

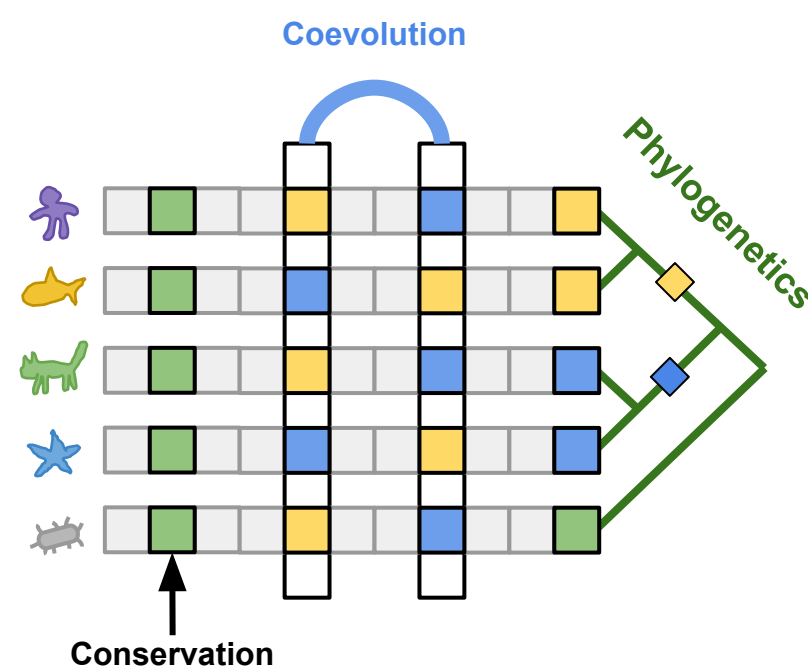
# Unified framework for modeling multivariate distributions in biological sequences

Justas Dauparas<sup>1</sup>, Haobo Wang<sup>1</sup>, Avi Swartz<sup>2</sup>, Peter Koo<sup>3</sup>, Mor Nitzan<sup>4</sup>, Sergey Ovchinnikov<sup>4</sup>

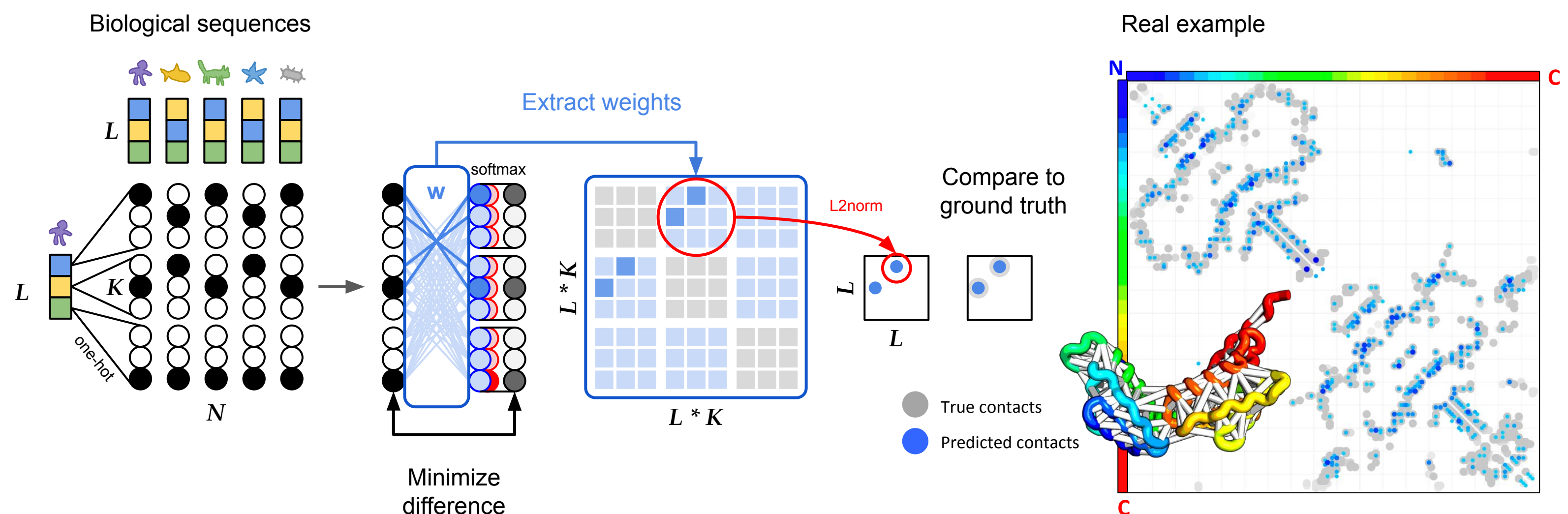
<sup>1</sup>FAS Division of Science <sup>2</sup>Harvard College <sup>3</sup>Howard Hughes Medical Institute <sup>4</sup>JHDSF Program, Harvard University

## Introduction

Revealing evolutionary **conserved**, **co-evolving** sites and **phylogenetics** (relationship between sequences) is a fundamental and challenging task in the protein field.



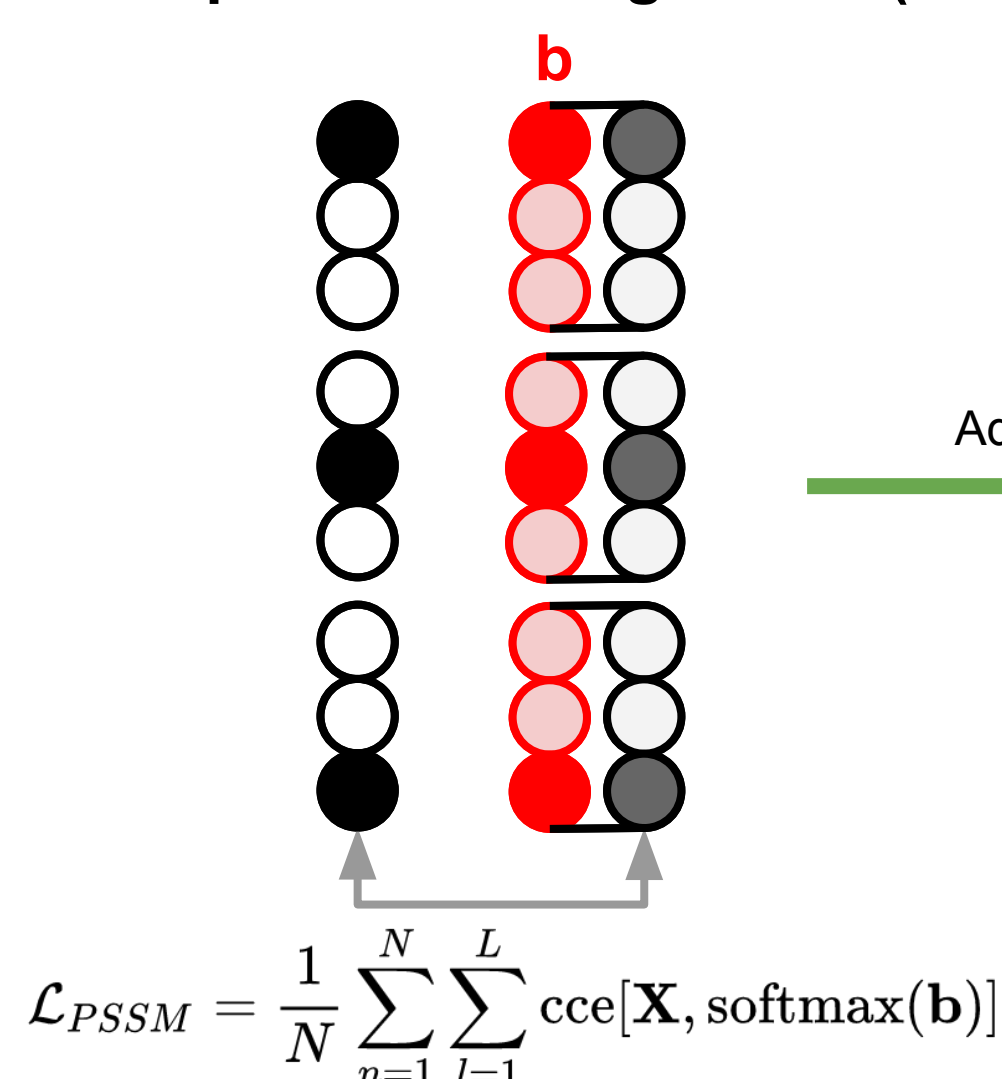
**Coevolution** in 1D biological sequences, captured by **weights** of a fully-connected (or dense) layer, is correlated to physical contacts in 3D structures!



## Methods

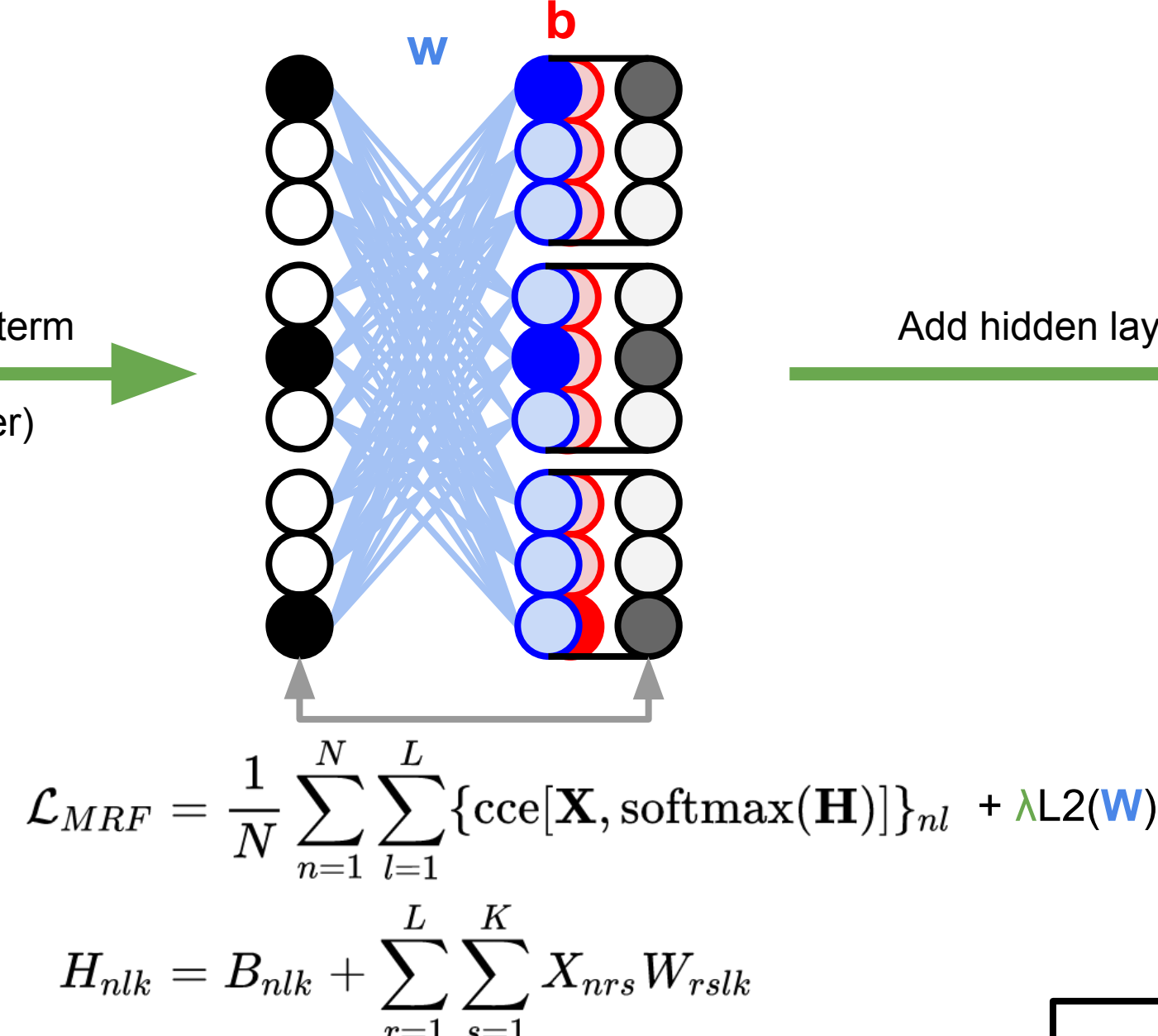
Here we propose a unified framework for the current models (Position-Specific Scoring Matrices, Markov Random Fields, Multivariate Gaussian models, and Autoencoders), that allows for interpretable transformations between these methods and naturally incorporates the advantages and insight gained individually in the different communities.

### Position-Specific Scoring Matrix (PSSM)

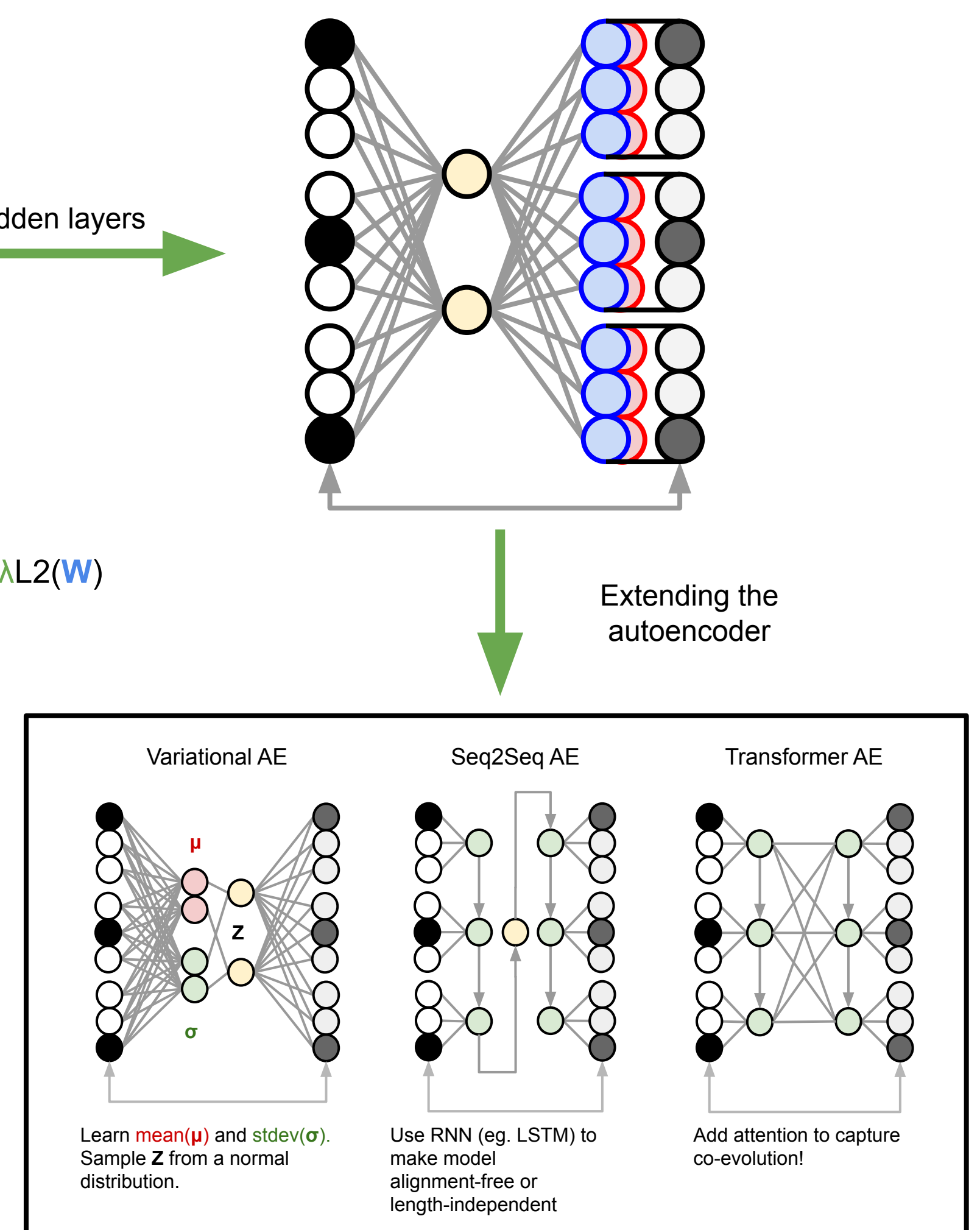


### Markov Random Field (MRF)

w/ pseudo-likelihood approx.

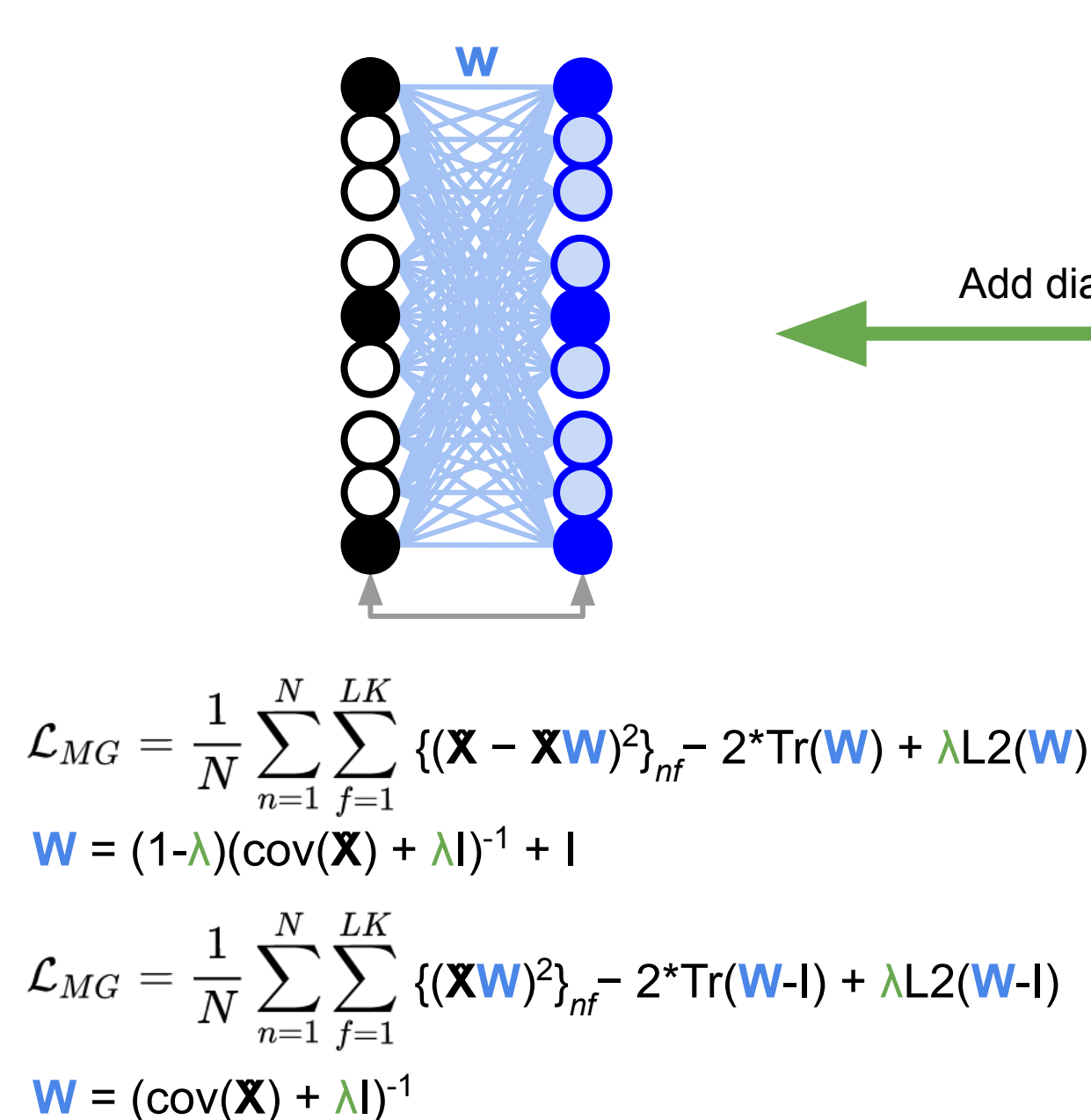


### Autoencoders (AE)

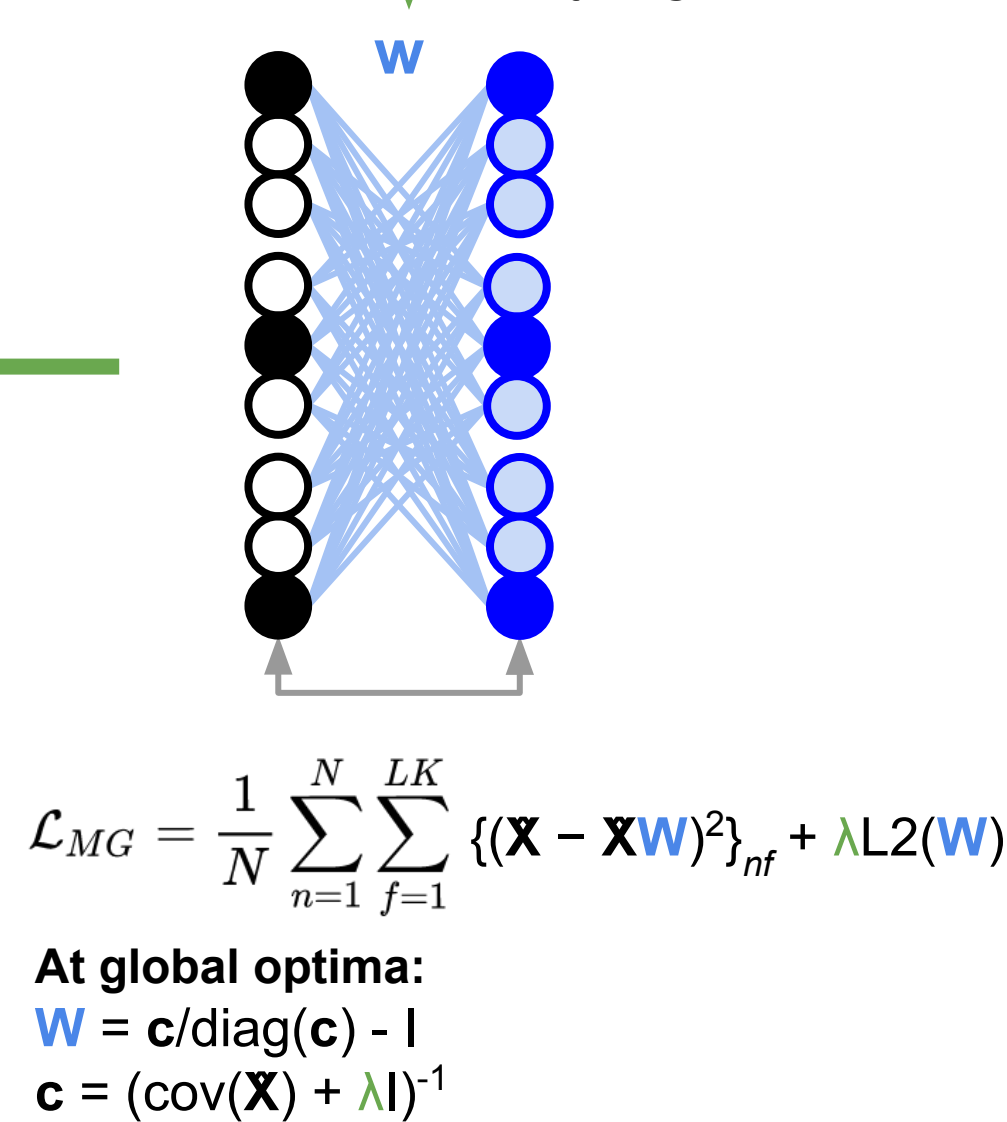


### Multivariate Gaussian (MG)

w/ mean-squared-error approx.



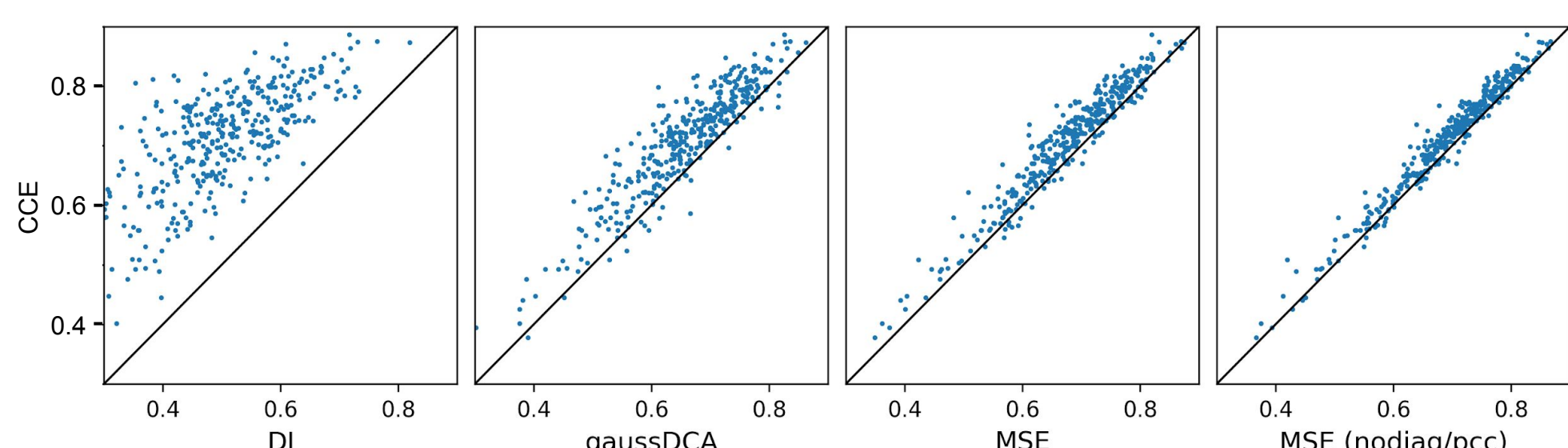
Replace softmax/CCE with MSE



- Interesting observations:**
- Pseudo-likelihood = Categorical cross-entropy (CCE)
  - Multivariate Gaussian = Mean square error (MSE)
  - L2 regularization = "shrinkage"
  - Zeroing diagonal = "Partial correlation coefficient"

## Results

### Contact prediction accuracy



For contact prediction, Categorical cross entropy (CCE) is significantly more accurate than traditional inverse-covariance methods (DI, GaussDCA). By applying L2 regularization to MSE, we find the inv-cov method to be almost as accurate as CCE. Each point corresponds to a different protein.

## Conclusion

- Widely used models for multivariate distributions in biological sequences can all be expressed as a single fully-connected layer, where the weights and bias of the dense layer captures the co-evolution and conservation, respectively.
- The differences in the models comes down to the loss function used, where inverse-covariance methods are effectively minimizing the mean-squared error (more appropriate for continuous data) and markov-random-field methods are minimizing the categorical cross entropy (more appropriate for categorical data).

### Emails:

**We are looking for collaborators/postdocs/students!**

Mor Nitzan mornitzan@fas.harvard.edu  
Sergey Ovchinnikov so@fas.harvard.edu

<http://github.com/sokrypton/seqmodels>

### Github code

