

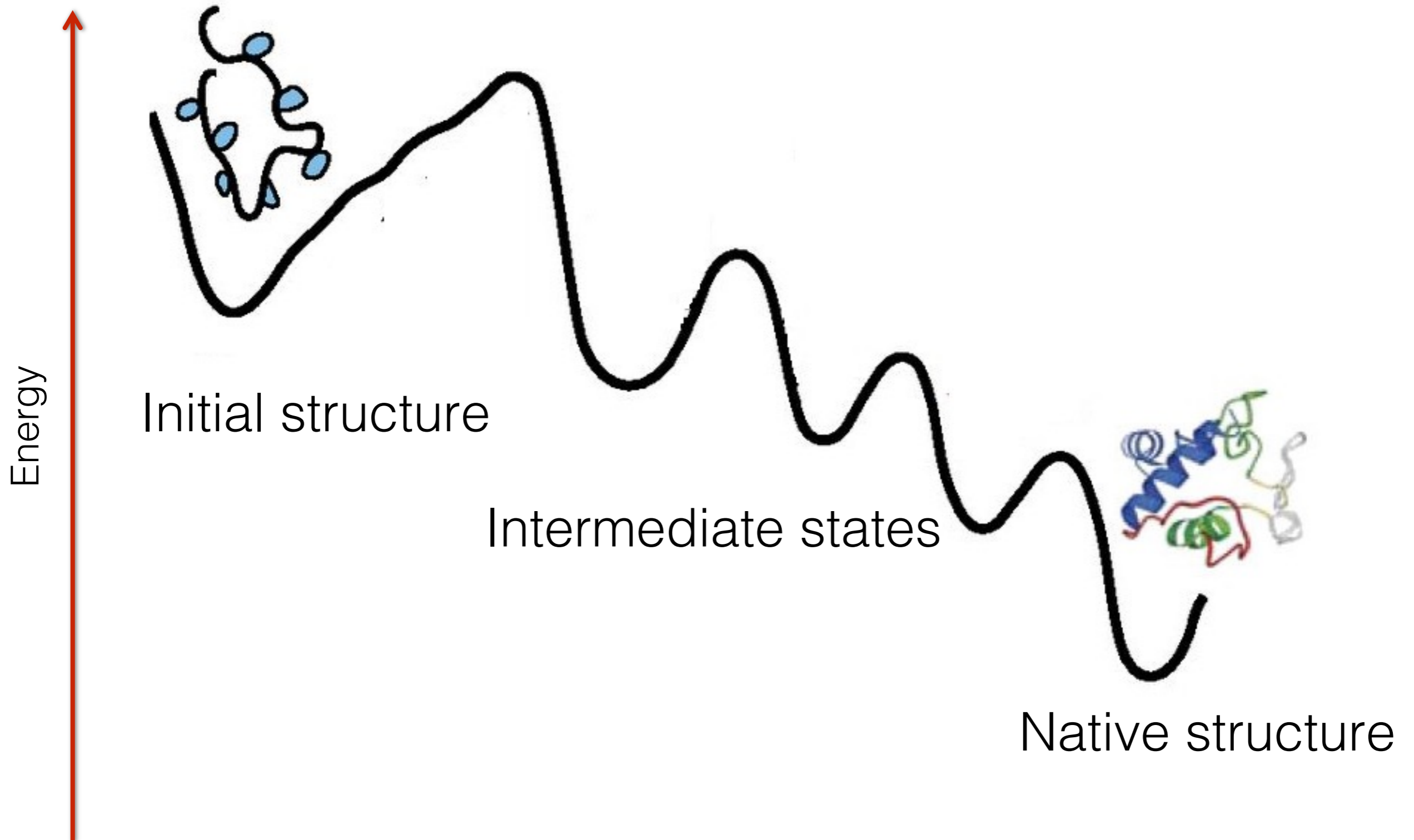


# Machine learning algorithms for structural and functional genomics

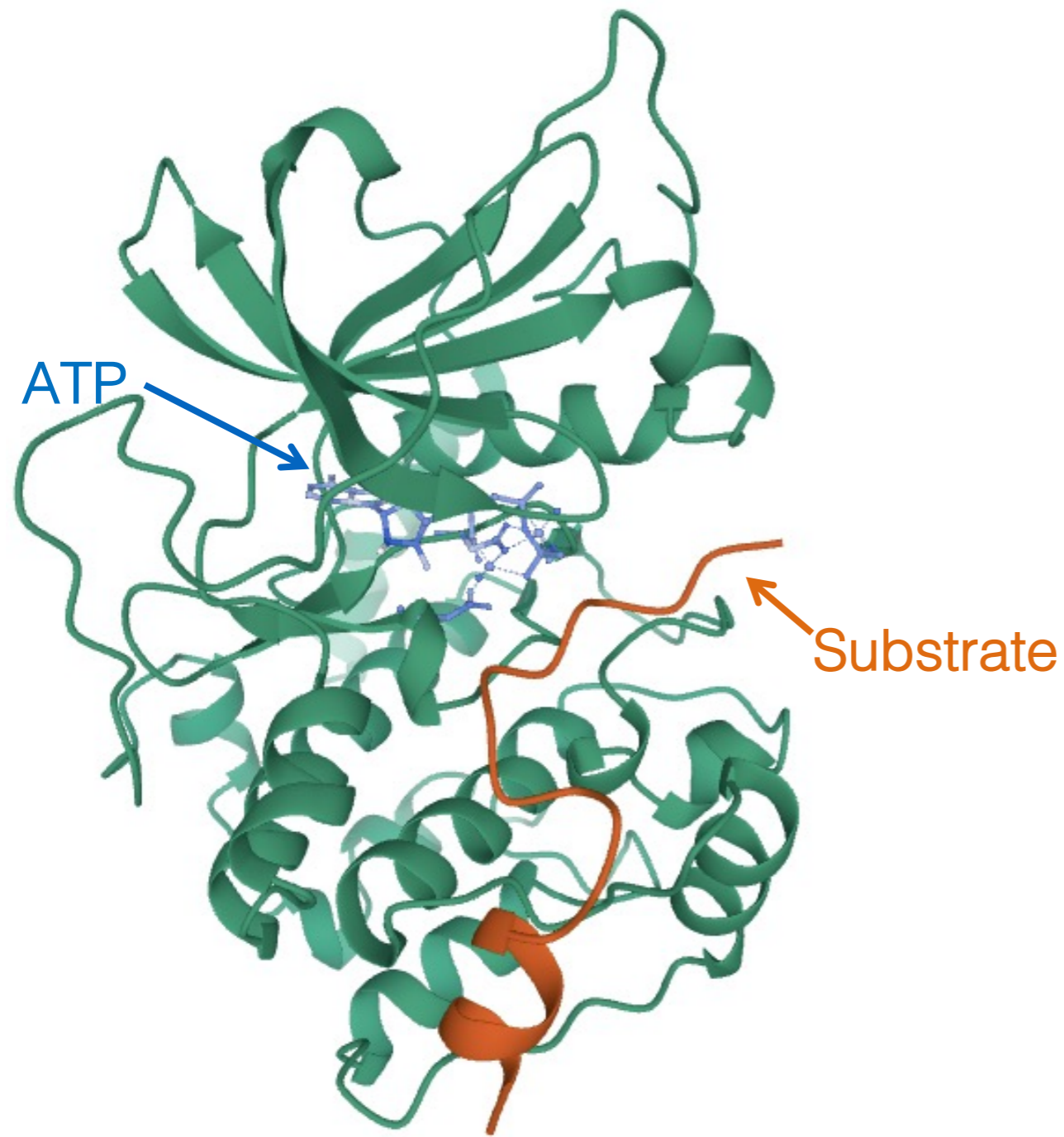
Jian Peng

Department of Computer Science  
College of Medicine, Institute of Genomic Biology  
University of Illinois at Urbana-Champaign

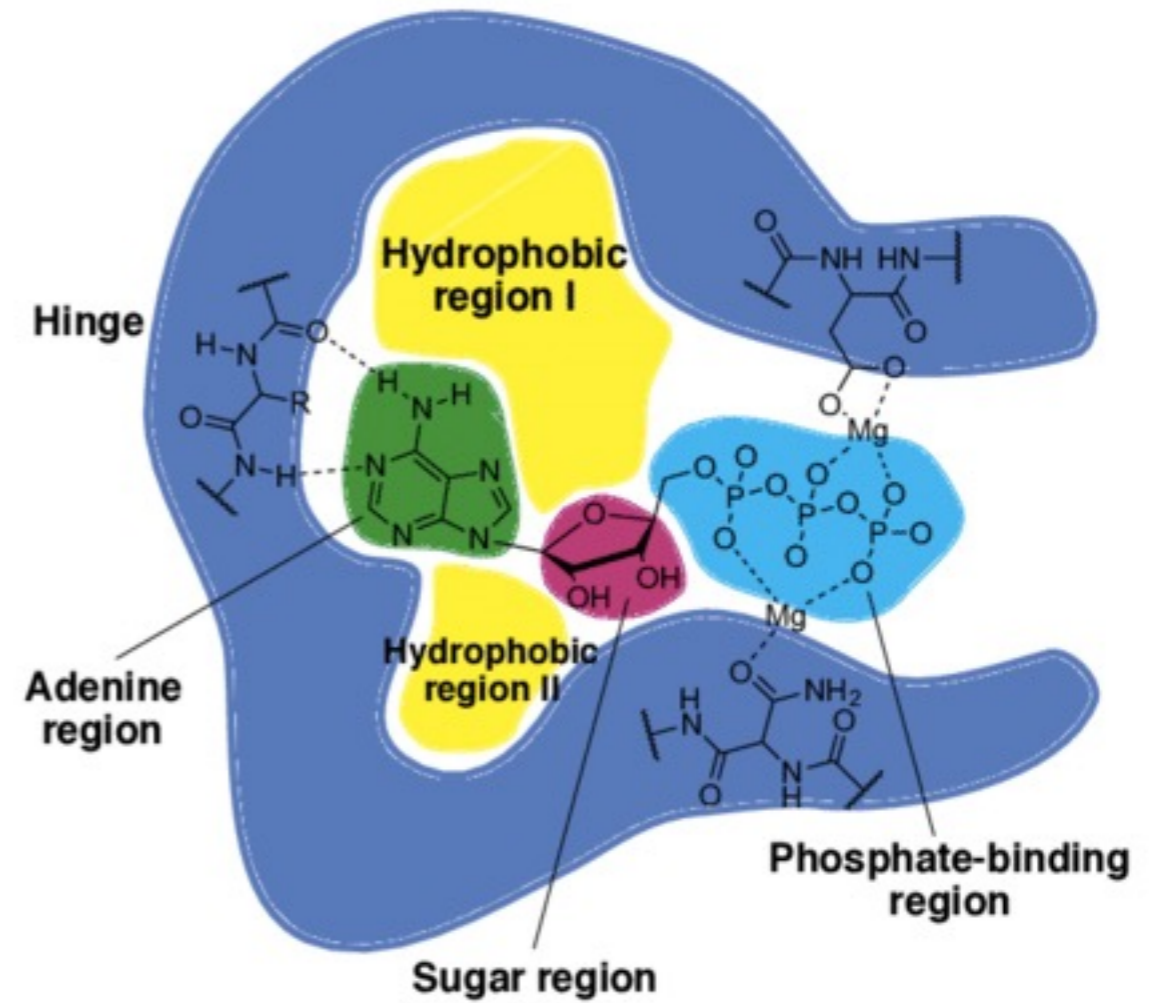
# The protein folding problem



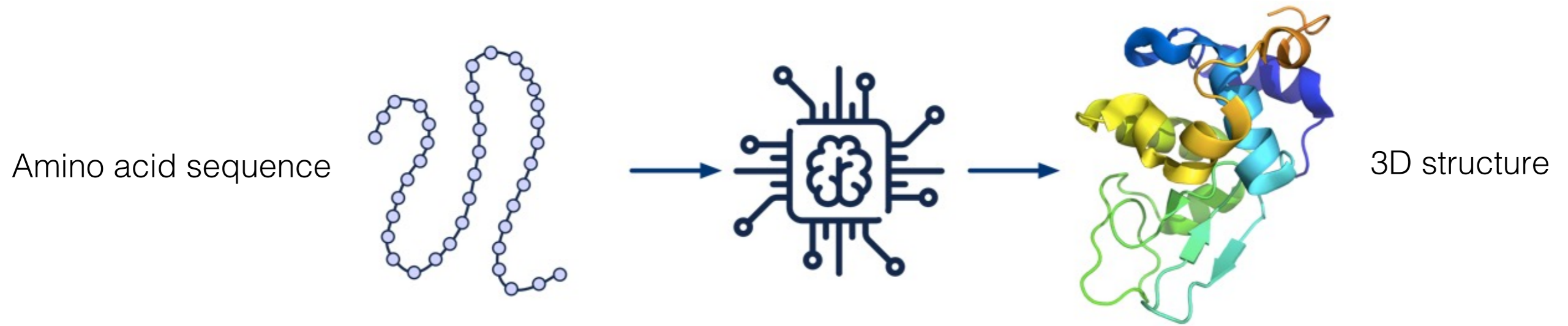
# Structure provides insights on function



Protein Kinase



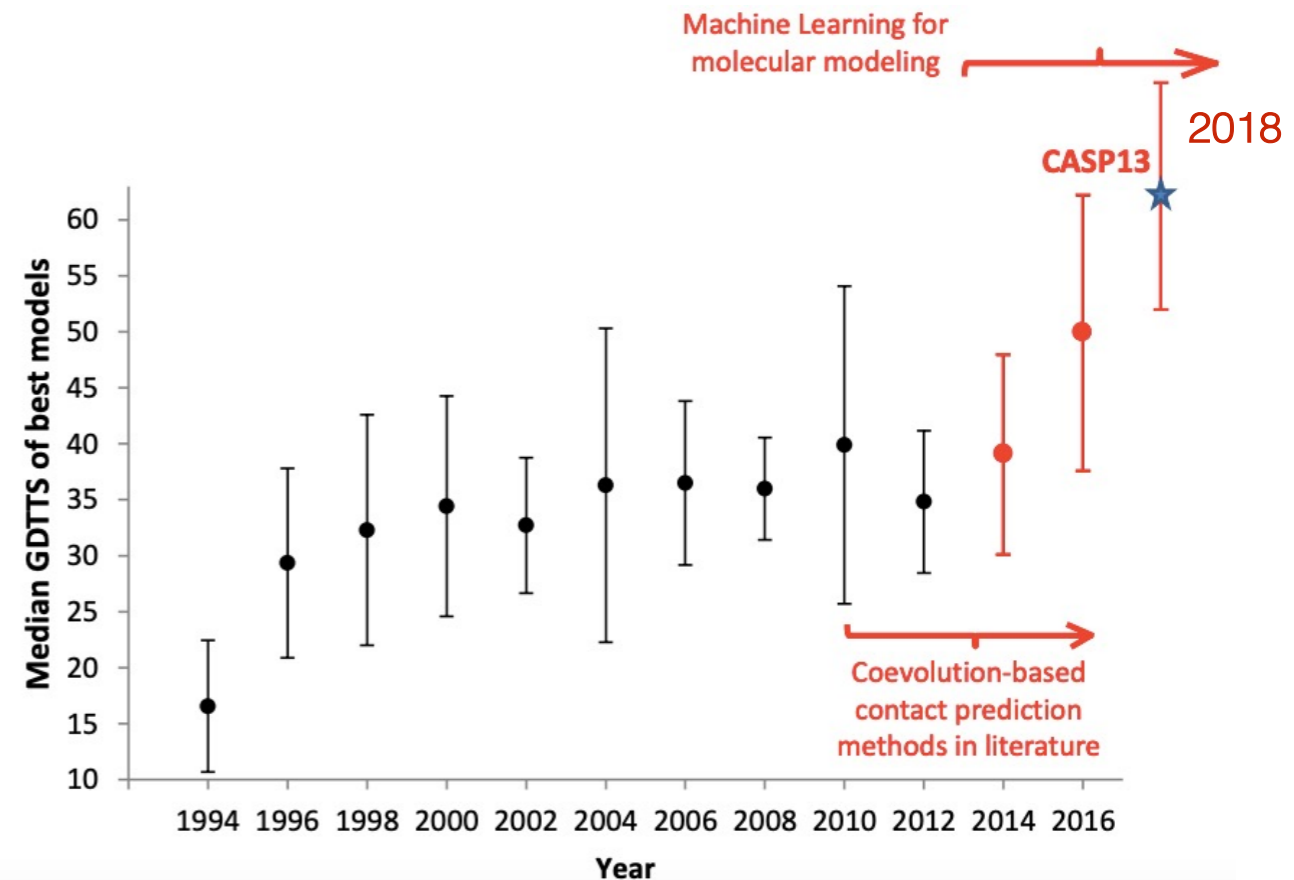
# Protein structure prediction



## Successful prediction algorithms



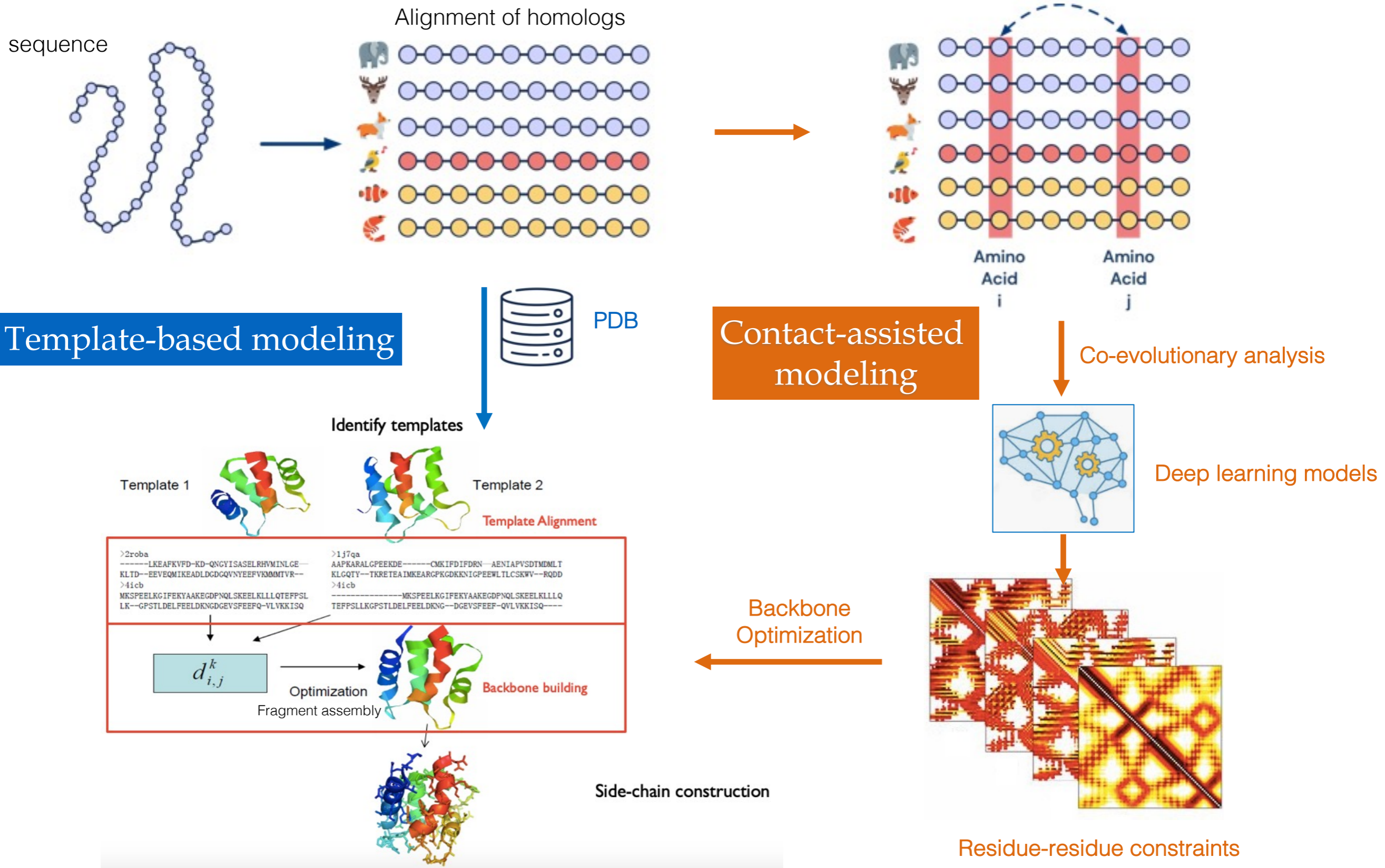
## Progress in recent CASPs





# Structure prediction

# Current status: Protein structure prediction

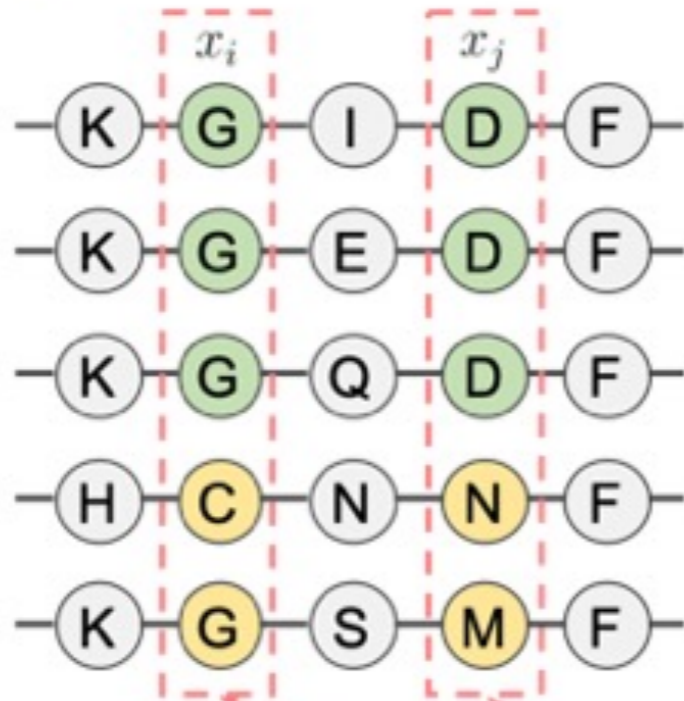


# Exploiting co-evolution for contact prediction

sequence: GEELFTGKKGIDFLGDIVNGSV...

Search homologous sequences

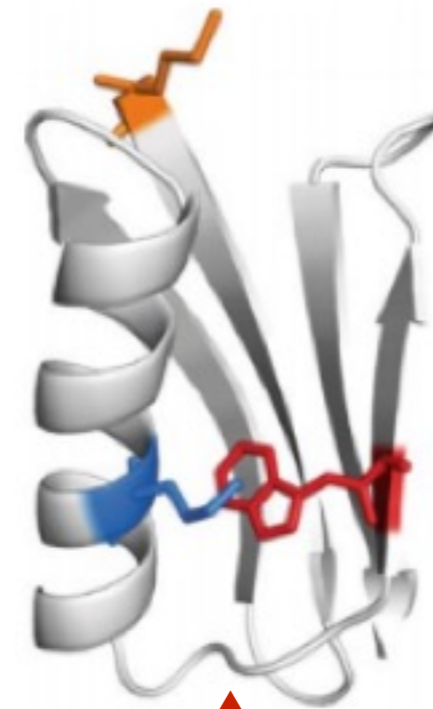
Multiple sequence alignment



Strength of co-evolution

$$E(x) = \sum_i e_i(x_i) + \sum_{i < j} e_{i,j}(x_i, x_j)$$

Local residue preference



Residue contact

Evolutionary and structural constraints

# Learning couplings from protein alignment

Capture independent sites

$$P(\mathbf{x}) = p_1(x_1)p_2(x_2)\cdots p_L(x_L)$$

or equivalently

$$P(\mathbf{x}) = \frac{1}{Z} \exp \left( \sum_i e_i(x_i) \right)$$

Single potentials

Local preference

Amino acid  $i$

$x_i$

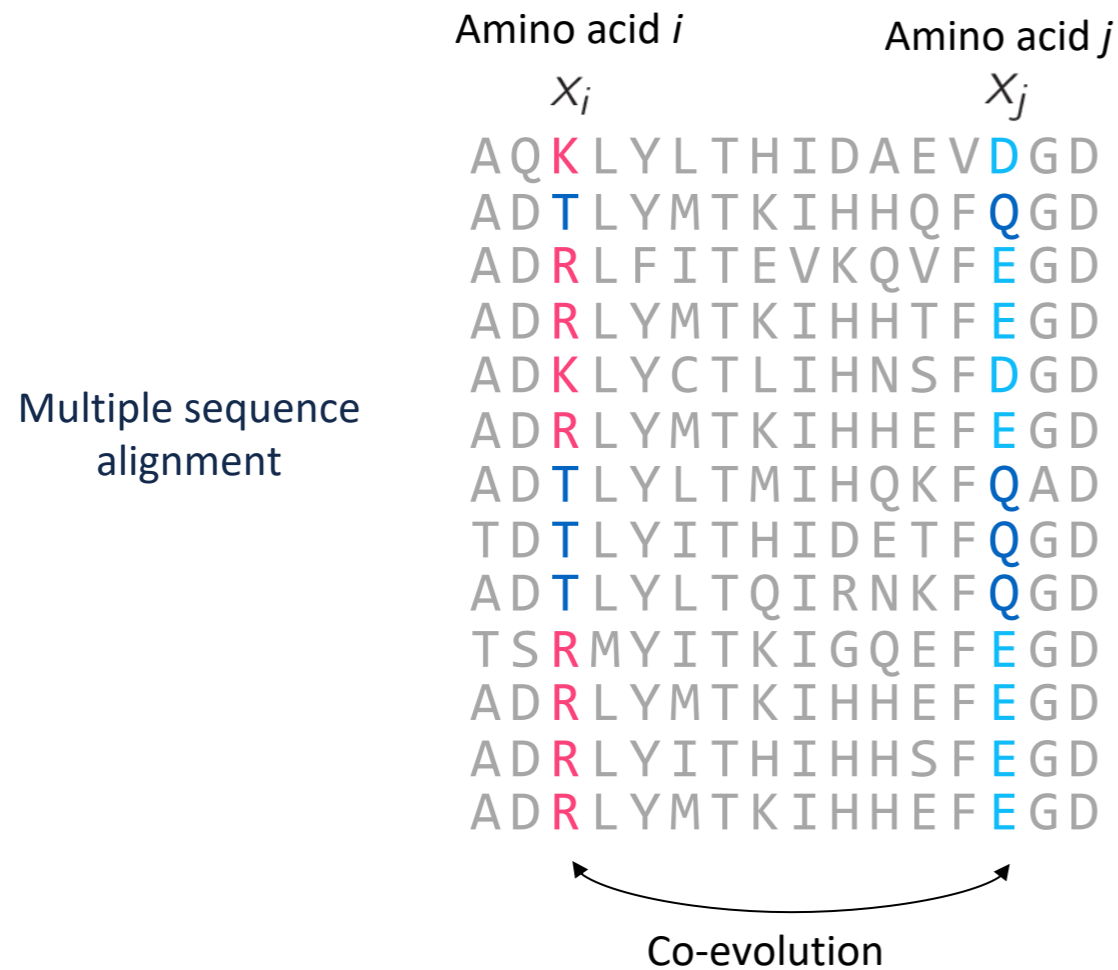
```
A Q K L Y L T H I D A E V D G D
A D T L Y M T K I H H Q F Q G D
A D R L F I T E V K Q V F E G D
A D R L Y M T K I H H T F E G D
A D K L Y C T L I H N S F D G D
A D R L Y M T K I H H E F E G D
A D T L Y L T M I H Q K F Q A D
T D T L Y I T H I D E T F Q G D
A D T L Y L T Q I R N K F Q G D
T S R M Y I T K I G Q E F E G D
A D R L Y M T K I H H E F E G D
A D R L Y I T H I H H S F E G D
A D R L Y M T K I H H E F E G D
```

Multiple sequence alignment



# Learning couplings from protein alignment

Capture pairwise interactions



$$P(\mathbf{x}) = \frac{1}{Z} \exp \left( \sum_i e_i(x_i) + \sum_{i \neq j} e_{ij}(x_i, x_j) \right)$$

Single potentials

Pairwise potentials

Local preference

Co-evolution strength

Markov random field  
Ising (Potts) model  
Undirected graphical model

# Learning with Markov Random Fields

$$L(e) = \prod_{n=1}^N \frac{1}{Z_e^{(n)}} \prod_i^L \exp \left[ e_i(x_i^n) + \sum_{j \neq i} e_{i,j}(x_i^n, x_j^n) \right]$$

Partition  
function

Singleton  
potentials

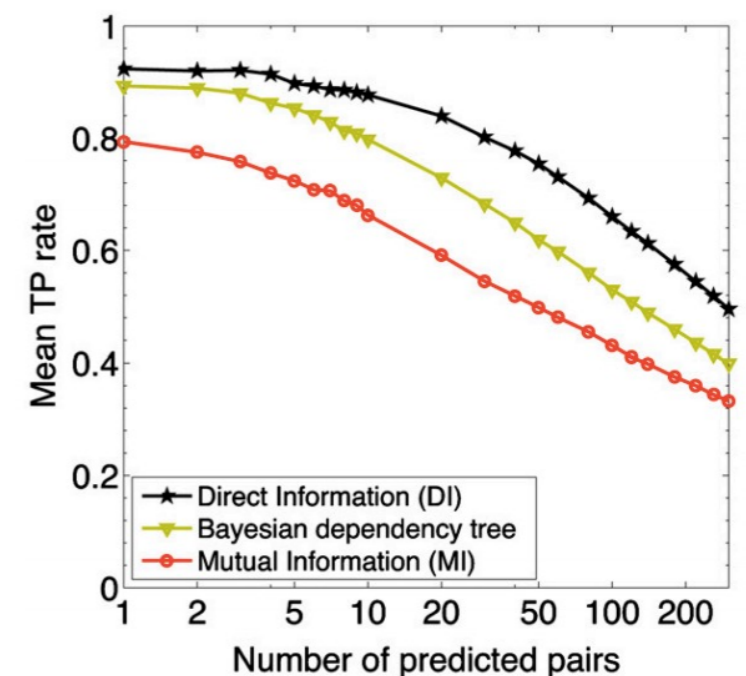
Pairwise  
potentials

Local AA preference

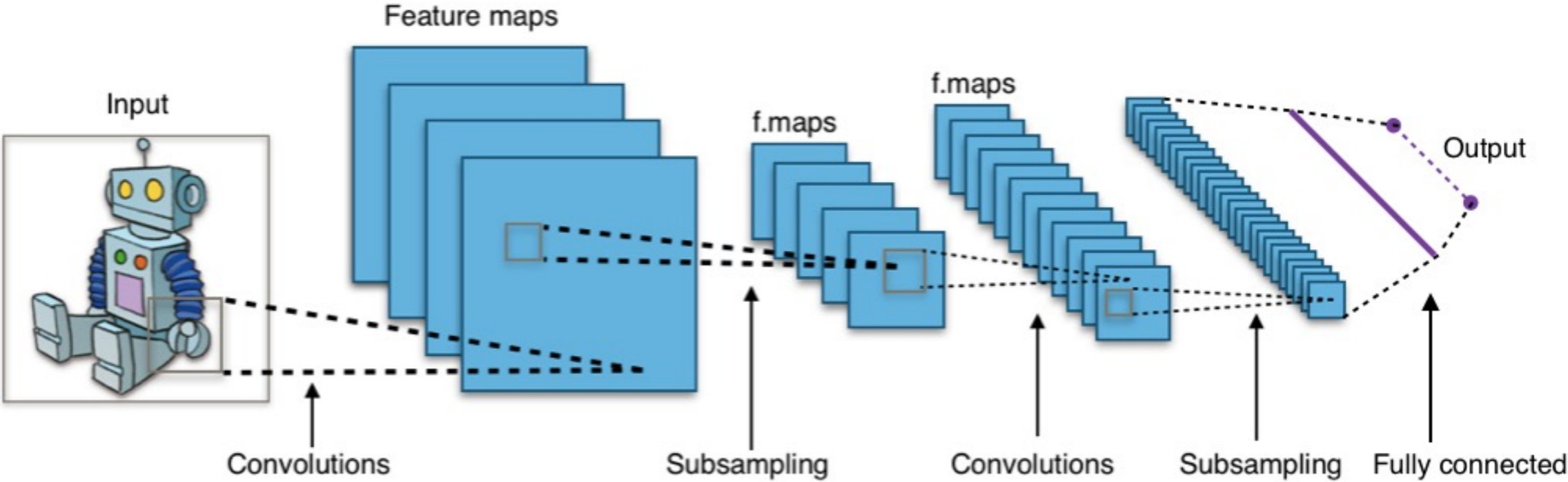
Pairwise AA couplings

Learning algorithms:

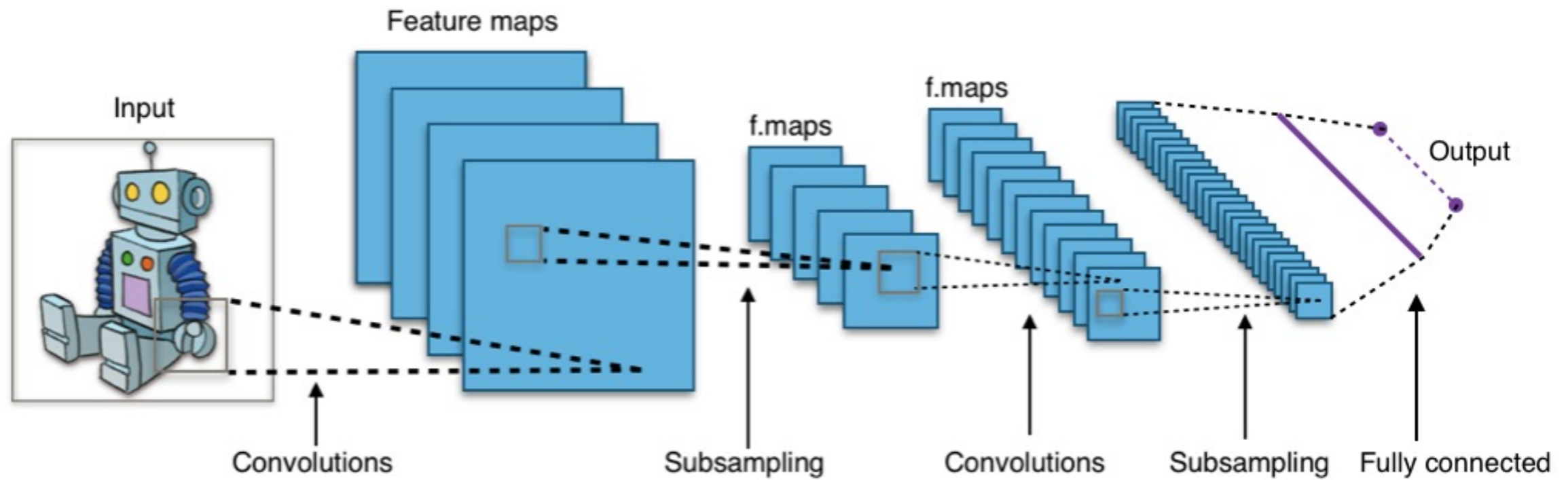
- Mean fields approximation: EVFold, DirectInfo
- Gaussian approximation: PSICOV
- Pseudolikelihood: GREMLIN, CCMpred



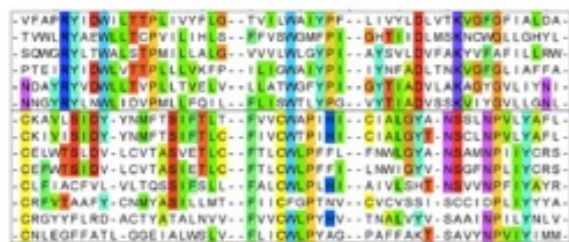
# Deep convolutional NNs recognize image patterns



# Deep convolutional NNs recognize coevolutionary patterns



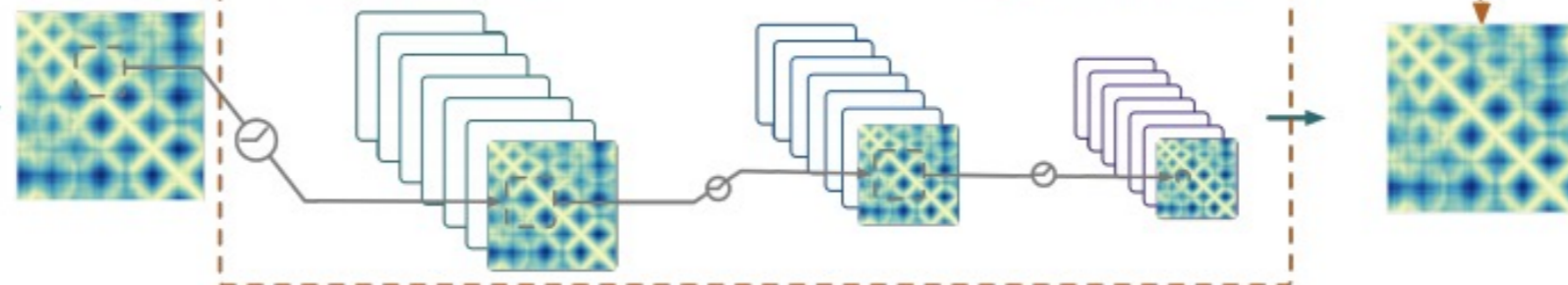
sequence alignment



evolutionary couplings

Low-level feature

High-level feature

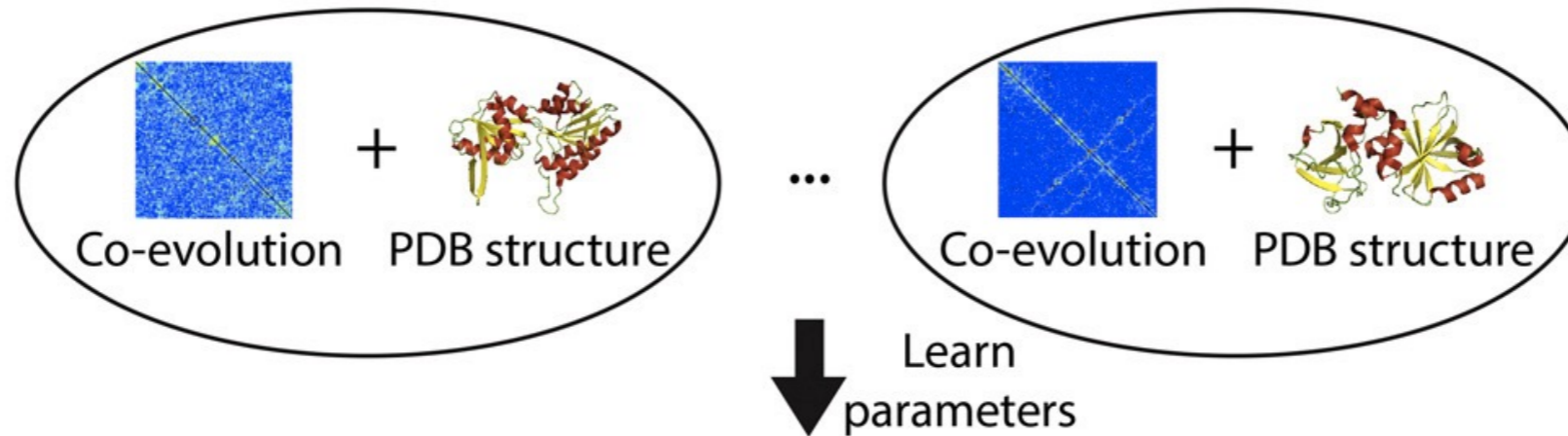


distance contact map

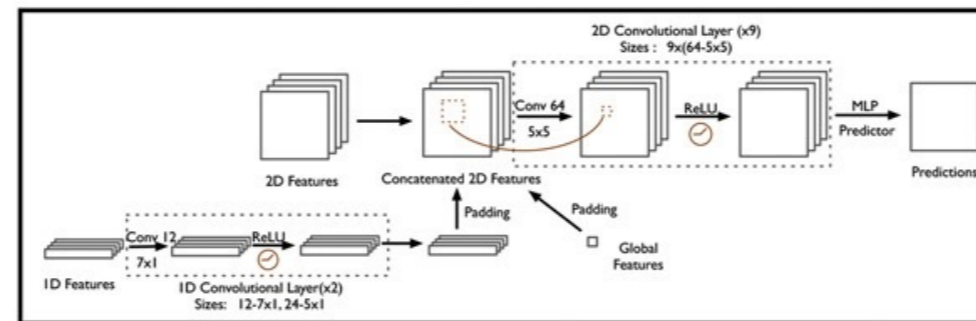


# DeepContact for contact prediction

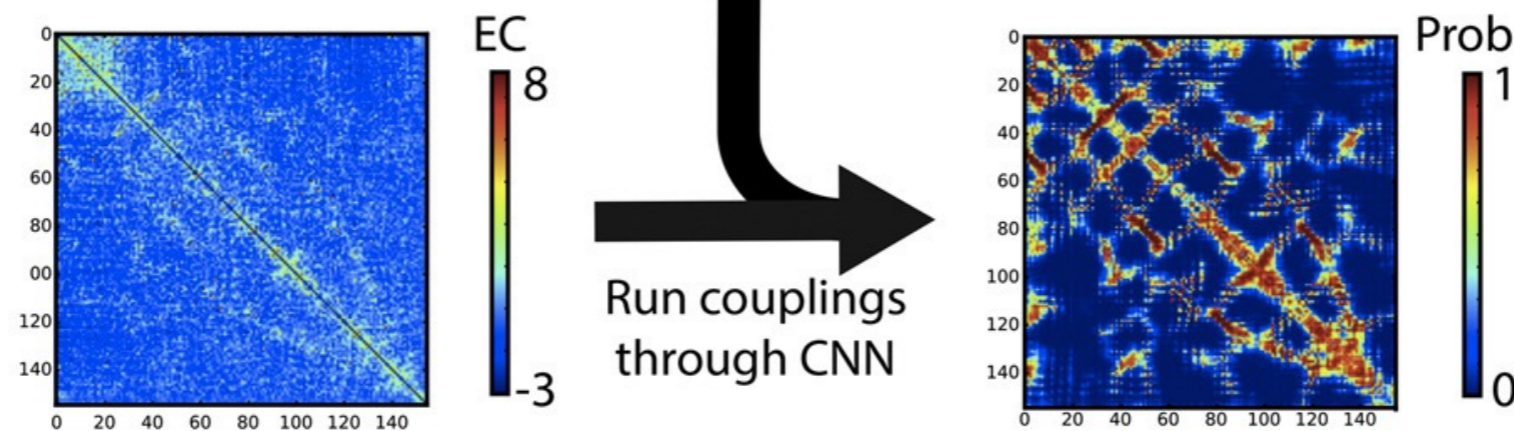
Train on set of inferred couplings and known structures:



DeepContact convolutional neural network (CNN):



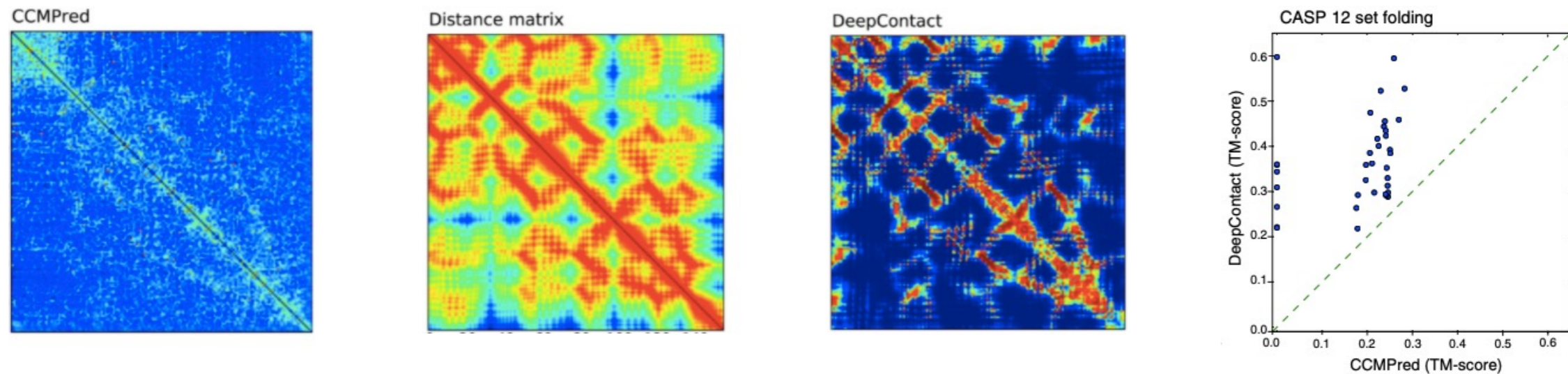
Novel prediction from sequence:



Co-evolution map from sequence

DeepContact contact map

# Deep learning improves coevolution-based contact prediction

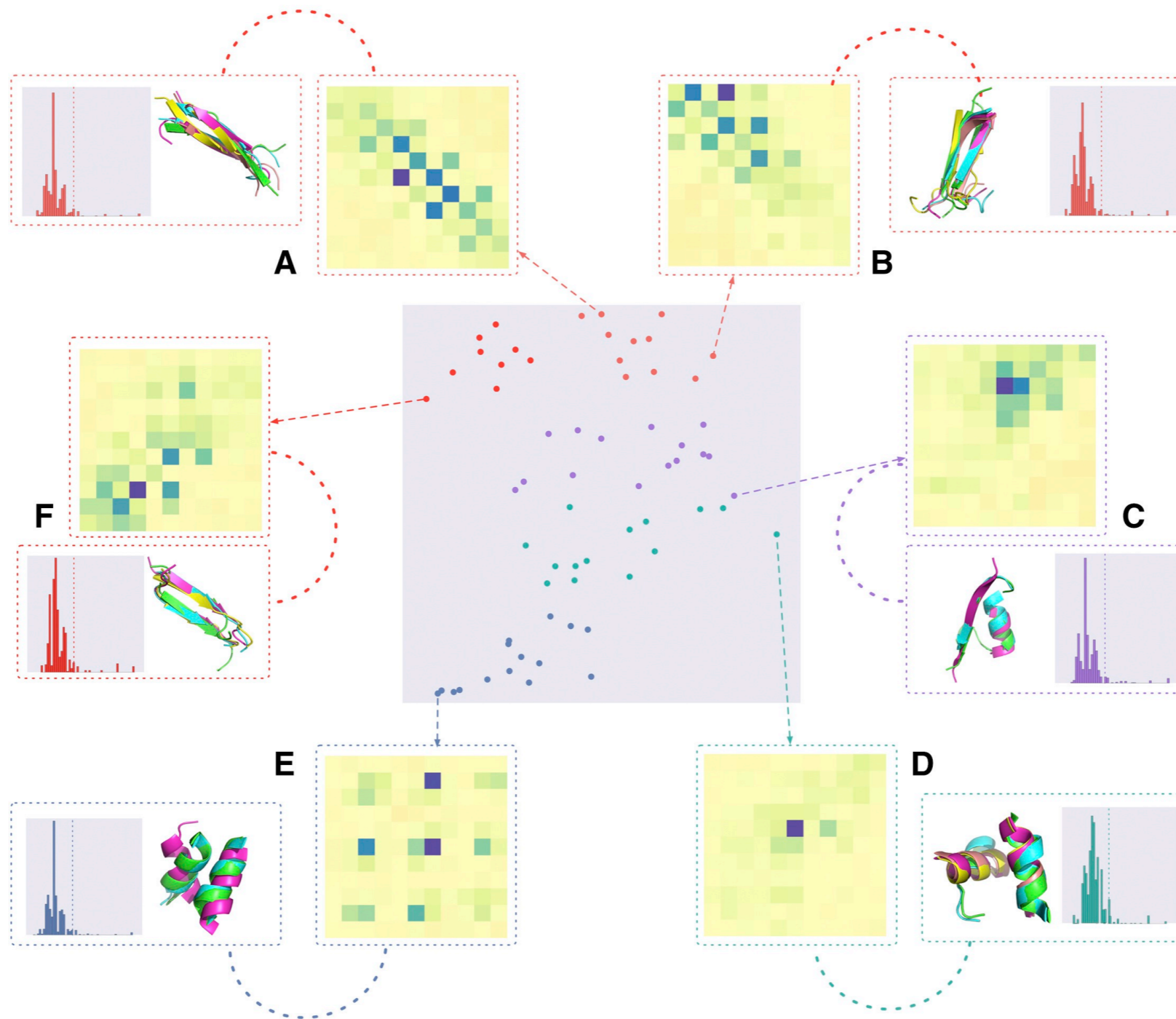


Ranked at the top in CASP12 in 2016 in Z-score ranking on par with two other deep-learning based methods (RaptorX-Contact and (Deep) MetaPSICOV) on other metrics

#	GR code	GR name	Domains Count	SUM Zscore (>-2.0)
1.	079	iFold_1	37	40.1644
2.	219	Deepfold-Contact	36	34.5989
3.	451	RaptorX-Contact	38	34.4778
4.	109	naive	36	27.8373
5.	431	Shen-Group	38	21.3752
6.	013	MetaPSICOV	38	17.6201
7.	287	MULTICOM-CLUSTER	38	15.1225
8.	345	MULTICOM-NOVEL	38	14.1530
9.	236	MULTICOM-CONSTRUCT	38	12.3499
10.	320	raghavagps	38	11.4462

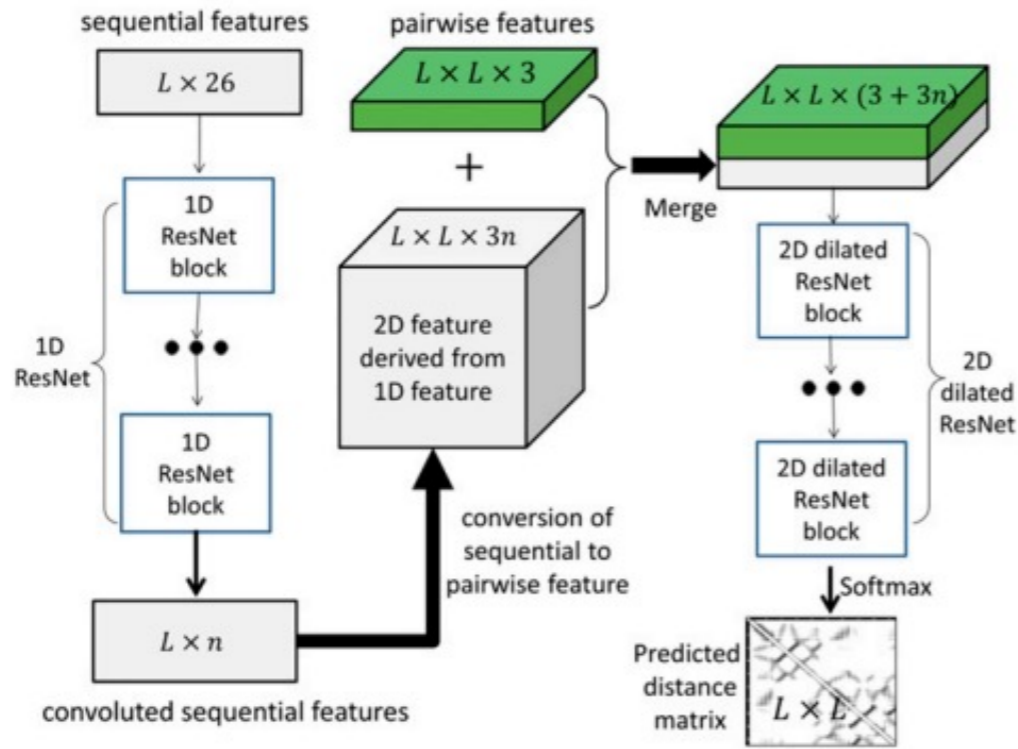
# Why is deep learning effective?

2D projection of 1st layer filters using tSNE

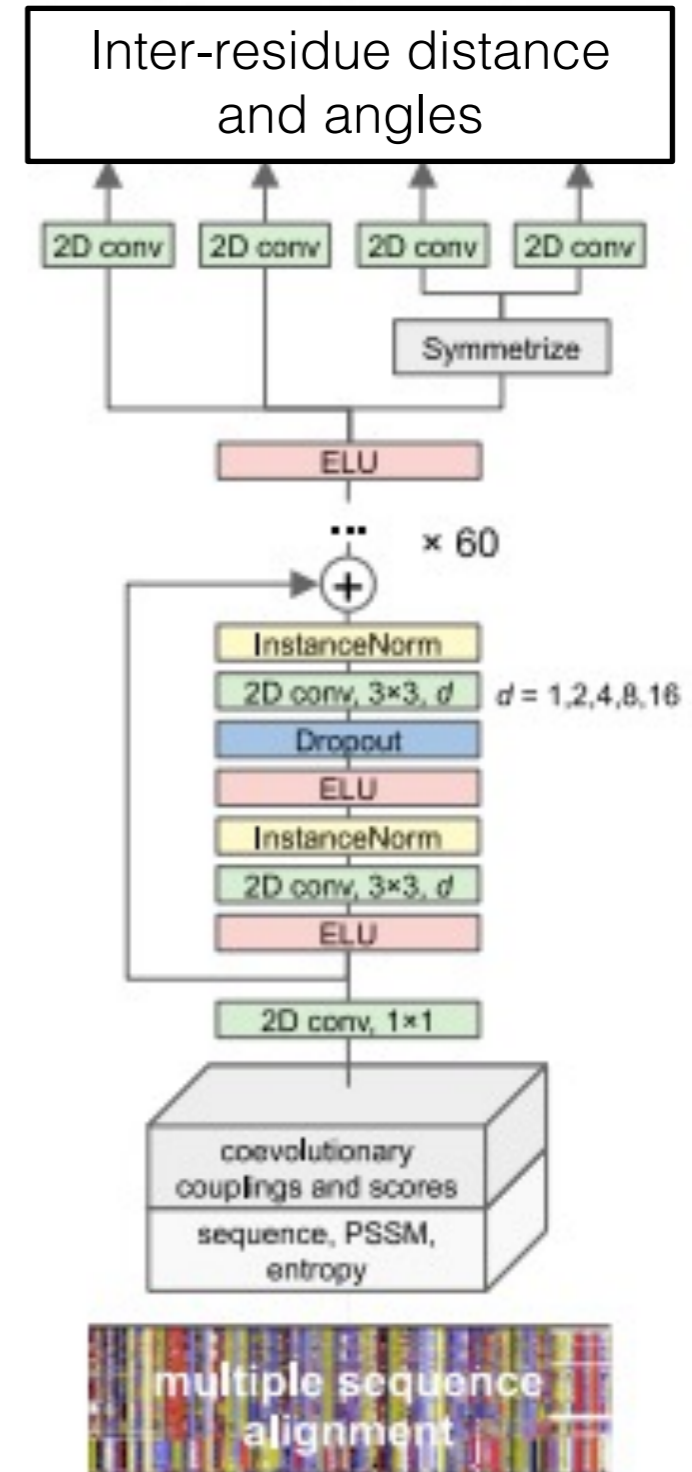


Deep neural network learns contact patterns

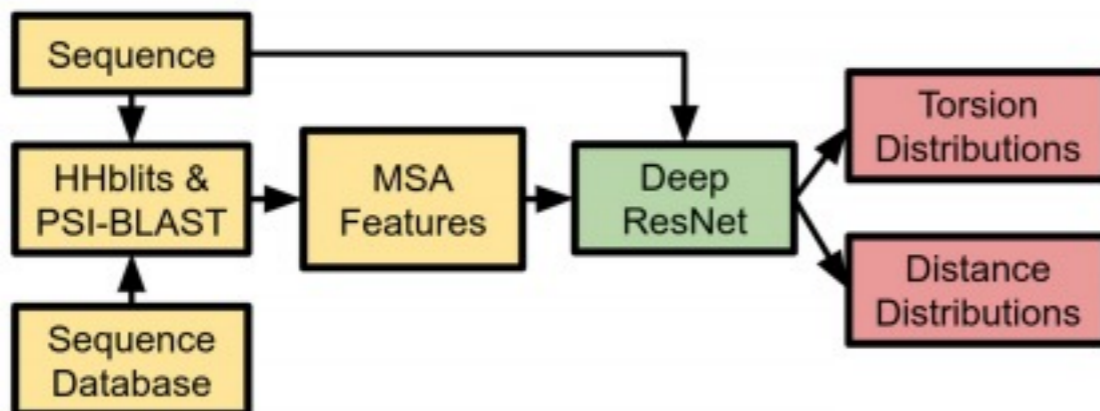
# Recent developments go beyond contact prediction



Rosetta

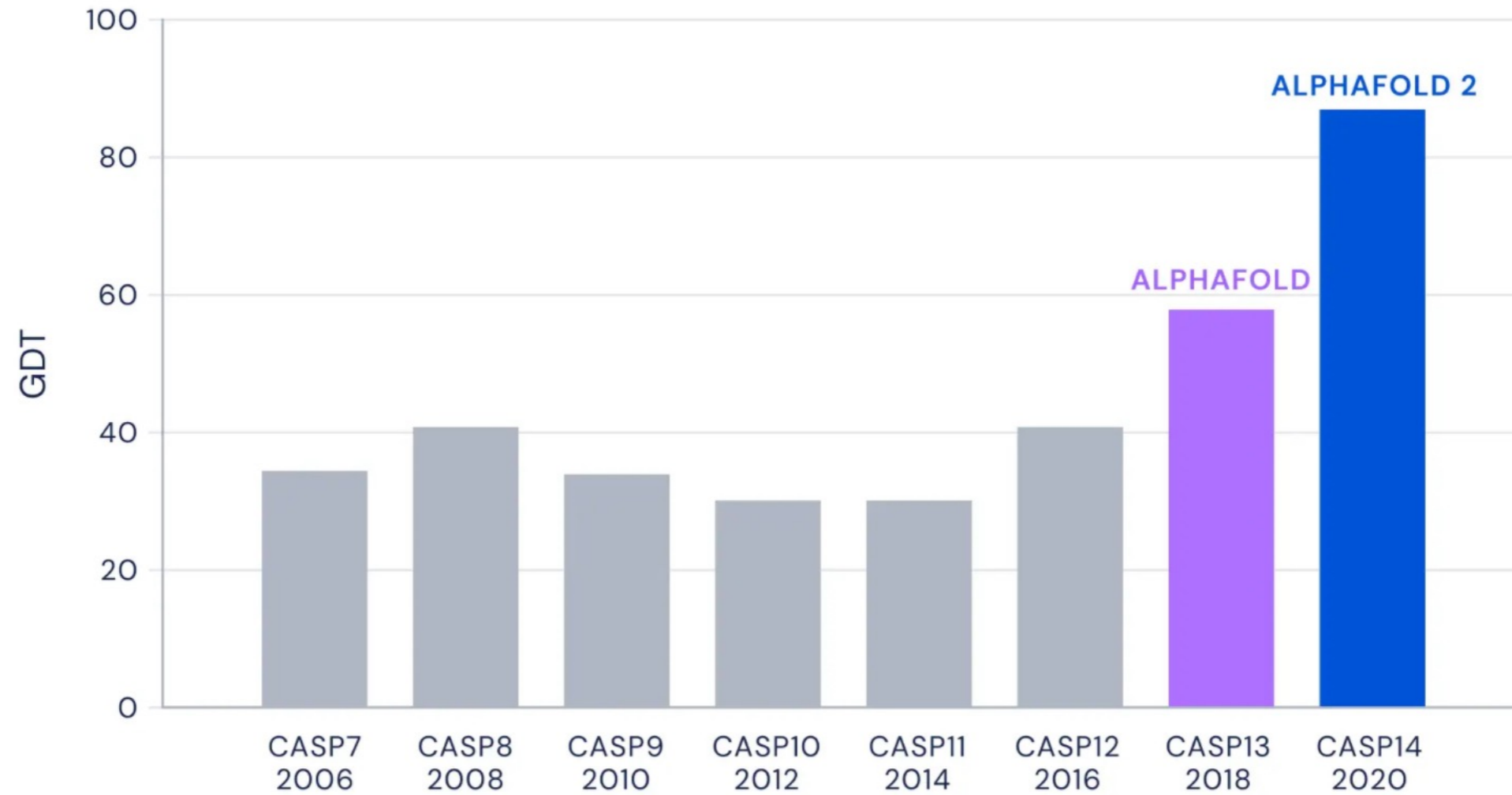


ALPHAFOLD



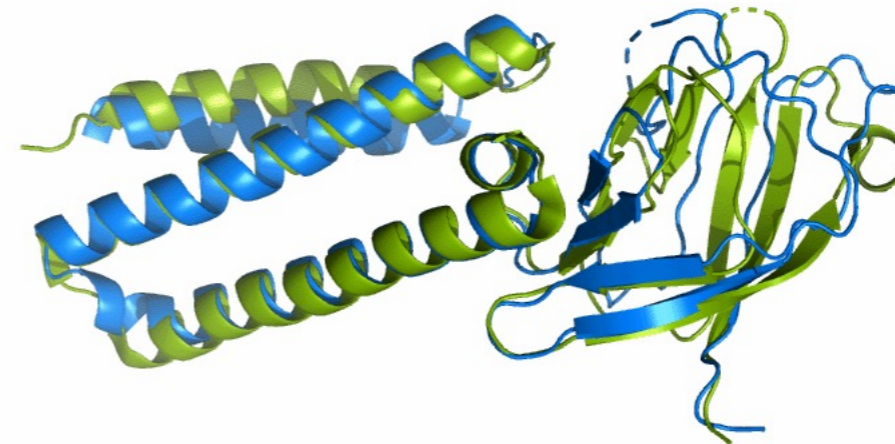
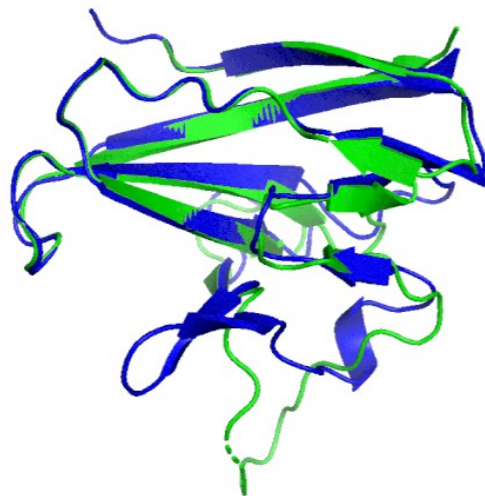


# CASP14: DeepMind's AlphaFold 2



Blue: Predicted  
Green: Actual

ORF8



ORF3a



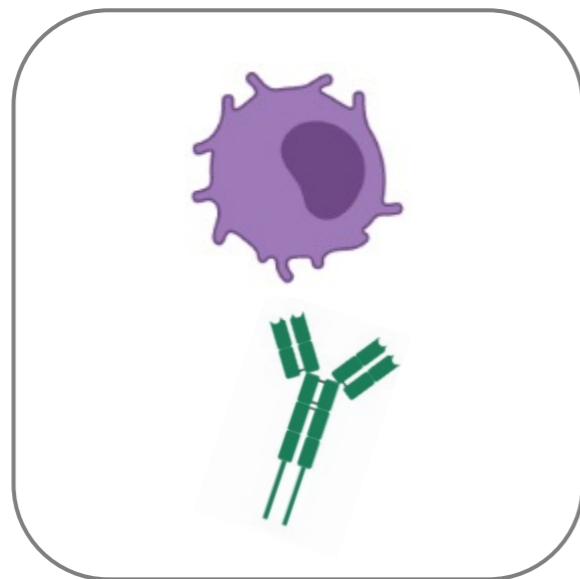
# Function Prediction

# Optimization of protein function



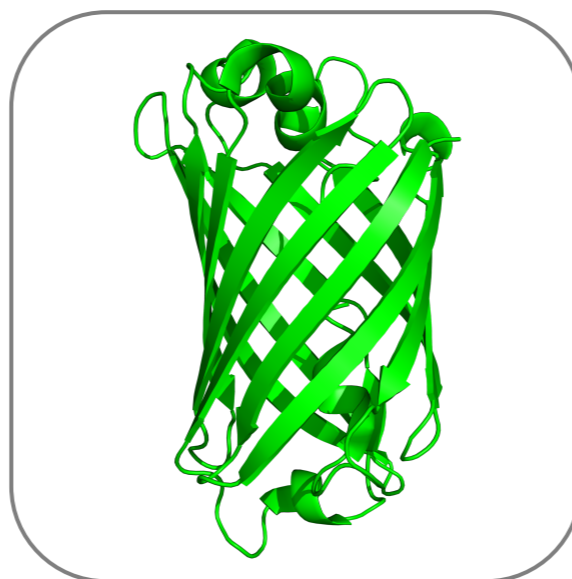
Sequence →

Function



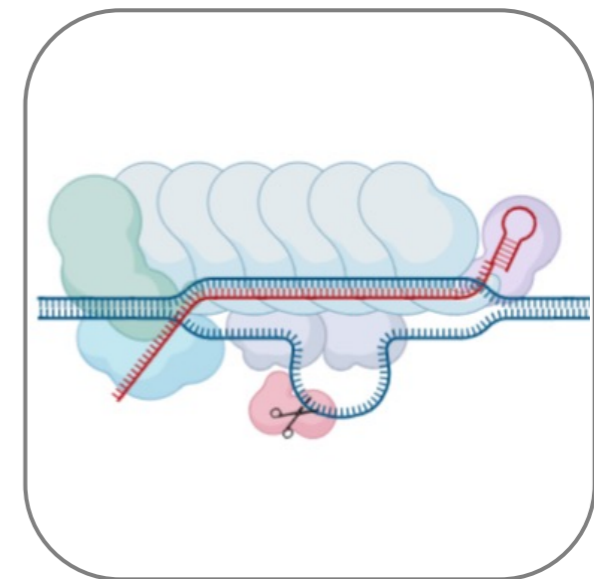
Antibody

*Binding affinity*



Fluorescent protein

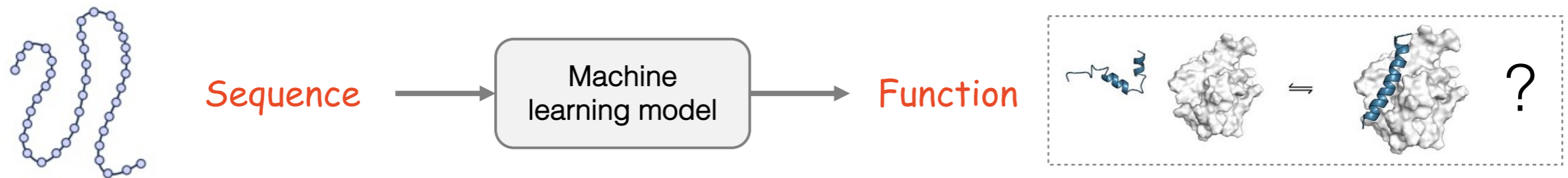
*Fluorescence*



CRISPR/Cas9

*Specificity*

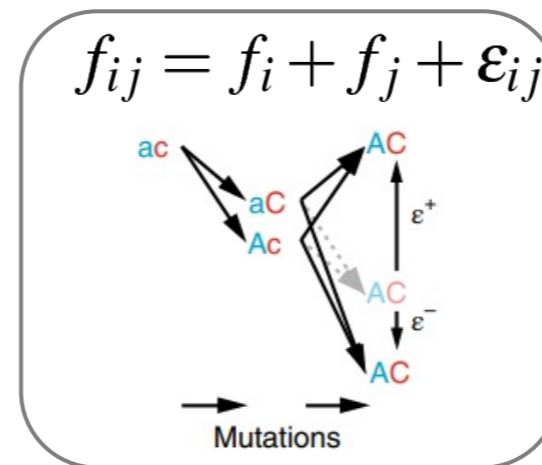
# Sequence-to-function modeling



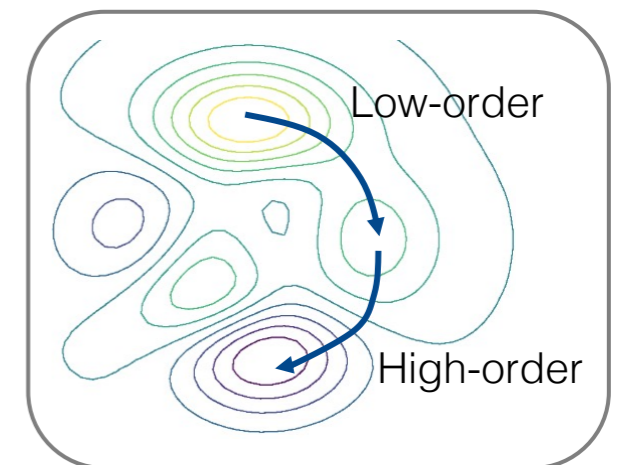
Need to differentiate function levels of closely related sequences

Sequence	Fitness
...DNGVDGEWTYDDATKTFTVTE	1.0
...DNG <b>C</b> DGEWTYDDATKTFTVTE	0.2
...DNGV <b>W</b> GEWTYDDATKTFTVTE	3.9
...DNGV <b>W</b> GEWTYDDATKTFT <b>F</b> TE	5.4
...DNGV <b>M</b> GEWTYDDATKTFT <b>D</b> TE	0.1

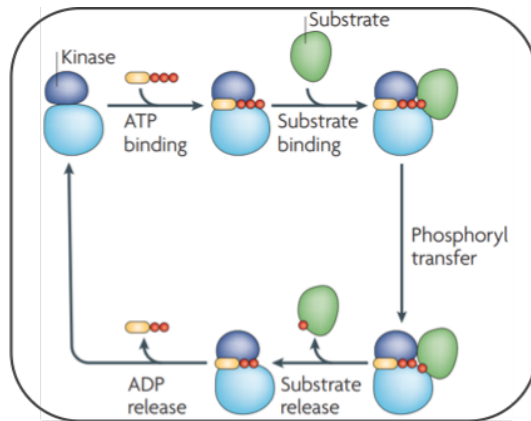
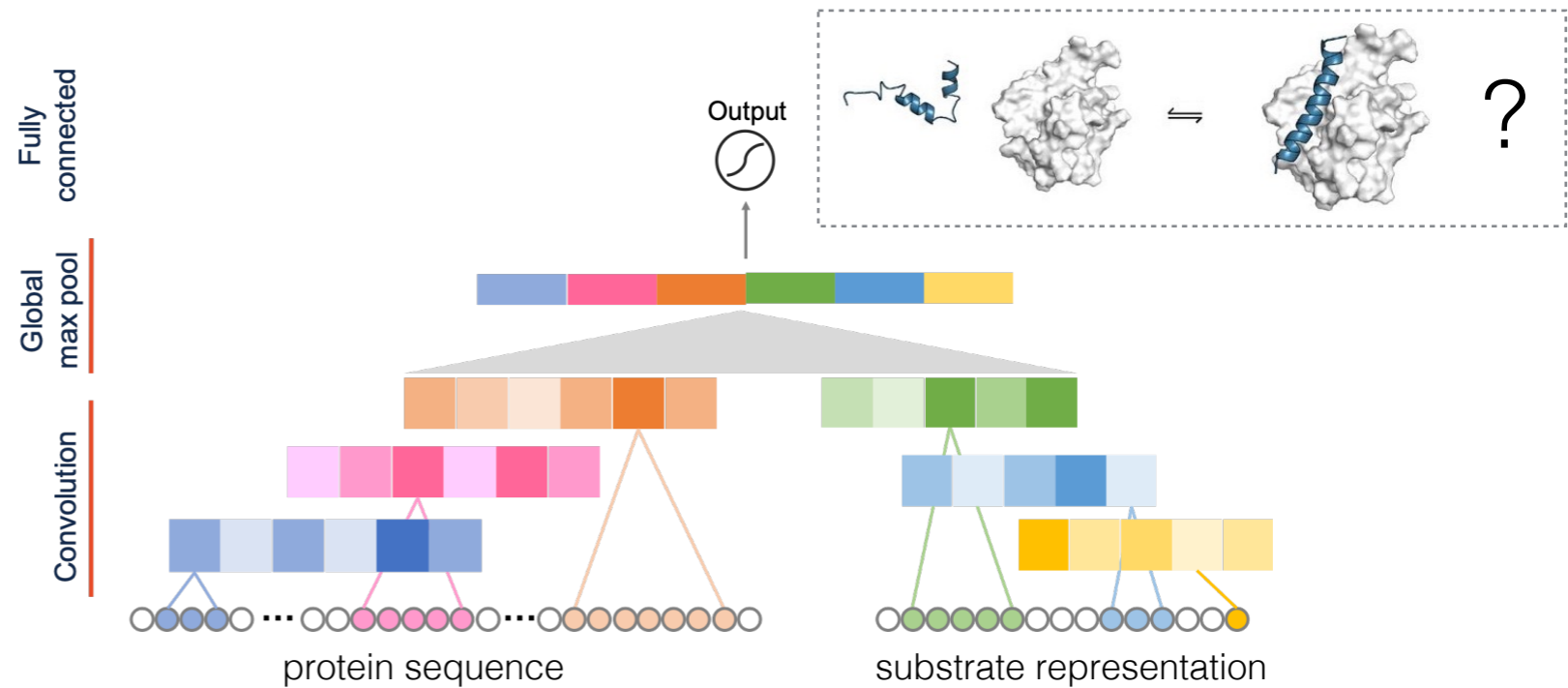
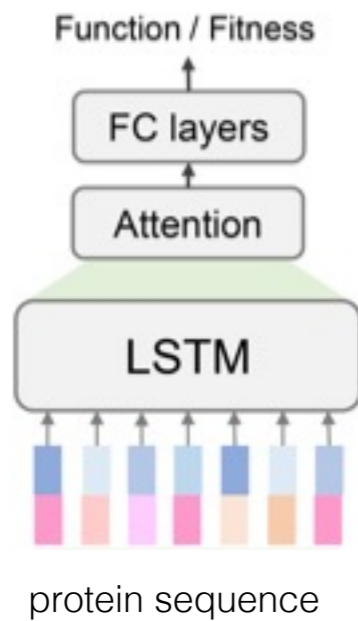
Need to model non-additive effect (epistasis)



Need to generalize to unseen sequences/mutations



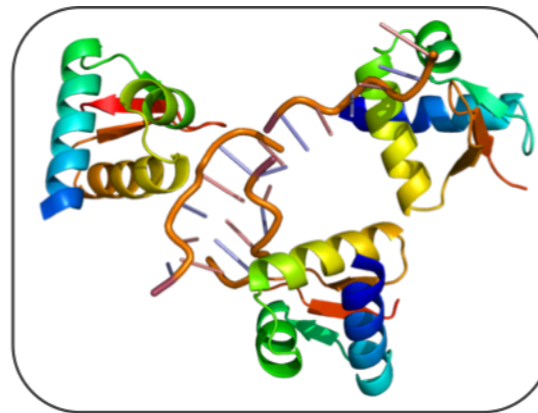
# Successful sequence-to-function models



**Kinase-peptide binding  
(protein phosphorylation)**

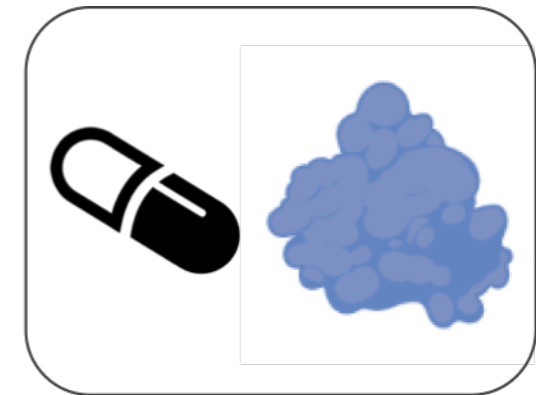
[ Luo\*, Ma\* et al., *RECOMB*, 2018 ]

[ Luo\*, Ma\* et al., *RECOMB*, 2019 ]



**Protein-RNA binding**

[ Su\*, Luo\* et al, *PLOS CB*, 2019 ]

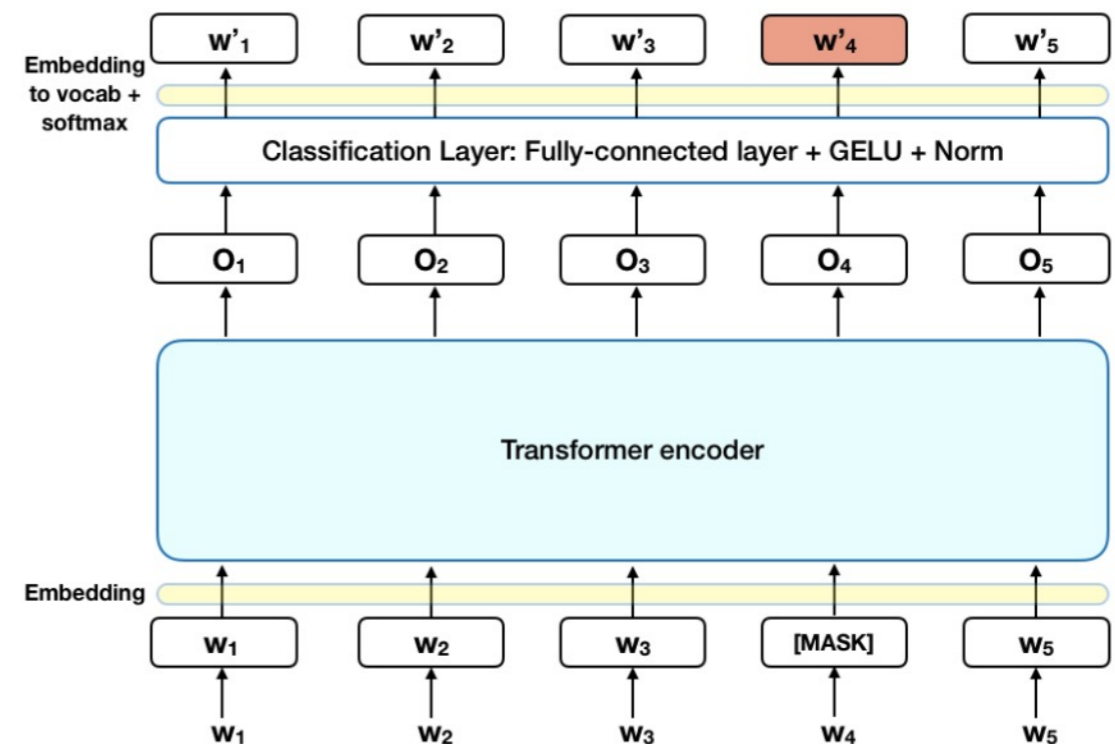
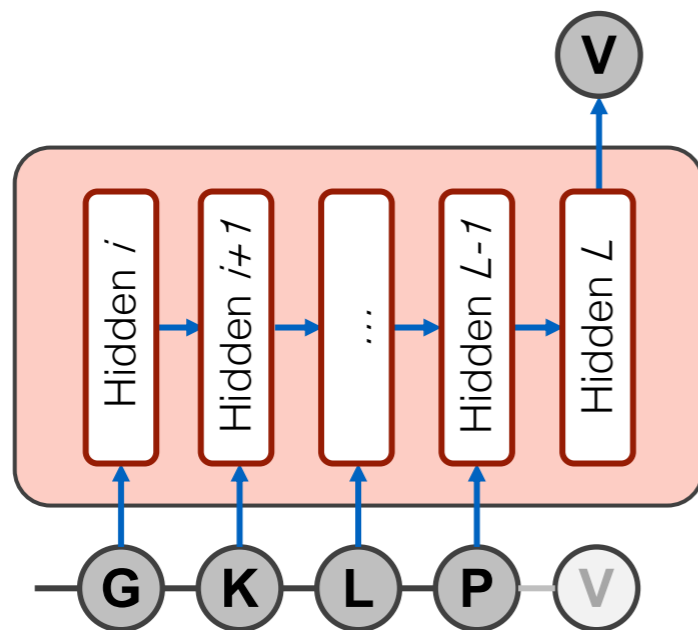


**Kinase-drug binding**

[ Winner of IDG-DREAM Challenge, 2019 ]

# Challenge: labeled data are expensive to get

Idea: unsupervised representation learning using language models using unlabeled data

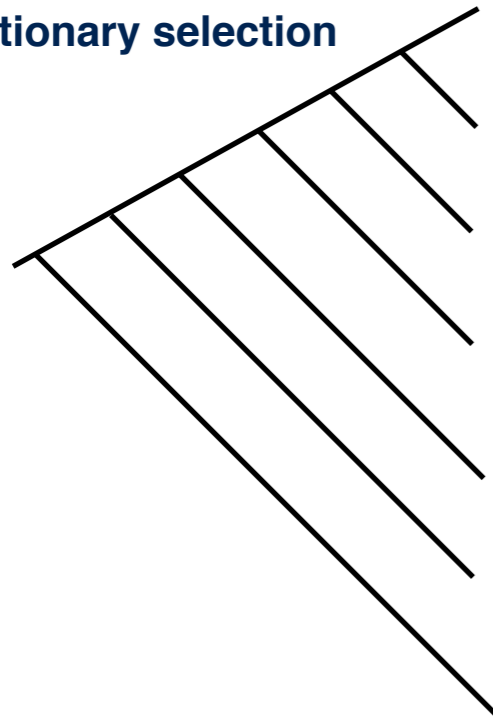


Trained on Pfam / UniProt database with unlabeled data

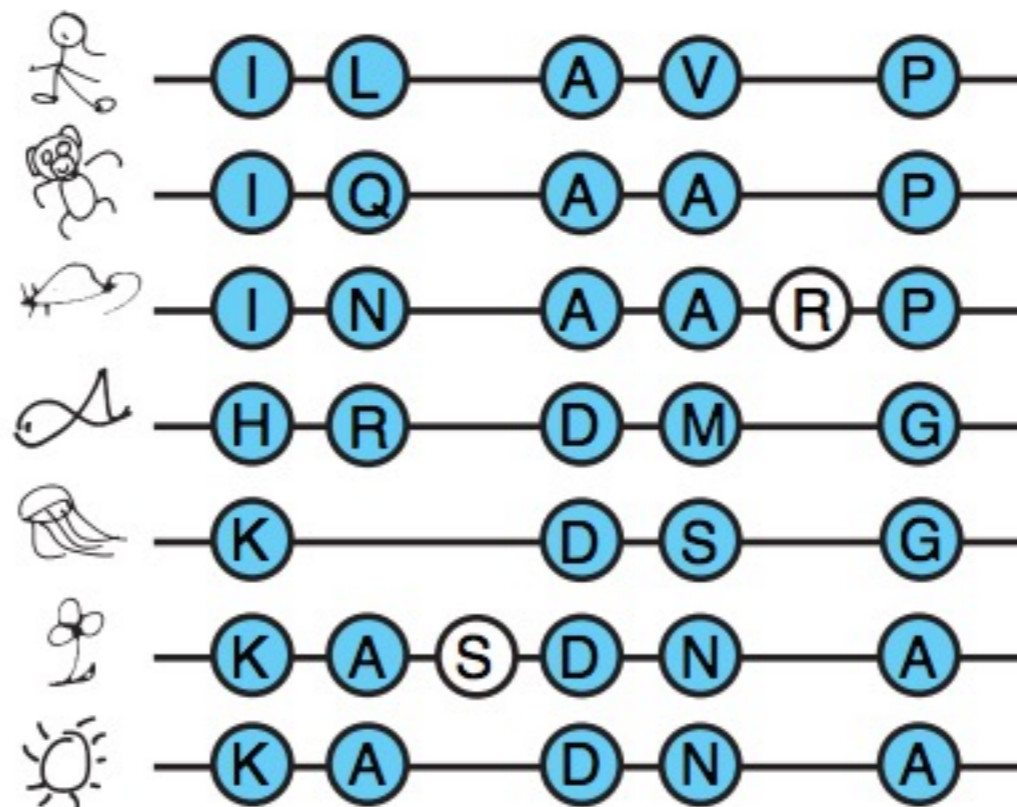
Too general, not specific/sensitive to a single or few changes in the sequence

# Another idea: Learning from natural evolution

evolutionary selection



homologous sequences

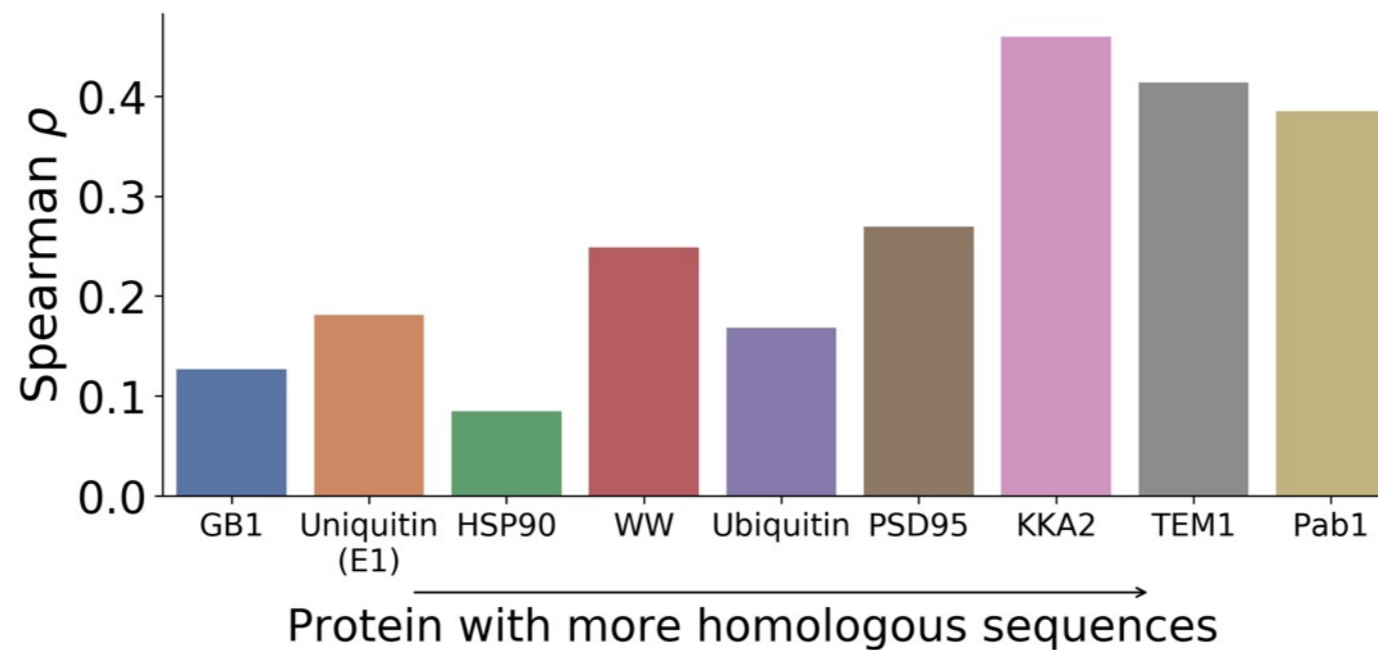
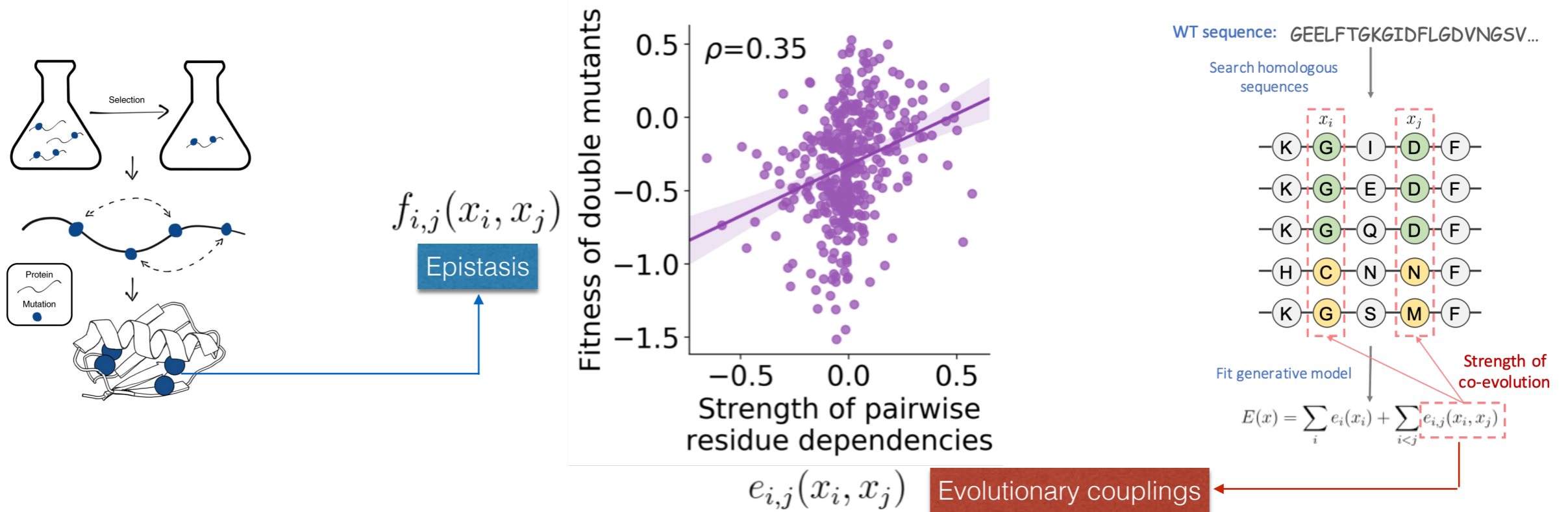


Survived?

+  
+  
+  
+  
+  
+  
+

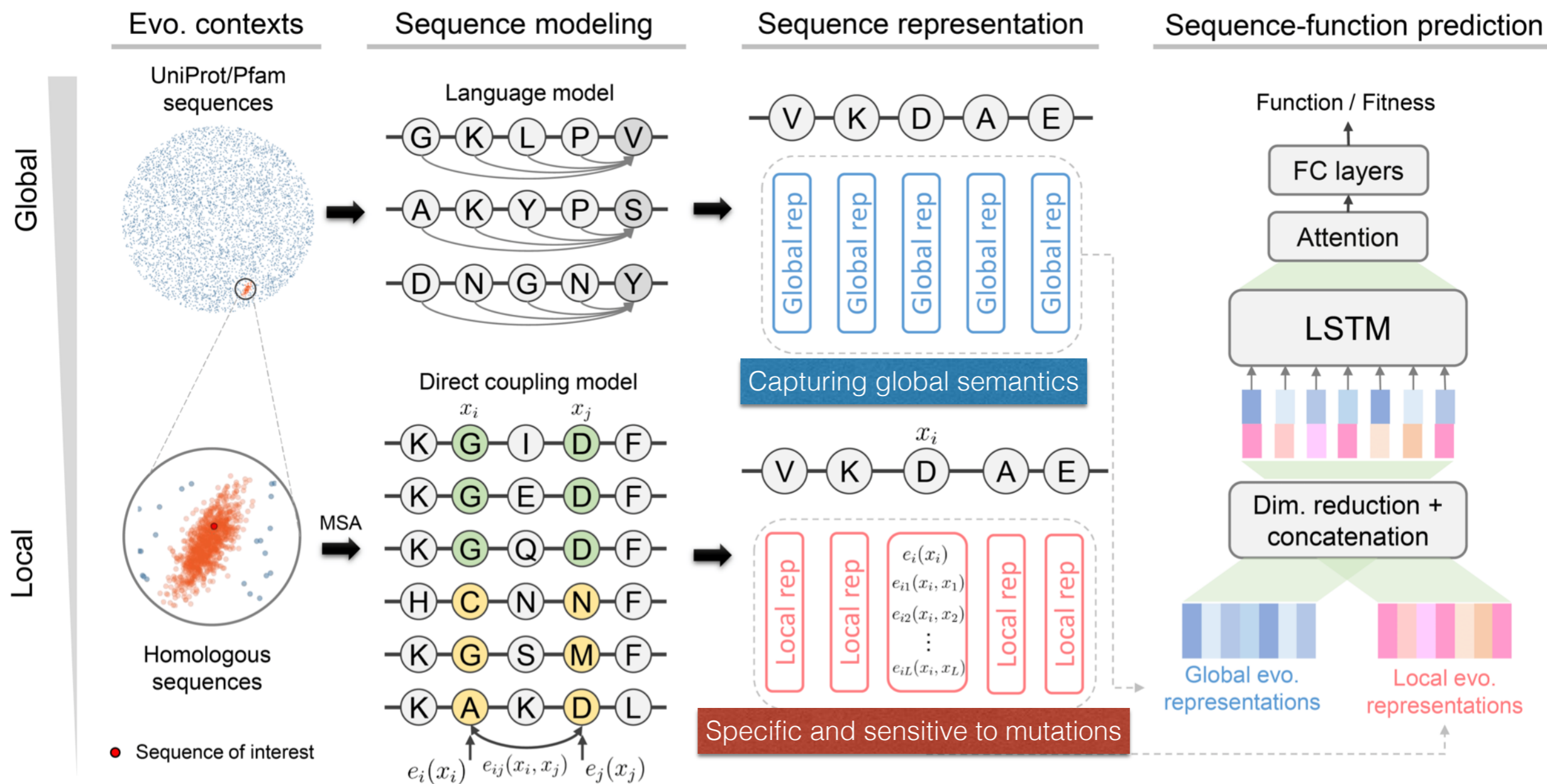
Not enough for fitting language models but enough for getting good features

# Co-evolution correlates with function

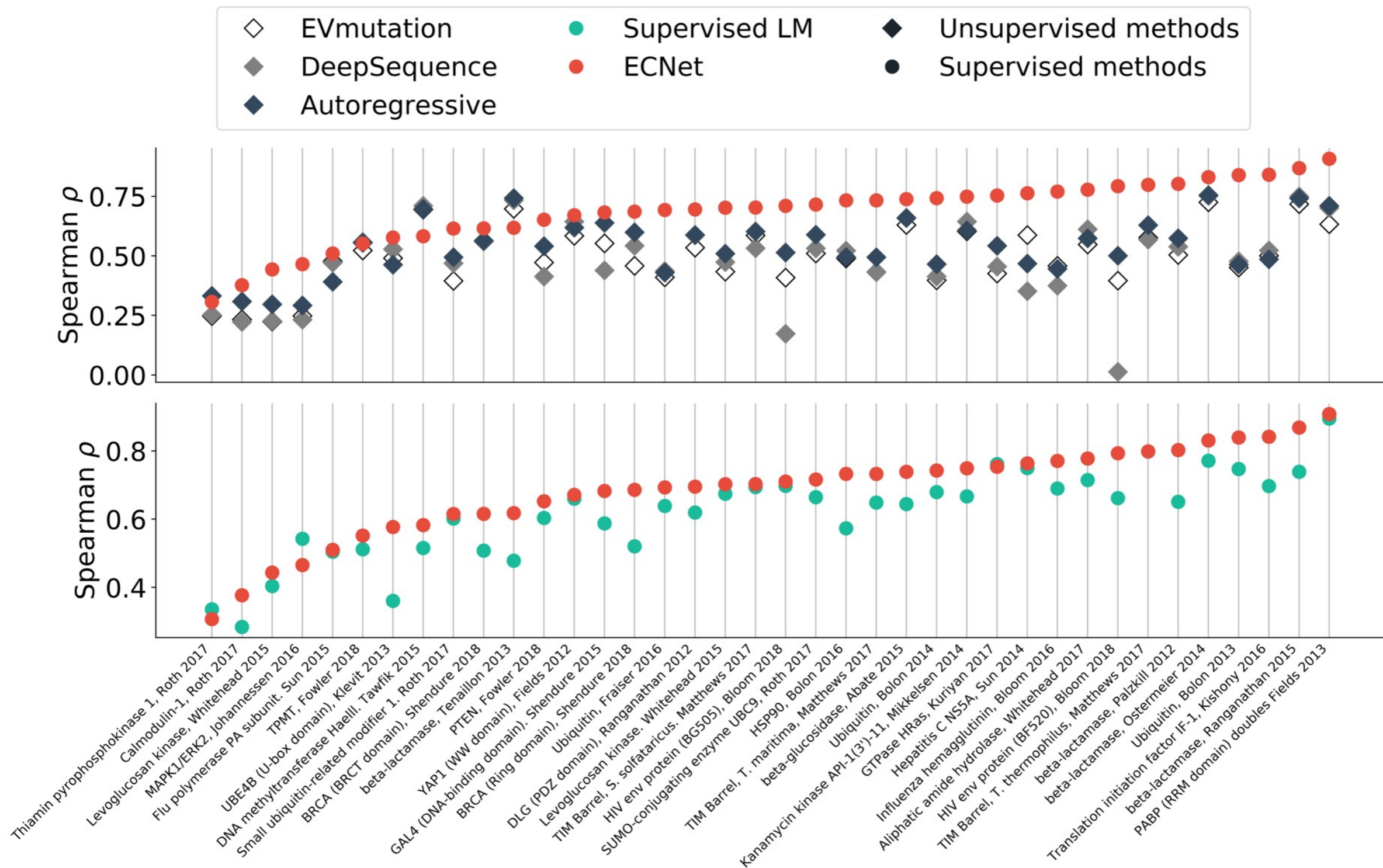




# ECNet: integrating evolutionary contexts for protein function prediction

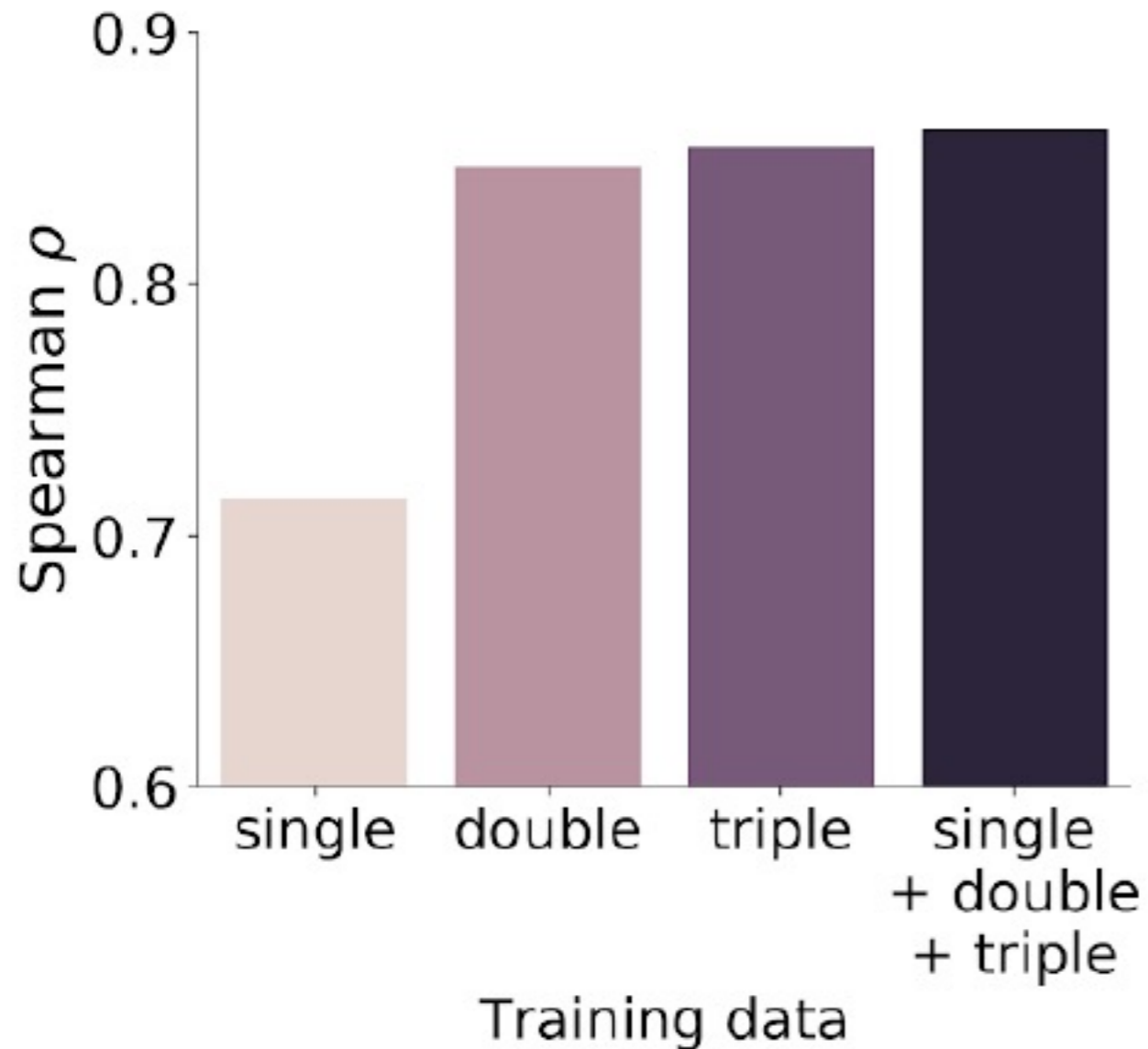


# Evaluation on single-mutation datasets

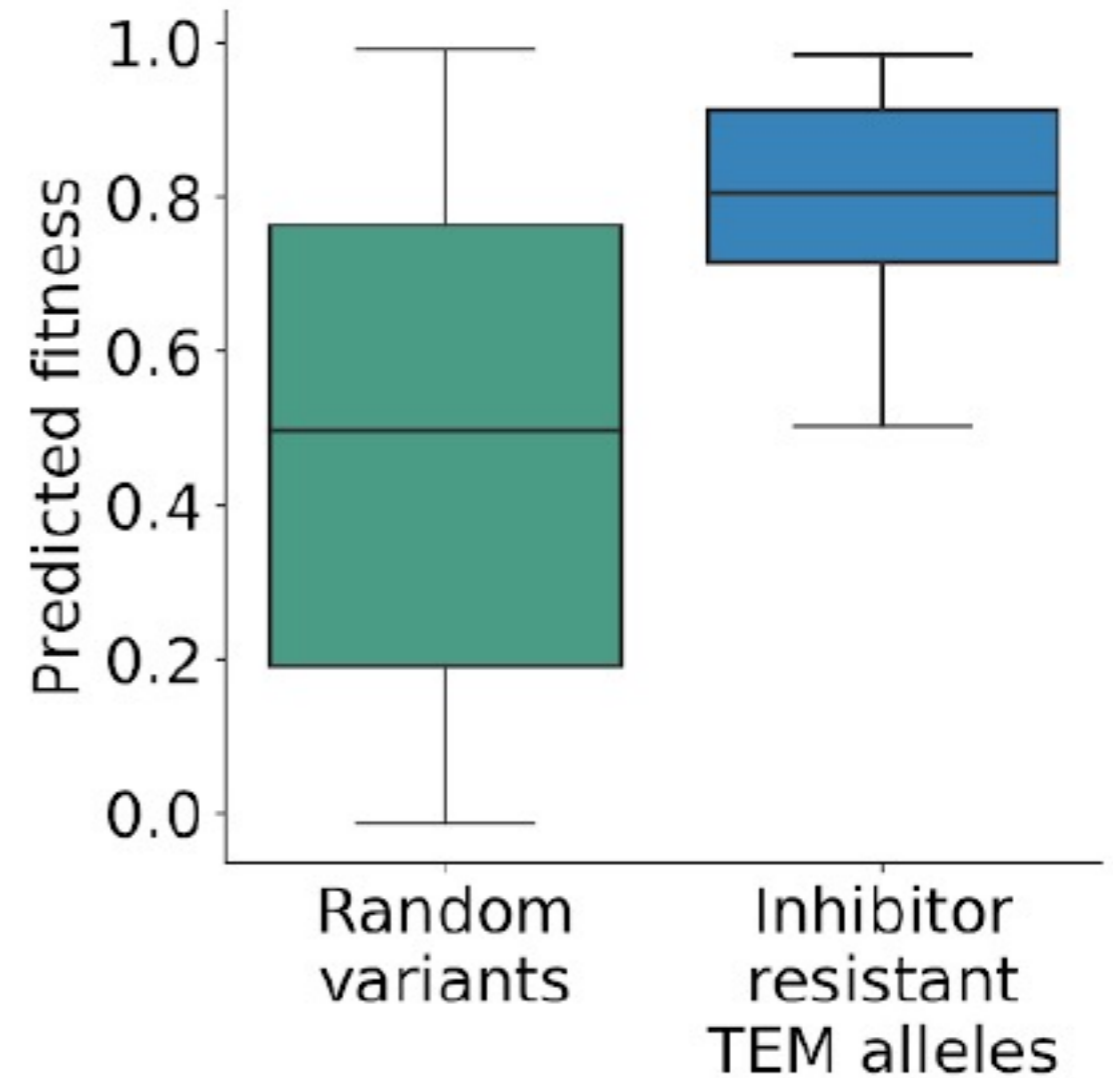


# Generalization to high-order mutations

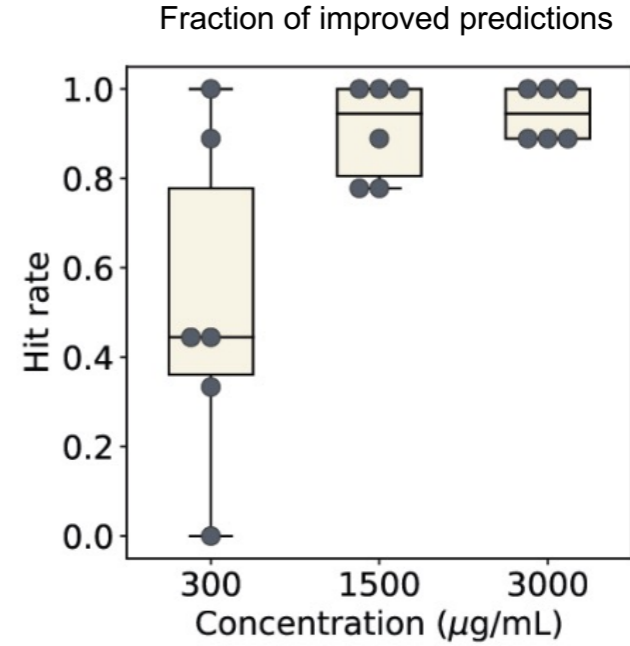
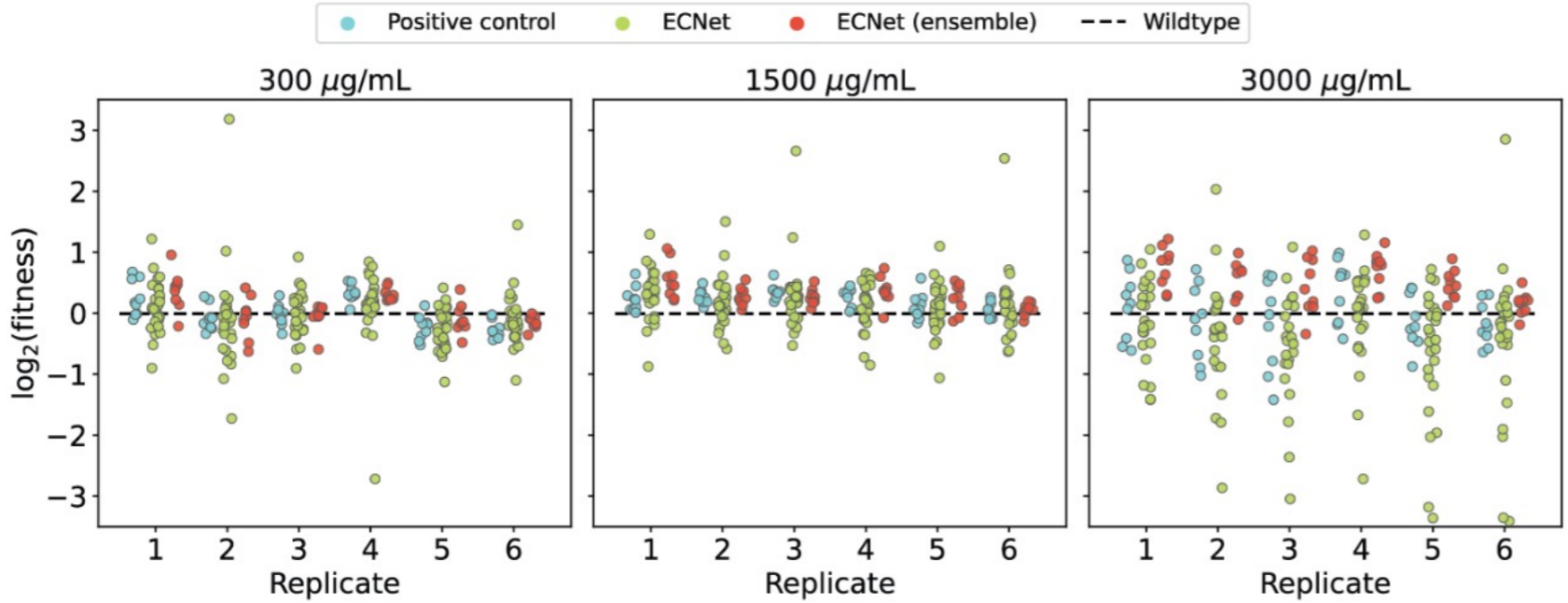
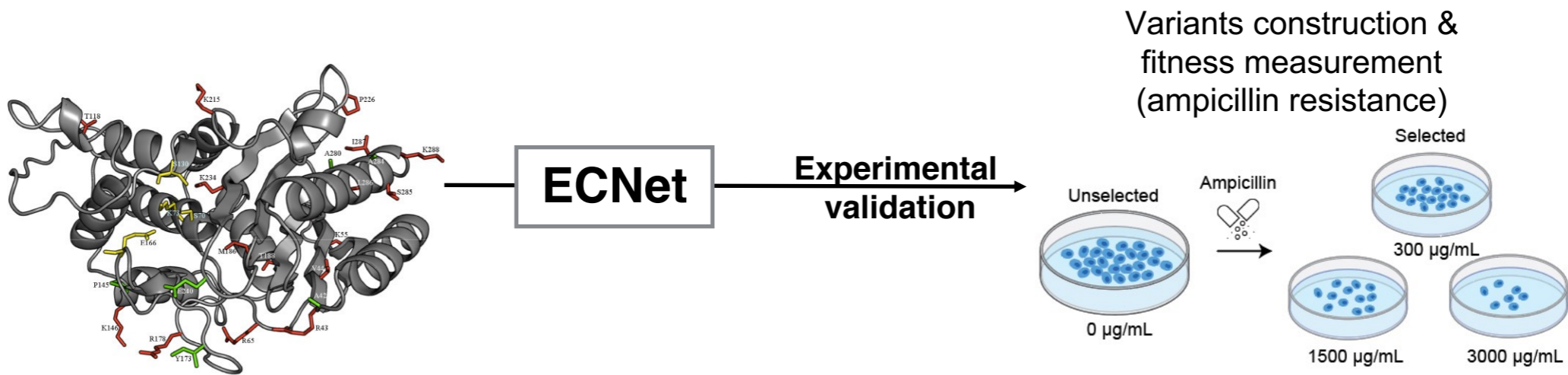
Test on high-order (4~11) GFP variants



High-order (3~15) resistant TEM



# Engineering of inhibitor-resistant TEM1 beta-lactamase



# Thank you & acknowledgements

**Students:** Sheng Wang, Yunan Luo, Hoon Cho, Nate Russell, Wesley Qian, Yang Liu, Yufeng Su, Palash Sashittal,

**Collaborators:** Bonnie Berger, Vik Khurana, Trey Ideker, Mikko Taipale, Jianzhu Ma, Jianyang Zeng, Qiang Liu, Mohammed El-Kebir, Huimin Zhao, Martin Burke, David Heckerman, ChengXiang Zhai, Jiawei Han, Aditya Parameswaran, Jadwiga Bienkowska, Lenore Cowen, Alex Schwing, Jonathan Carlson, Sumaiya Nazeen, Dina Zielinski



Alfred P. Sloan  
FOUNDATION



U.S. DEPARTMENT OF  
**ENERGY**



Microsoft

