# Protein Design & Structural Prediction

Alena Khmelinskaia

Department Chemie
Butenandtstr. 5-13, Haus B
81377 München
akhmelin@cup.lmu.de

LMU

LUDWIG-
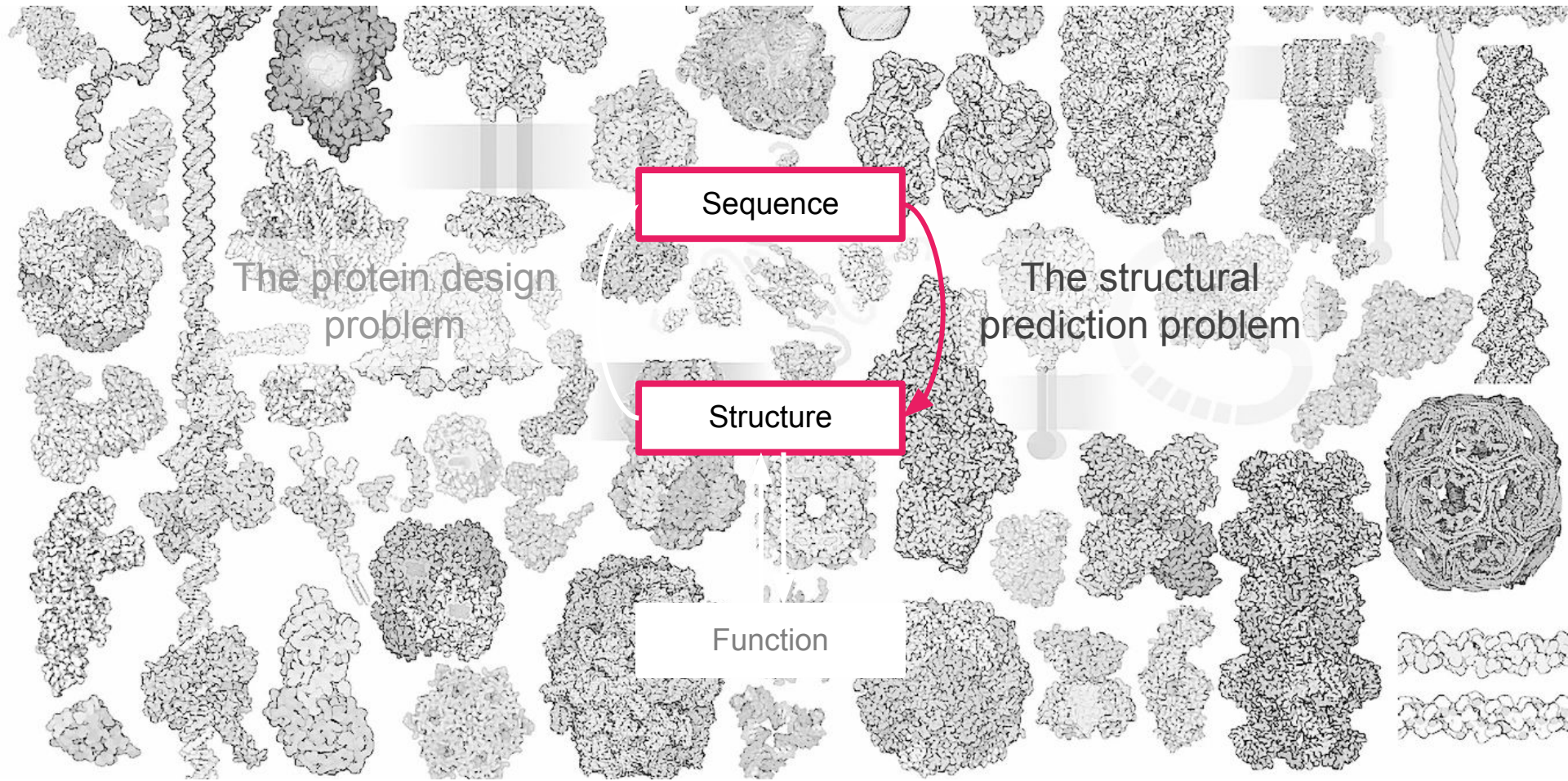MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

# Overview

# But first a recap !

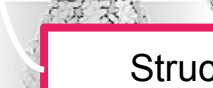# Protein structural prediction (with AI)

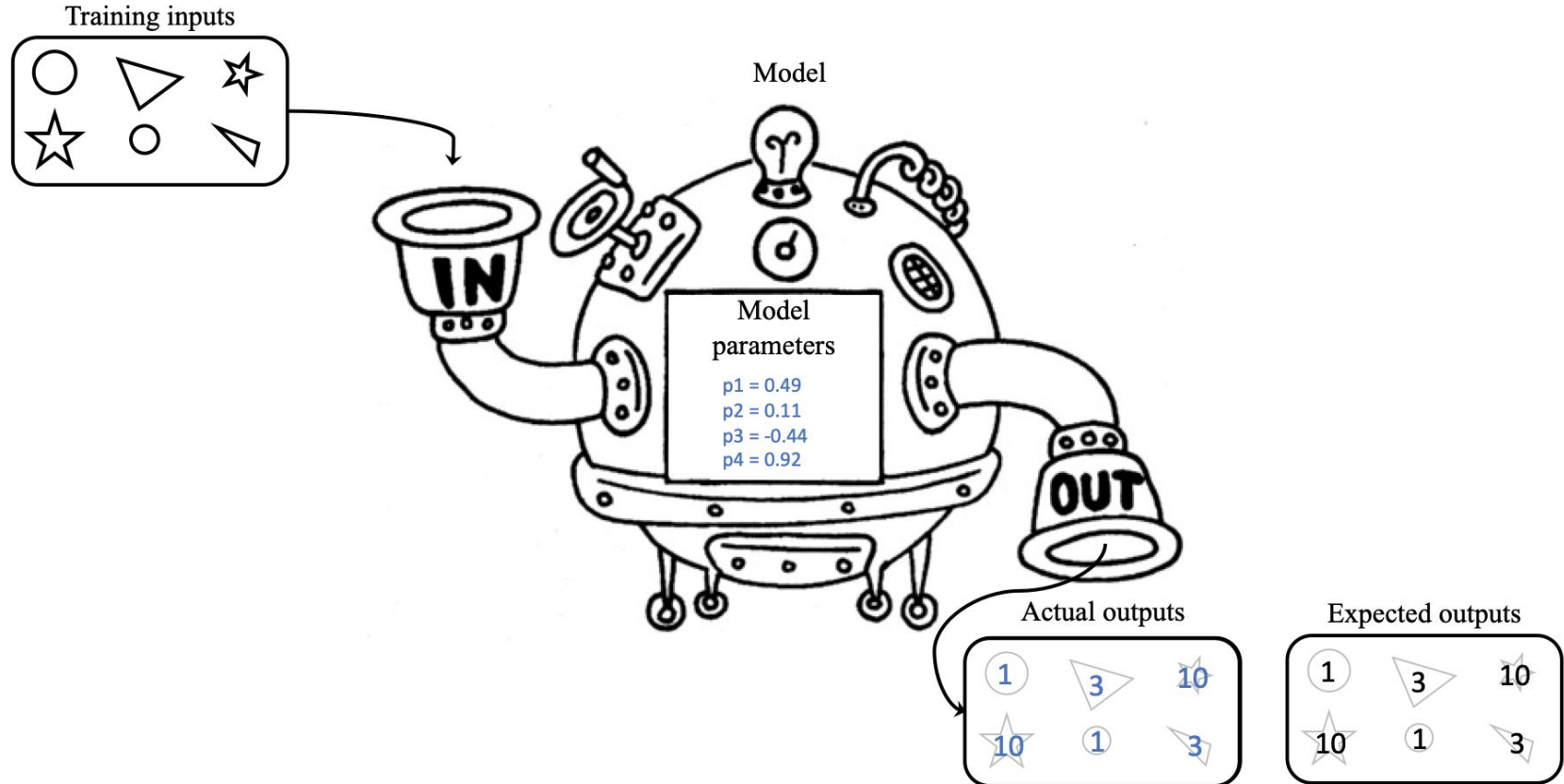The protein design problem

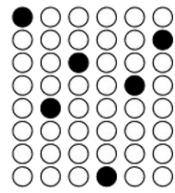The structural prediction problem

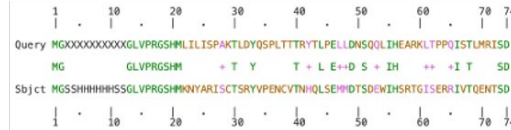Sequence

Structure

Function

# Supervised machine learning



Training inputs

Model

Model parameters

p1 = 0.49
p2 = 0.11
p3 = -0.44
p4 = 0.92

IN

OUT

Actual outputs

1    3    10
10    1    3

Expected outputs

1    3    10
10    1    3

# Common data representations for proteins in machine learning



Gao, Mahajan, Sulam & Gray *Patterns* 2020
https://doi.org/10.1016/j.patter.2020.100142

# AlphaFold model architecture consists of an MSA module (Evoformer) and a Structure module

# Protein language models for structural prediction

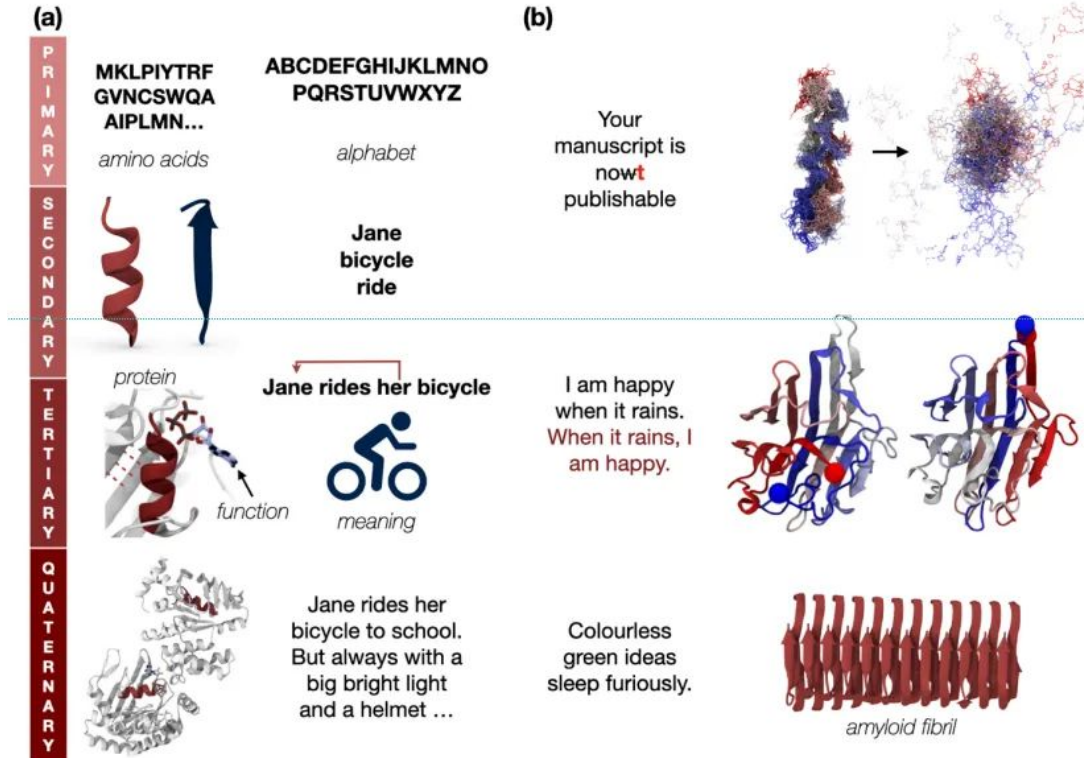# Protein language



**(a)**

PRIMARY — MKLPIYTRF GVNCSWQA AIPLMN… *amino acids* | ABCDEFGHIJKLMNO PQRSTUVWXYZ *alphabet*

SECONDARY — *protein* | Jane bicycle ride

TERTIARY — *function* | Jane rides her bicycle *meaning*

QUATERNARY — Jane rides her bicycle to school. But always with a big bright light and a helmet …

**(b)**

Your manuscript is now**t** publishable

I am happy when it rains. When it rains, I am happy.

Colourless green ideas sleep furiously.

*amyloid fibril*

https://doi.org/10.1038/s42256-022-00499-z

# ESMfold - a faster alternative to AF2

# ESMatlas

one-hot

Sequence

x

one-hot

Sequence

X

representation

Y

# How are protein language models trained?

# Unsupervised



Minimize difference between

# Masked language modeling (or self-supervised)

# "Masked language modeling" is an approximation of "Pseudolikelihood"



Mask

Mask

Minimize difference between

$$\mathcal{L}_{PL}(\theta; x) = \sum_{i=1}^{L} \log p_\theta(x_i | x_{\setminus i}) \qquad \mathcal{L}_{MLM}(\theta; x, M) = \sum_{i \in M} \log p_\theta(x_i | x_{\setminus M})$$

# Single-head model



Minimize difference between

# When pLM on a single protein family (MSA) we find:

- You only need **ONE** layer.
- You can replace the **positional embedding** with an **identity matrix** (encoding the exact positional information).
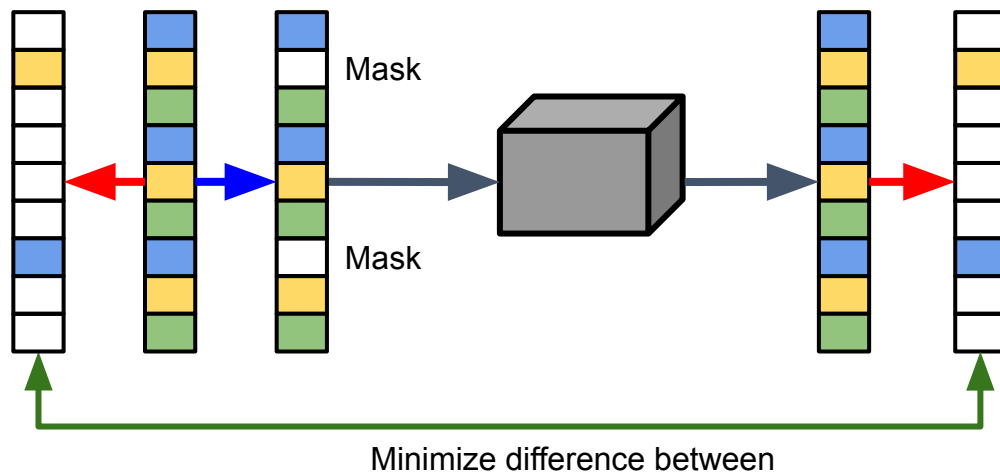- The weights of the **Query** and **Key** layers encode the contact map!



(**query** @ **key**)

Bhattacharya, N., Thomas, N., Rao, R., Dauparas, J., Koo, P.K., Baker, D., Song, Y.S. and Ovchinnikov, S., 2020. Single layers of attention suffice to predict protein contacts.

# Multi-head model



Minimize difference between

# Multi-layer/Multi-head model



Minimize difference between

# Extract attention matrices

GSHMPEEEKAARLFIEALEKGDPELMRKVISPDTRMEDNGREFTGDEVVEYVKEIQKRGEQWHLRRYTKEGNSWRFEVQVDNNGQTEQWEVQIEVRNGRIKRVTITHV
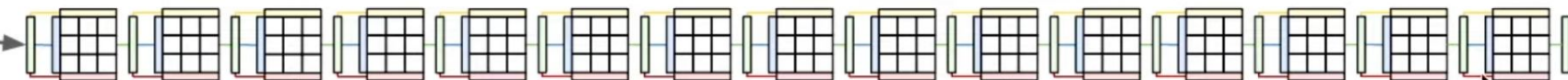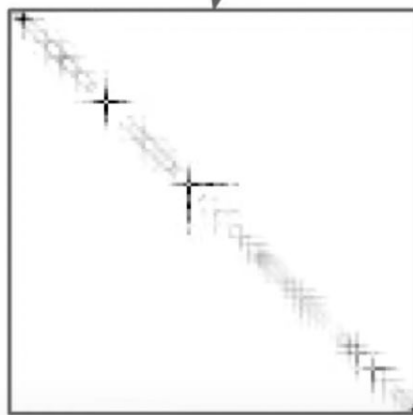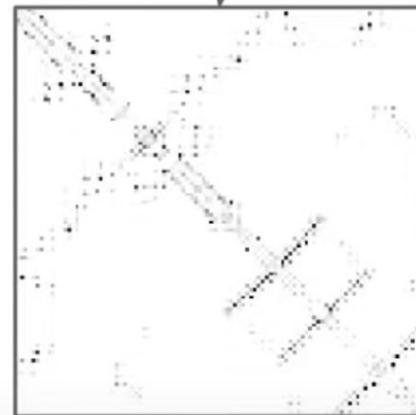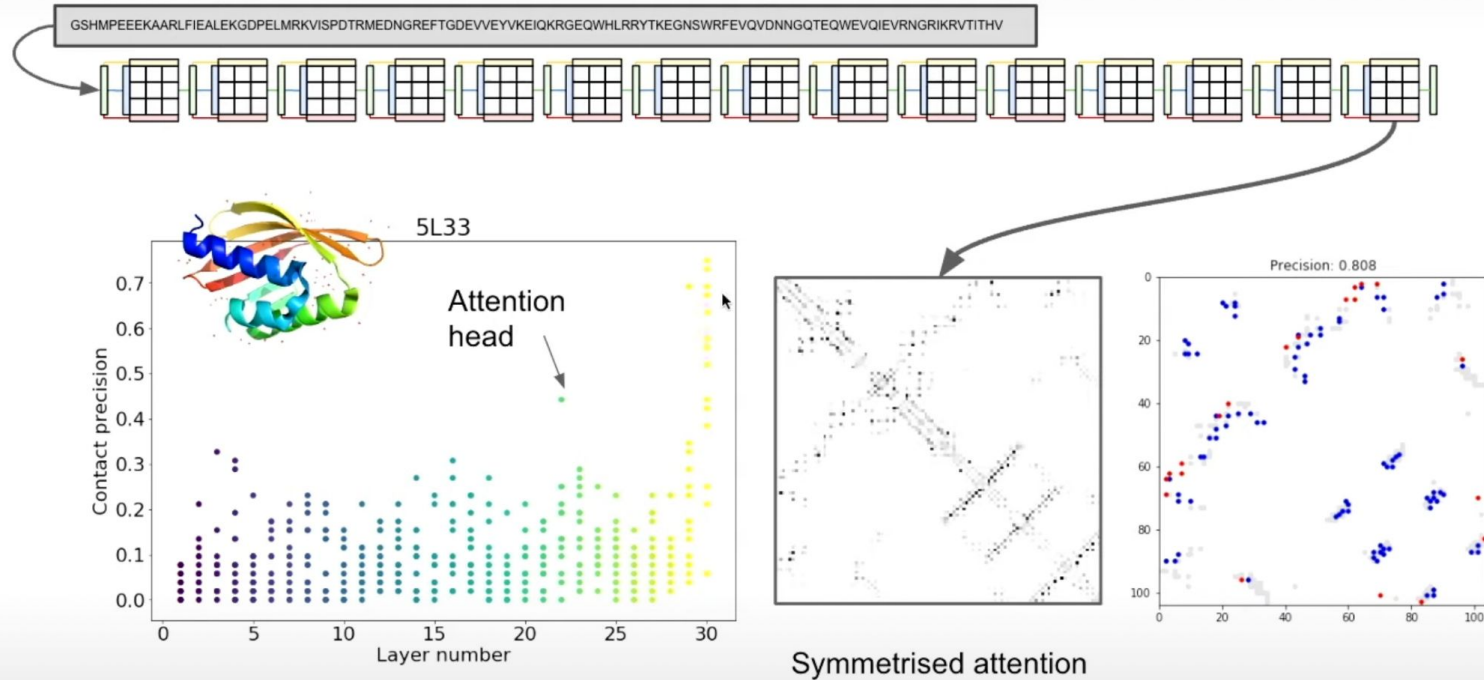


Symmetrised attention

Symmetrised attention

Symmetrised attention
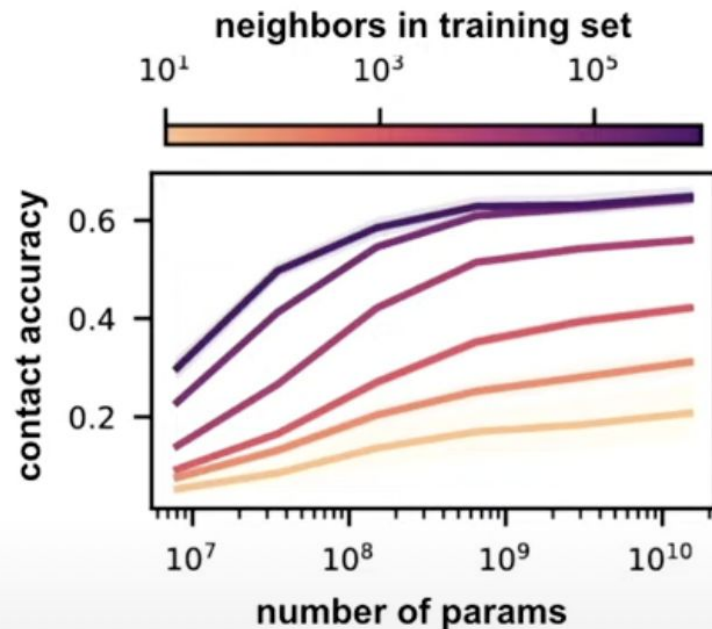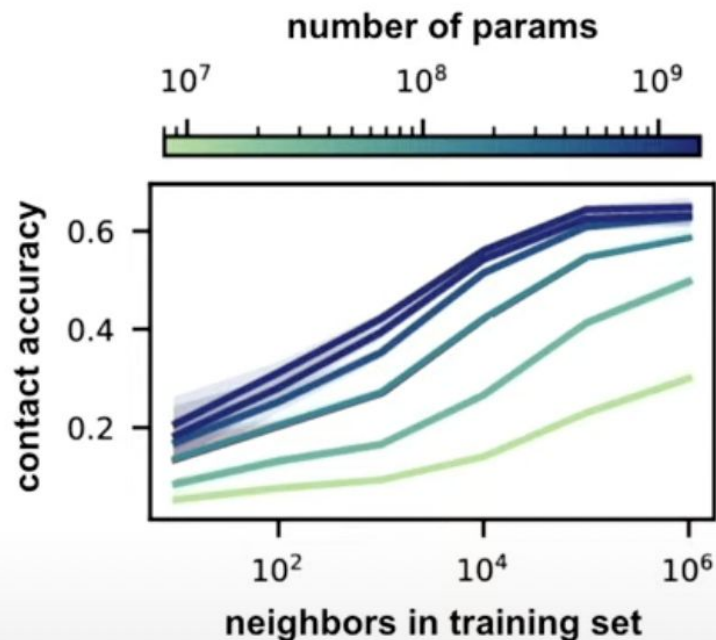
# Layers towards the end tend to capture contacts!

GSHMPEEEKAARLFIEALEKGDPELMRKVISPDTRMEDNGREFTGDEVVEYVKEIQKRGEQWHLRRYTKEGNSWRFEVQVDNNGQTEQWEVQIEVRNGRIKRVTITHV



5L33

Attention head

Contact precision

Layer number

Symmetrised attention

Precision: 0.808

Vig, Jesse, et al. "**BERTology Meets Biology: Interpreting Attention in Protein Language Models.**" (2020).
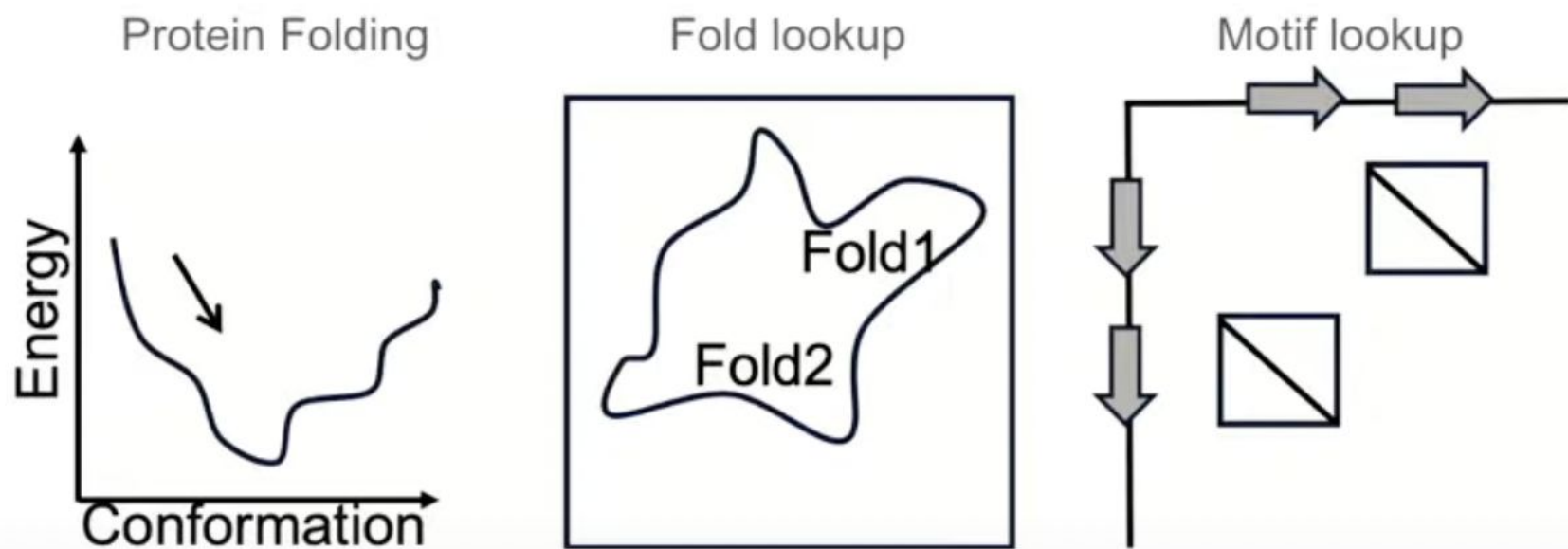
Bhattacharya, Nick, et al. **"Single Layers of Attention Suffice to Predict Protein Contacts"** (2020)

Rao, Roshan, et al. **"Transformer protein language models are unsupervised structure learners."** (2020)
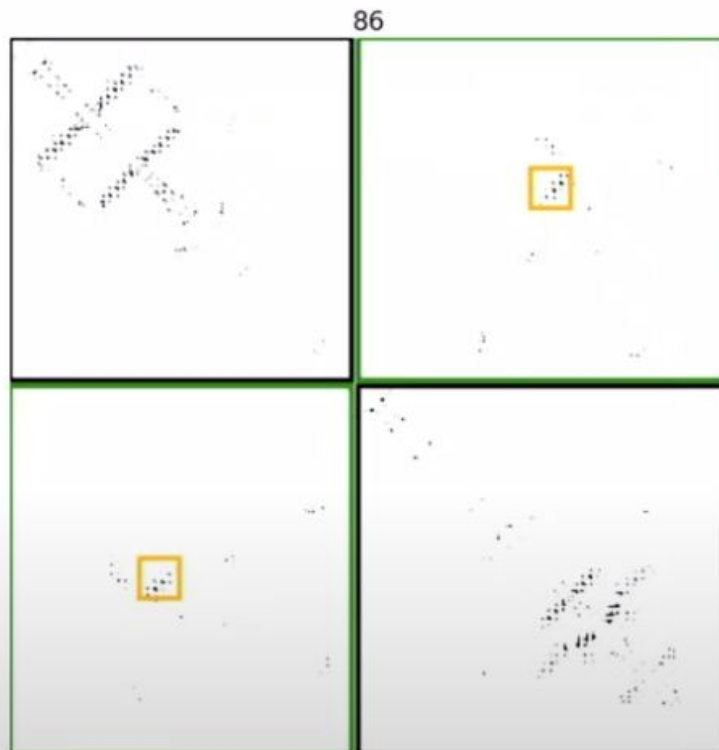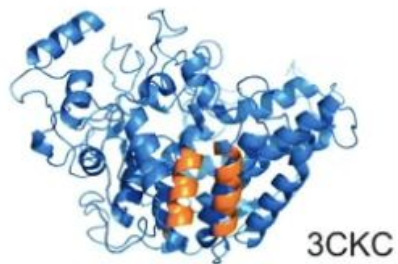
# ESM2: train models with different number of params



Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W., dos Santos Costa, A., Fazel-Zarandi, M., Sercu, T., Candido, S. and Rives, A., 2022. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*.

# How do protein language models "store" coevolution statistics?
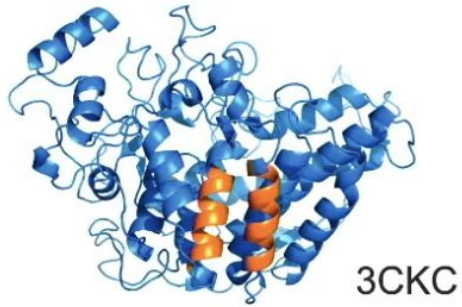


Protein Folding      Fold lookup      Motif lookup

# Masking majority of the sequence recovers the motif

SusD starch-
binding protein

3CKC

86

# Masking majority of the sequence recovers the motif
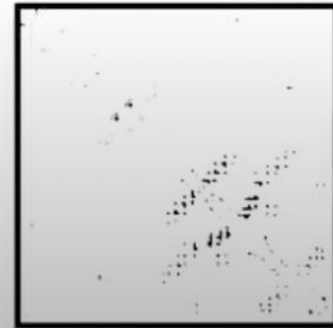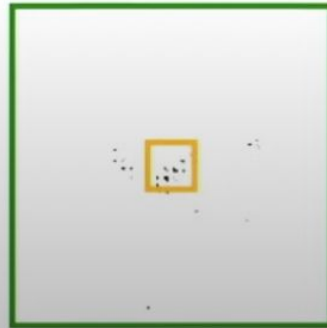
SusD starch-
binding protein



3CKC

58

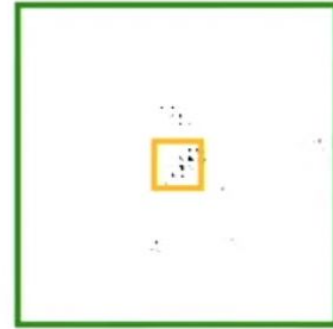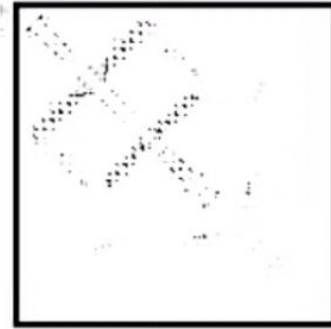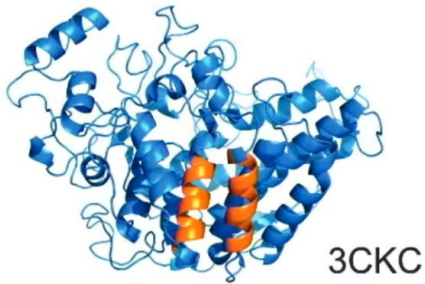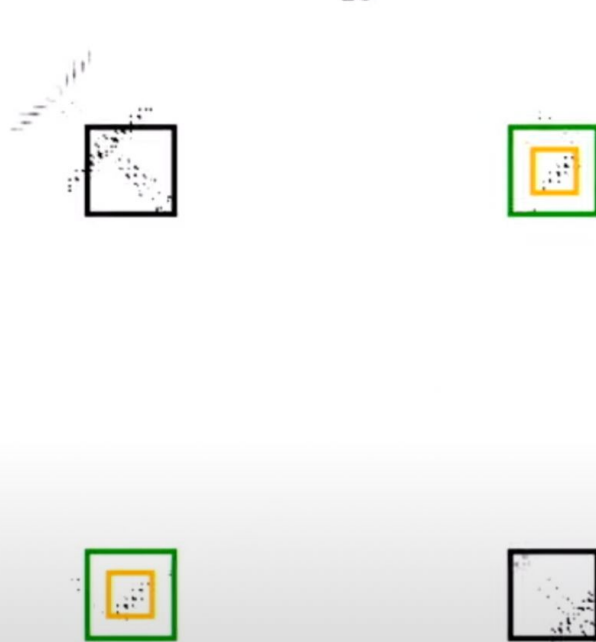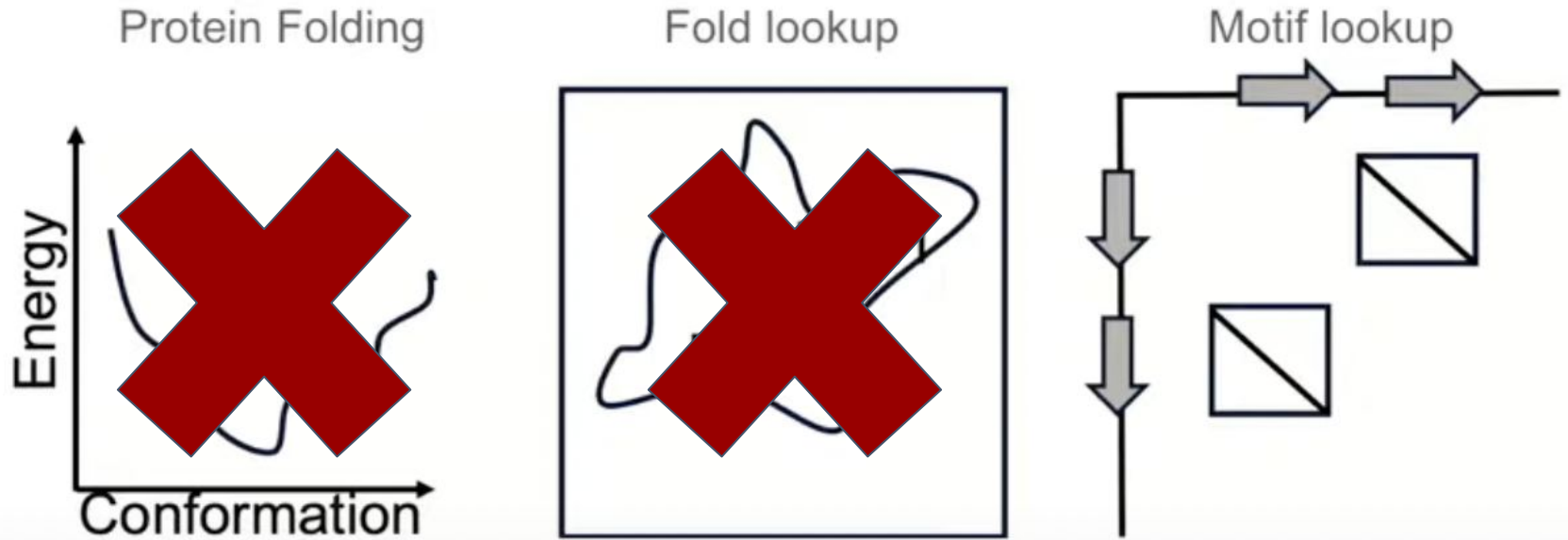# Masking majority of the sequence recovers the motif



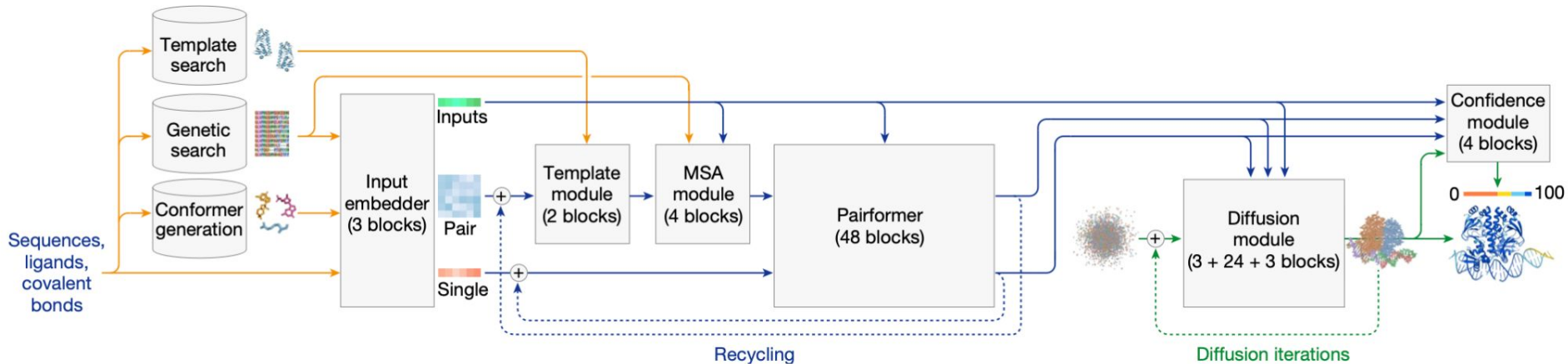SusD starch-binding protein

3CKC

10

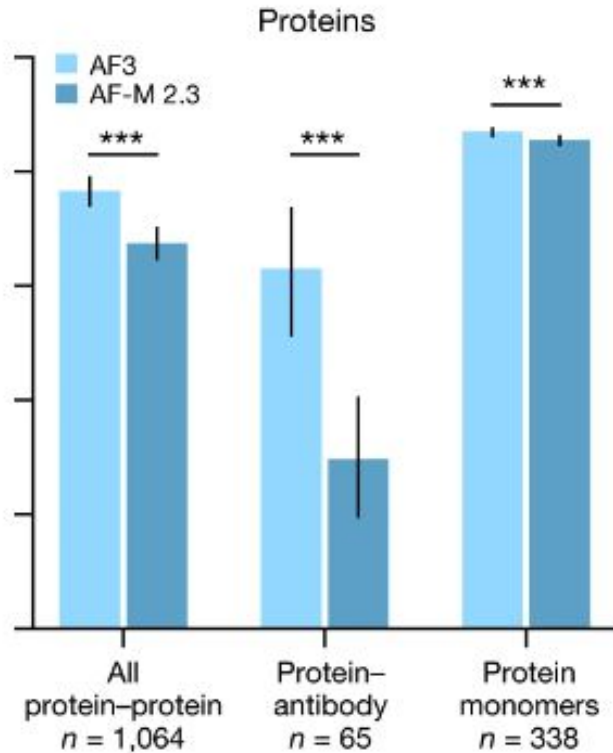# How do protein language models "store" coevolution statistics?

# The field continuously changes...

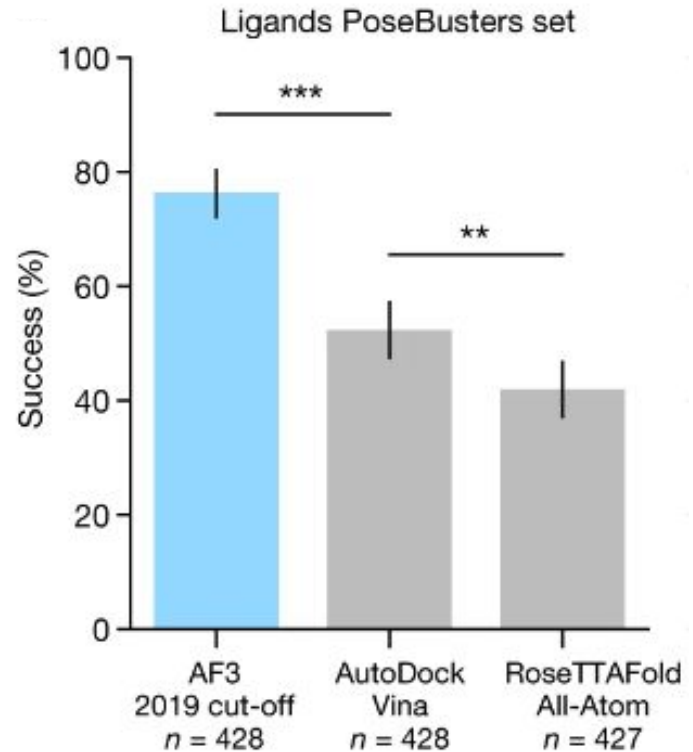# AF3 - making a deterministic problem not so...

# AF3 - making a deterministic problem not so...



Proteins

AF3
AF-M 2.3

***          ***          ***

All
protein–protein
n = 1,064

Protein–
antibody
n = 65

Protein
monomers
n = 338

# AF3 - making a deterministic problem not so…

https://www.nature.com/articles/s41586-024-07487-w

# ESM3 - combining language models with structural and functional information

# ESM3 - combining language models with structural and functional information



https://www.biorxiv.org/content/10.1101/2024.07.01.600583v1