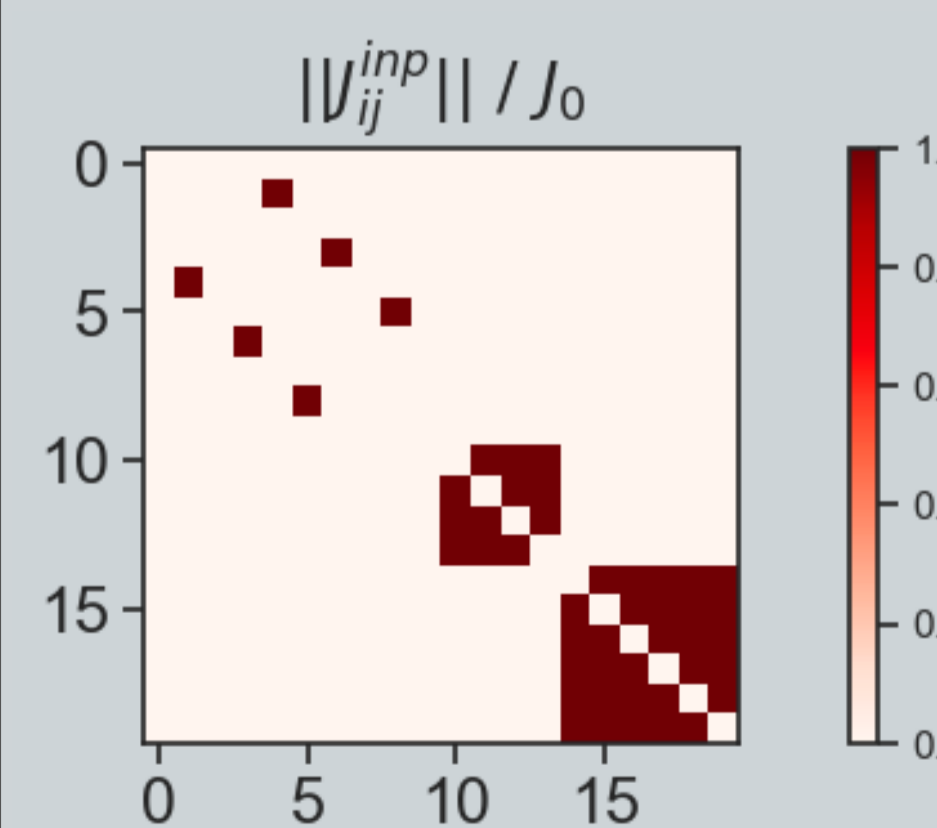
## Bias in the inference of heterogeneous couplings

### Toy Model [2]

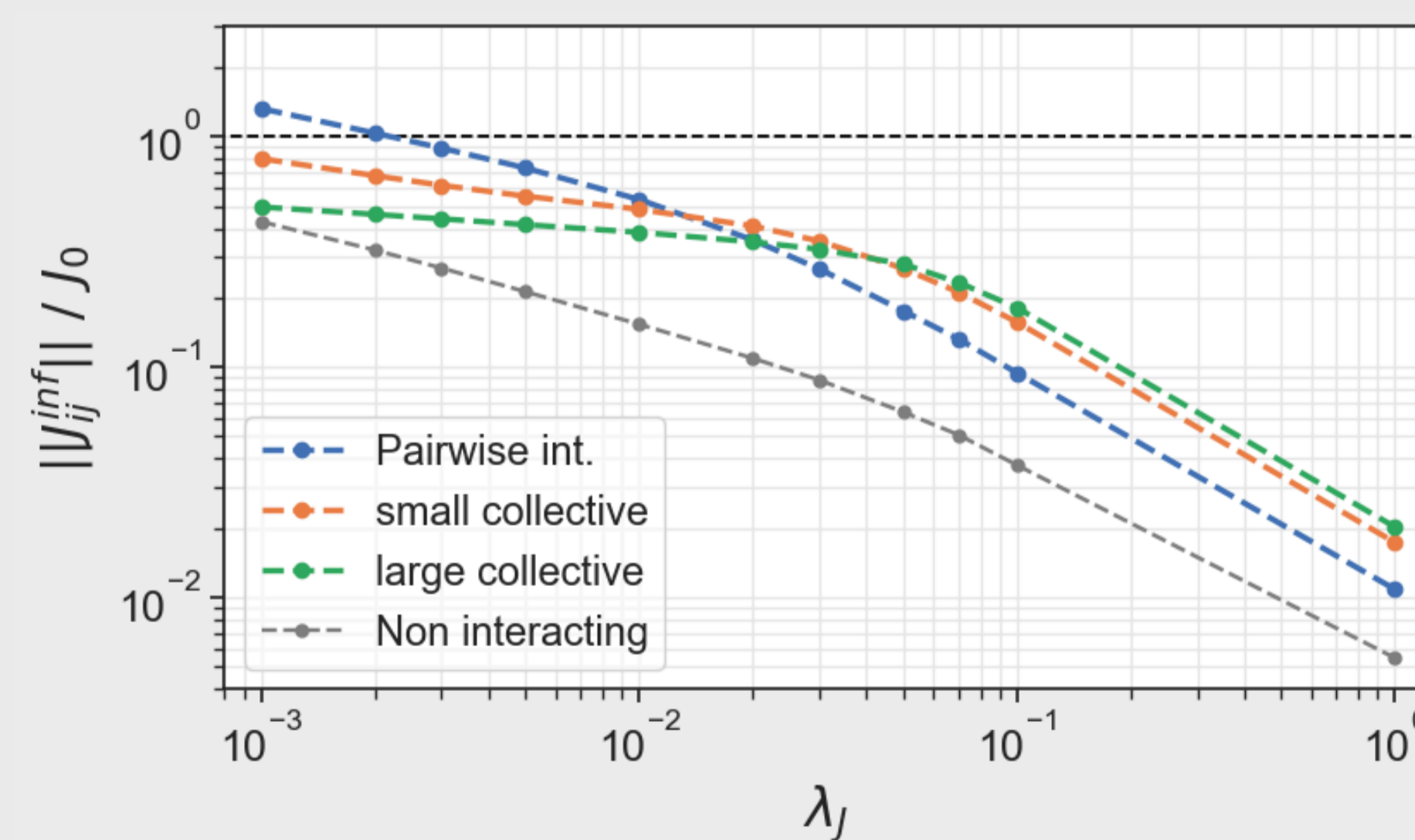
$q = 10, L = 20$

$J_0 = 2, h_i = 0 \forall i$

$M = 300$

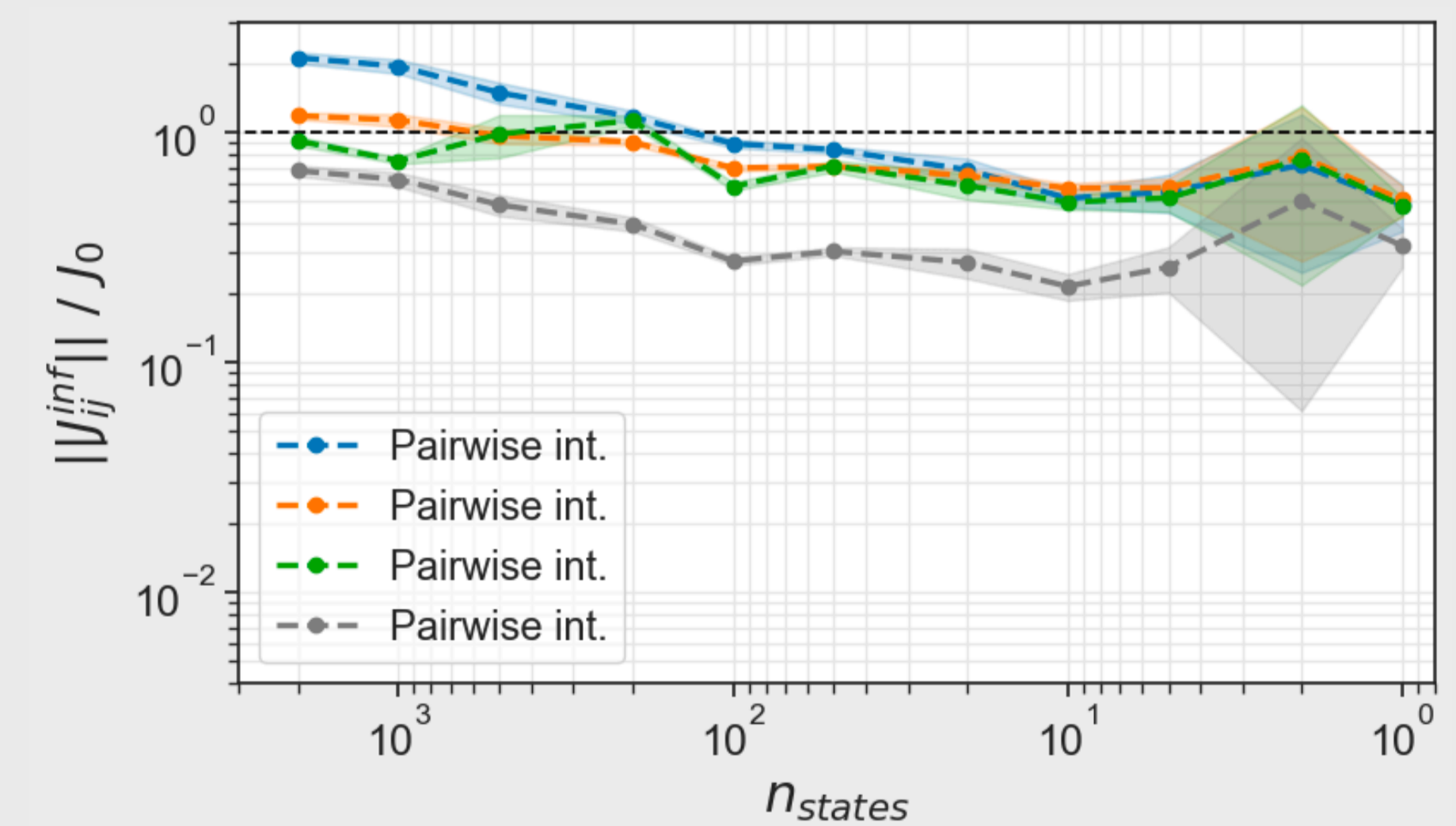


### BM ( $\lambda_h = 0.01, N_{chains} = 1000$ )



- $\lambda_J \sim 10^{-3}$ : large collective features have lower magnitude
- ➔ L2-regularization of BM leads to a bias in the coupling inference

### SBM ( $N_{iter} = 500, m = 1$ )



- $n_{states} \sim 200$ : features well inferred with same magnitude
- ➔ Implicit regularization of SBM does not lead to a bias in the coupling inference

## Discussion & Future work

- Computation Time: SBM method converges much faster than the BM method
- Open question: How to score without comparing to the true model?

Use this approach to infer Generative models for Bacterial & Archaeal genomes

## References

- [1] M. Figliuzzi, P. Barrat-Charlaix, and M. Weigt, 2018.
- [2] Y. Kleeorin, W. P. Russ, O. Rivoire, and R. Ranganathan, 2021.
- [3] W. P. Russ, M. Figliuzzi, C. Stocker, P. Barrat-Charlaix, M. Socolich, P. Kast, D. Hilvert, R. Monasson, S. Cocco, M. Weigt et al., 2020.
- [4] D. C. Liu and J. Nocedal, 1989.