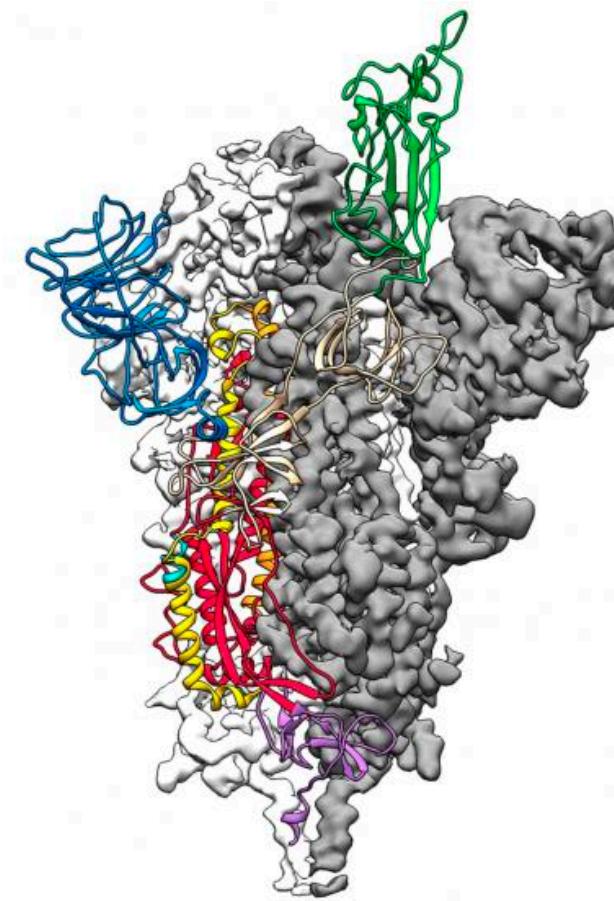


# Multimodal deep learning for protein engineering

Kevin Kaichuang Yang  
Microsoft Research New England  
 @KevinKaichuang

# Proteins are biology's actuators

- Human cells contain 1-3 billion proteins each
- Structure, metabolism, and signaling

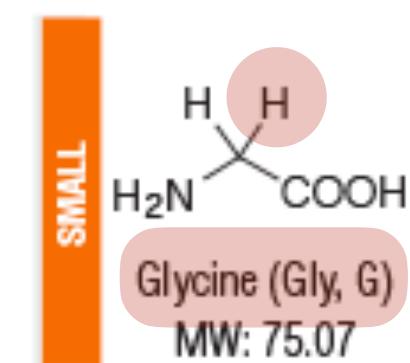
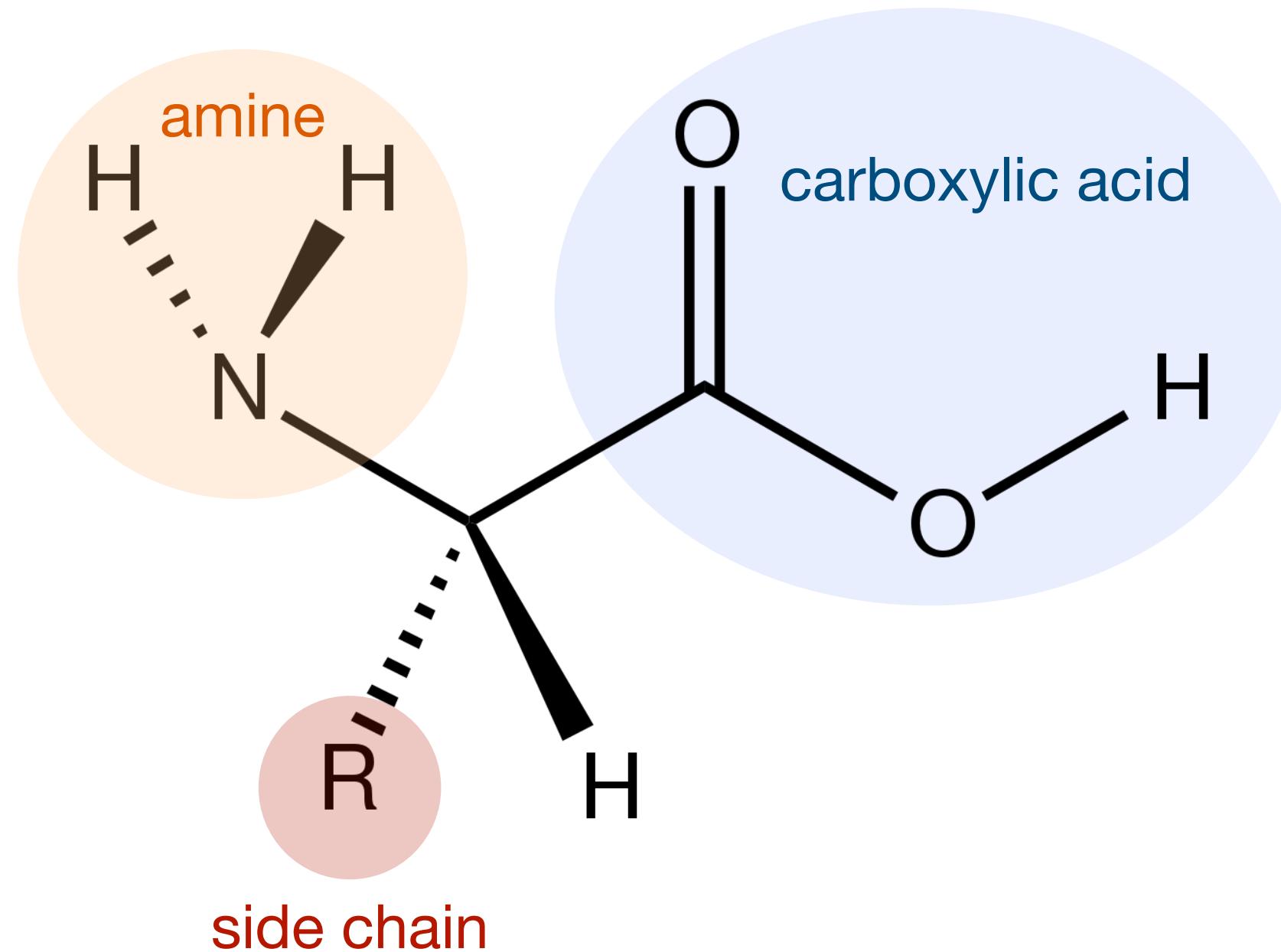


coronavirus spike protein

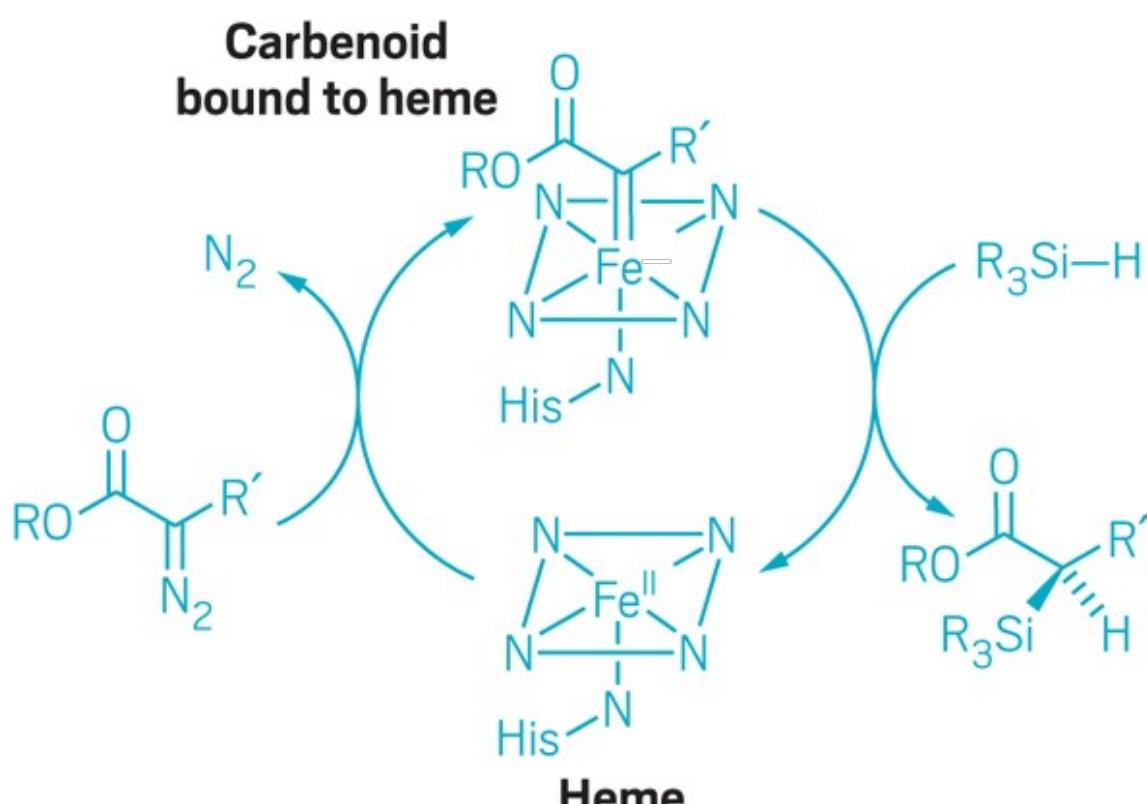


luciferase

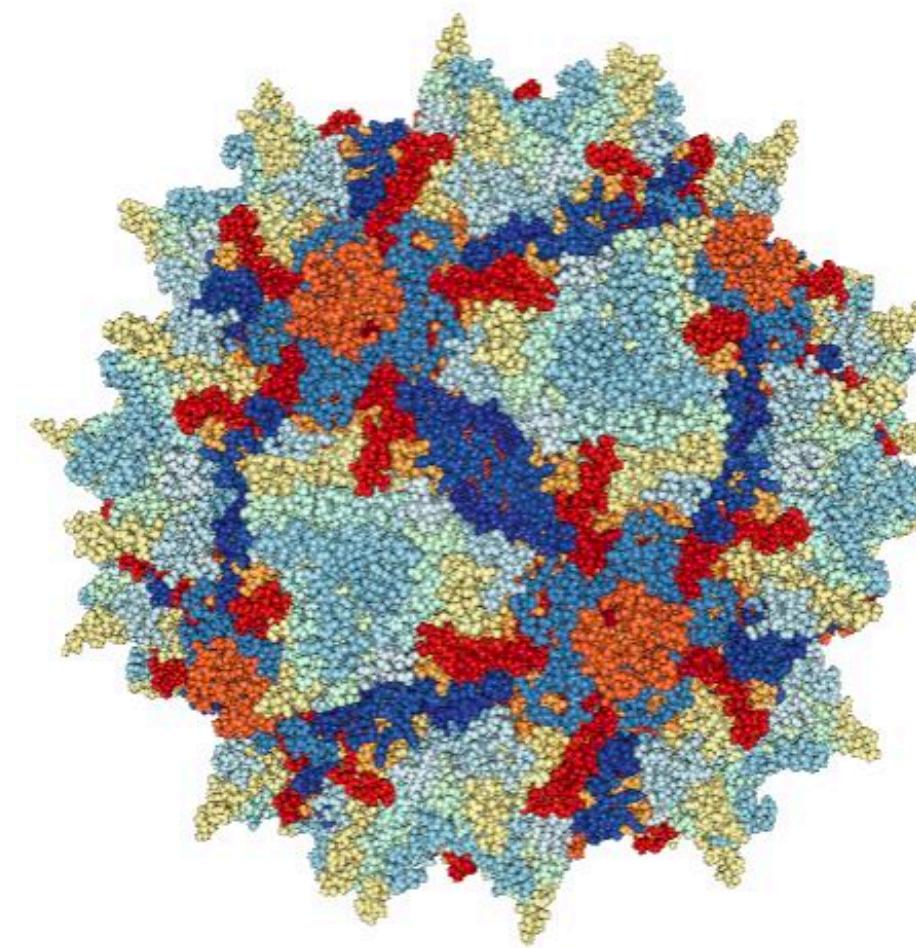
# Diversity arises from 20 building blocks



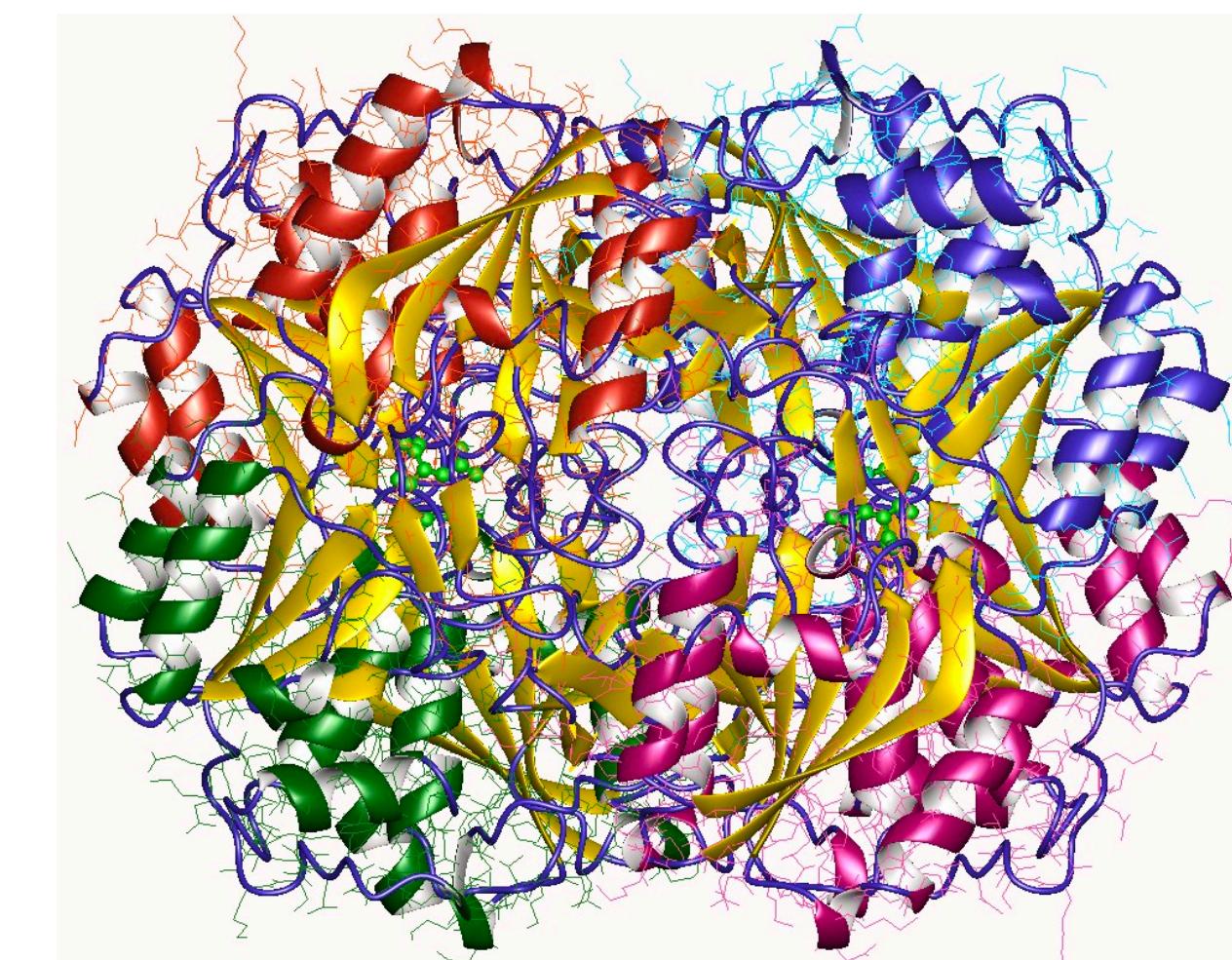
# Why design proteins?



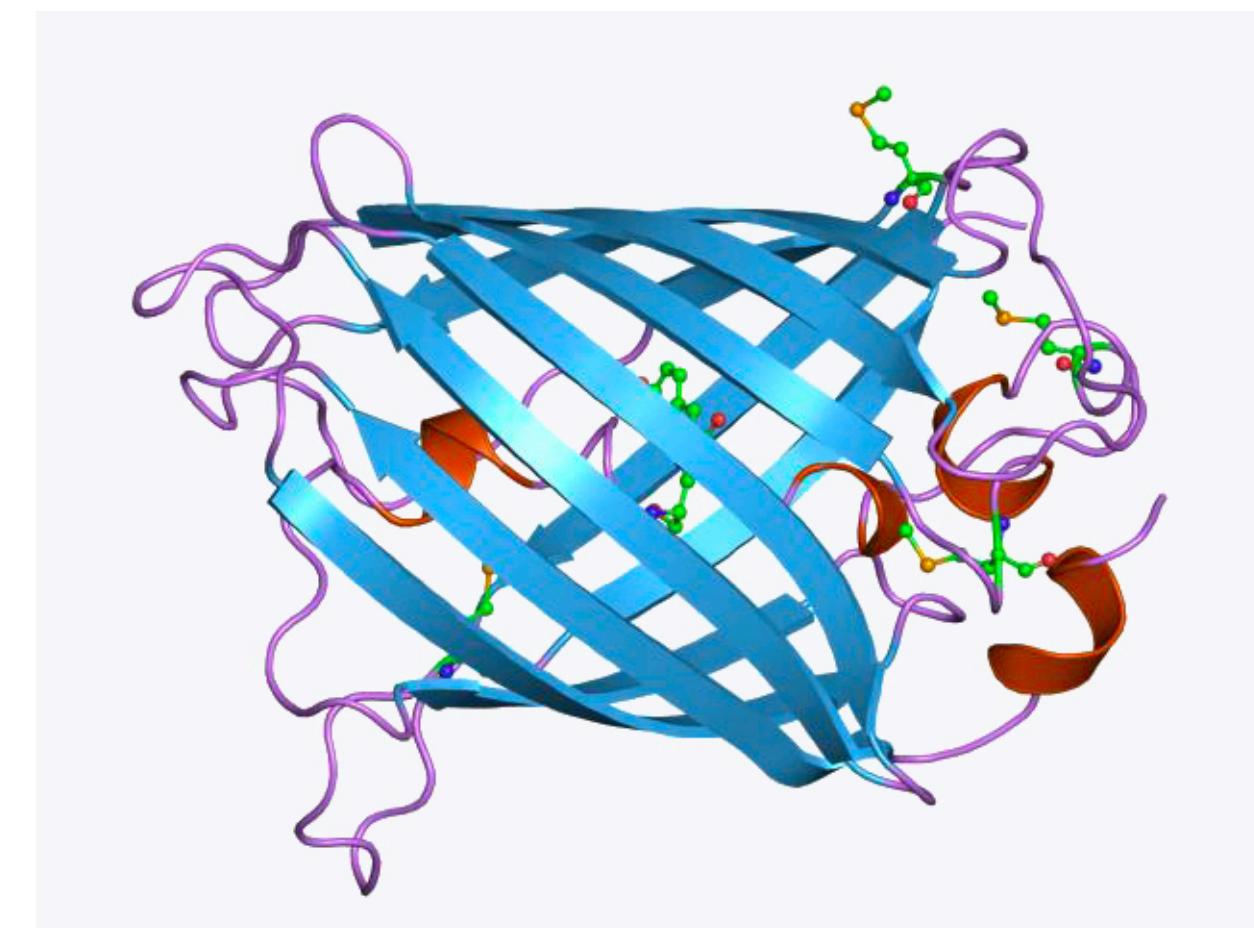
new chemistry



understand function

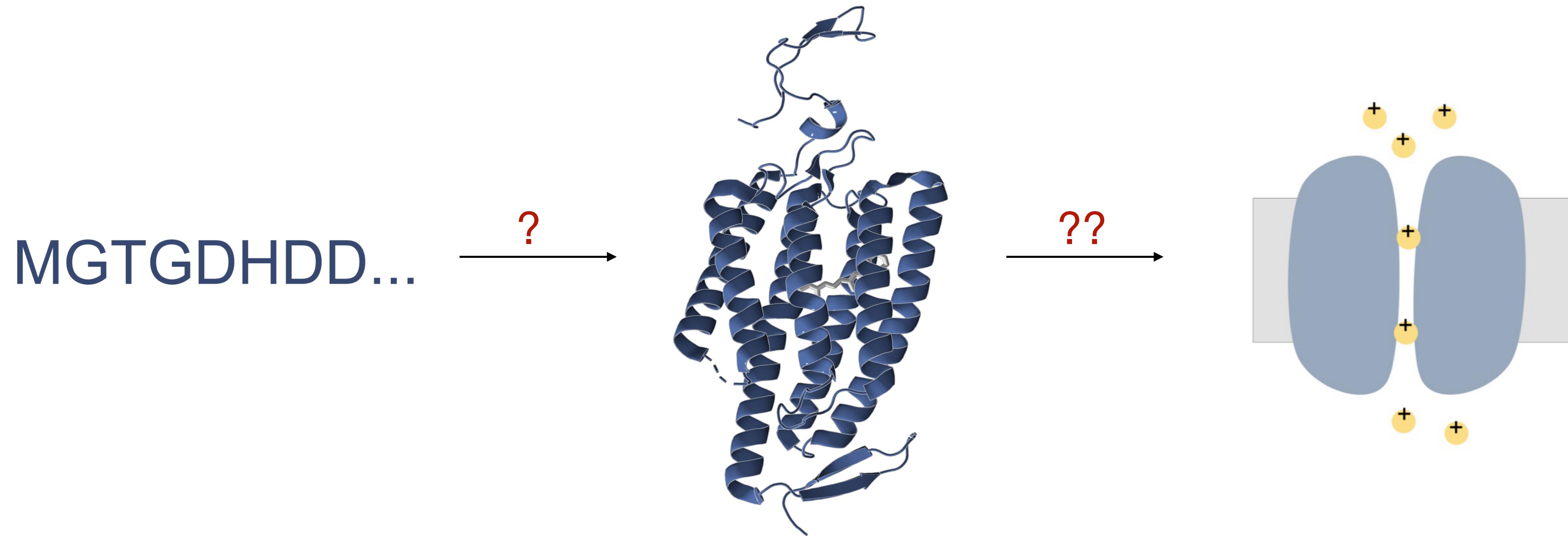


therapeutics



molecular tools

# The protein engineering problem



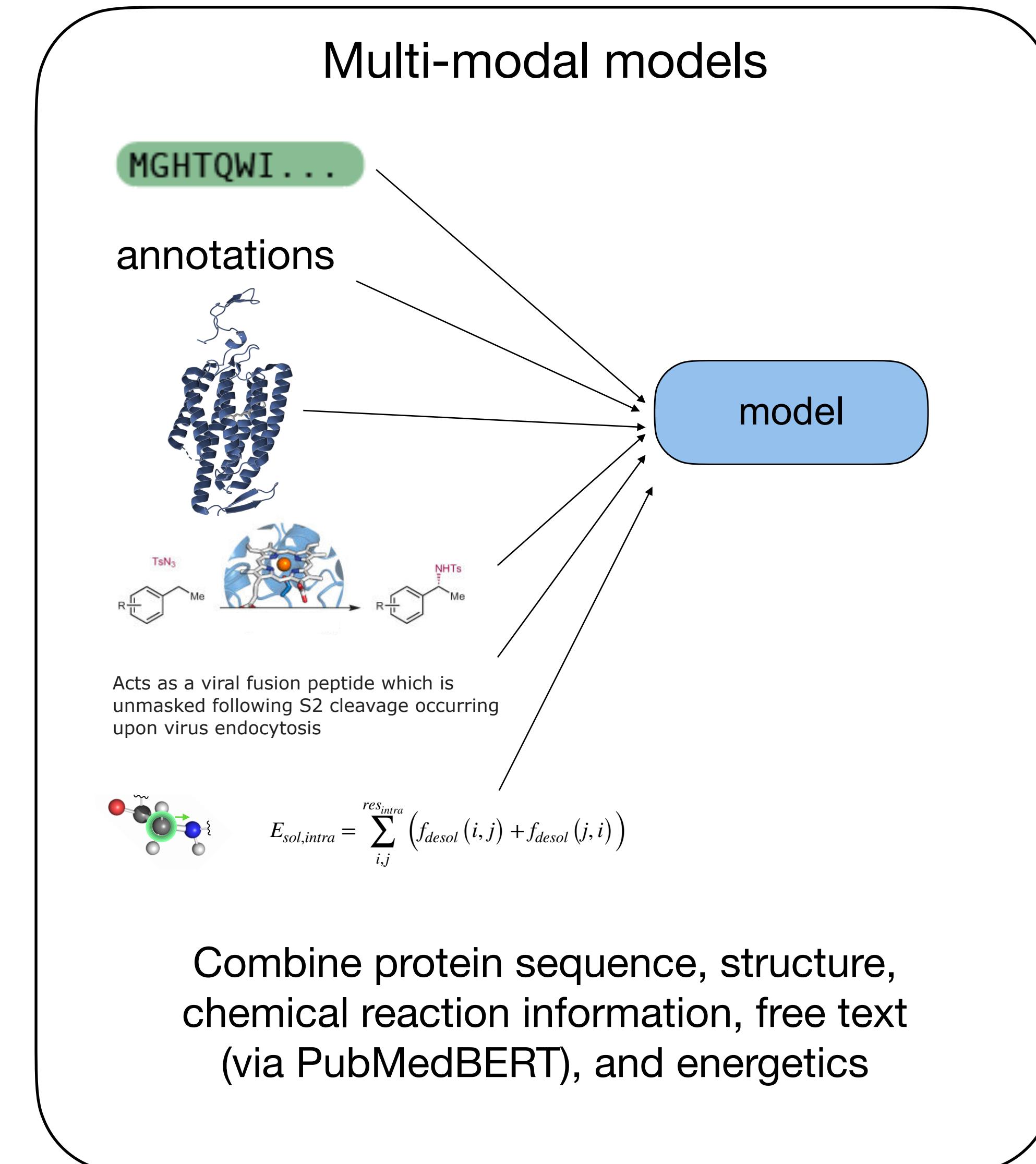
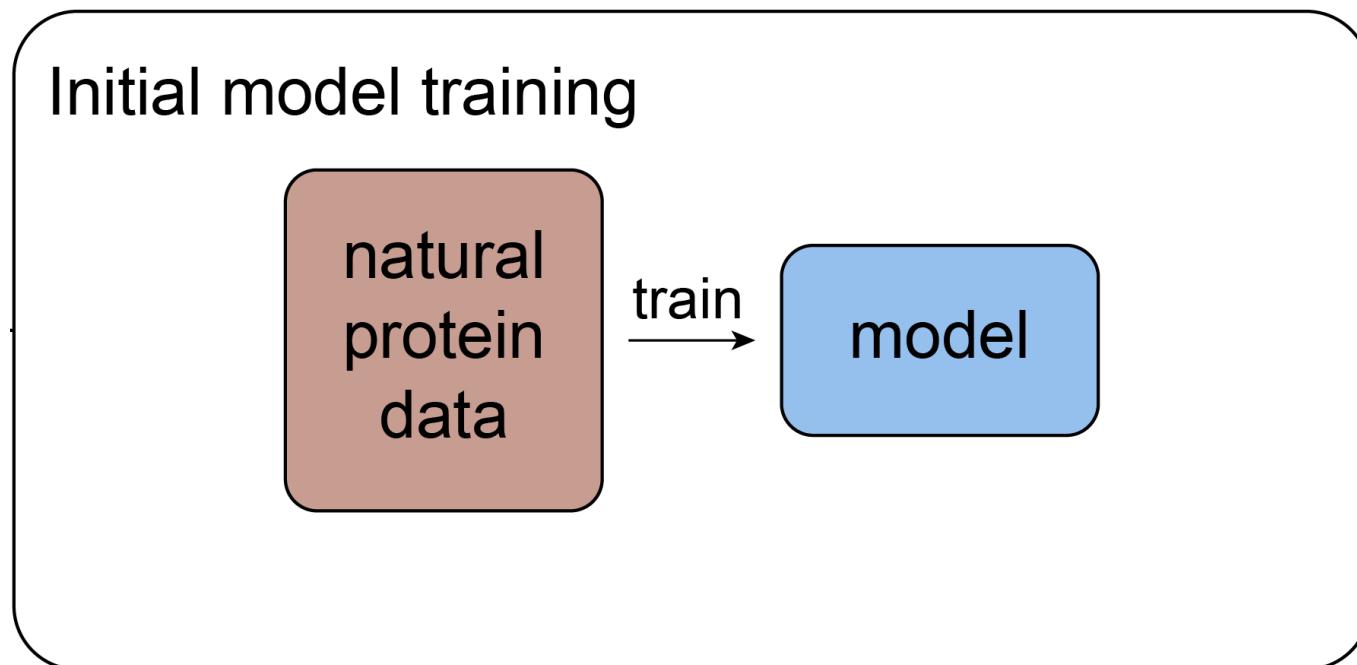
What sequence will give the desired function?

# The protein engineering problem

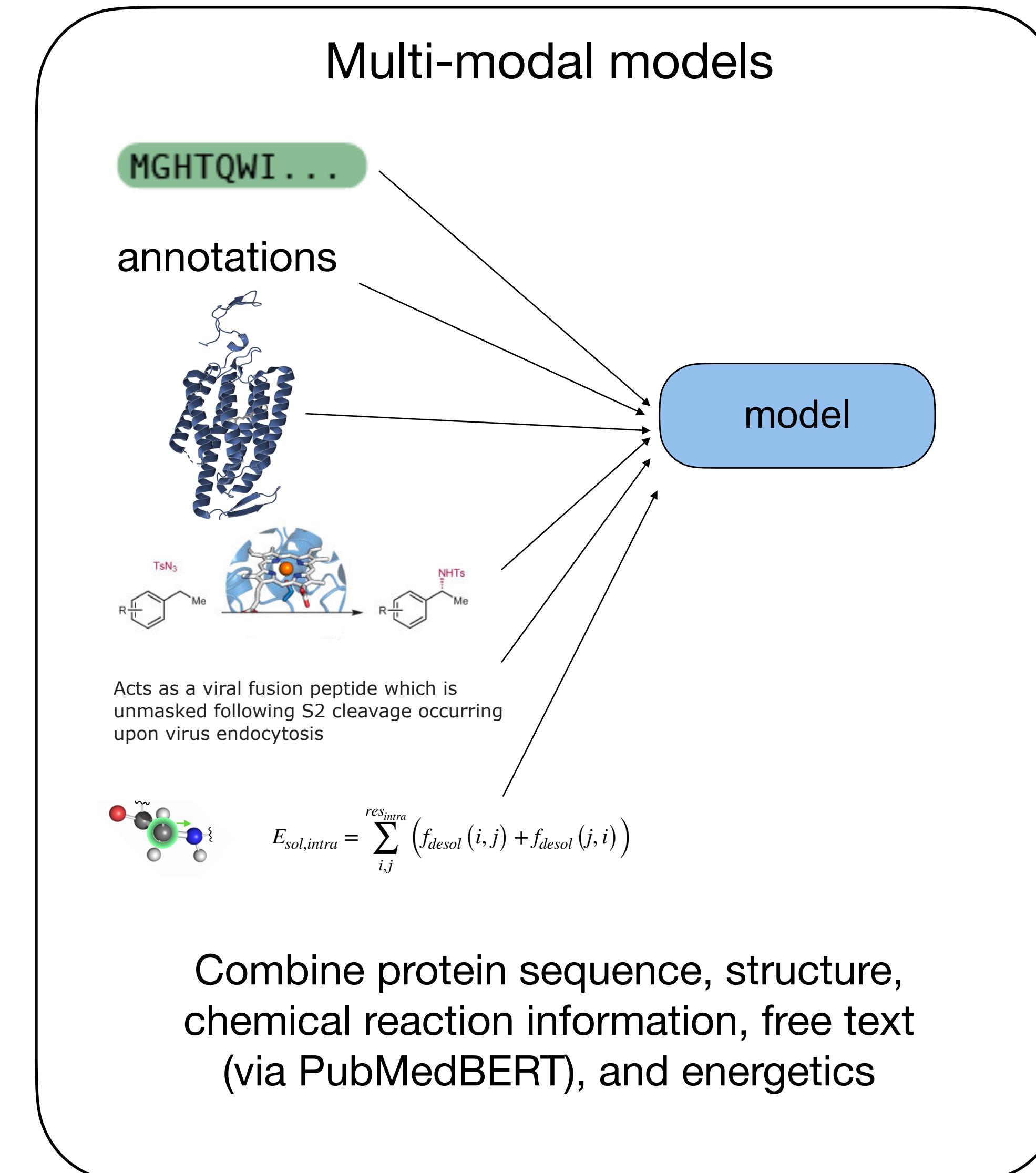
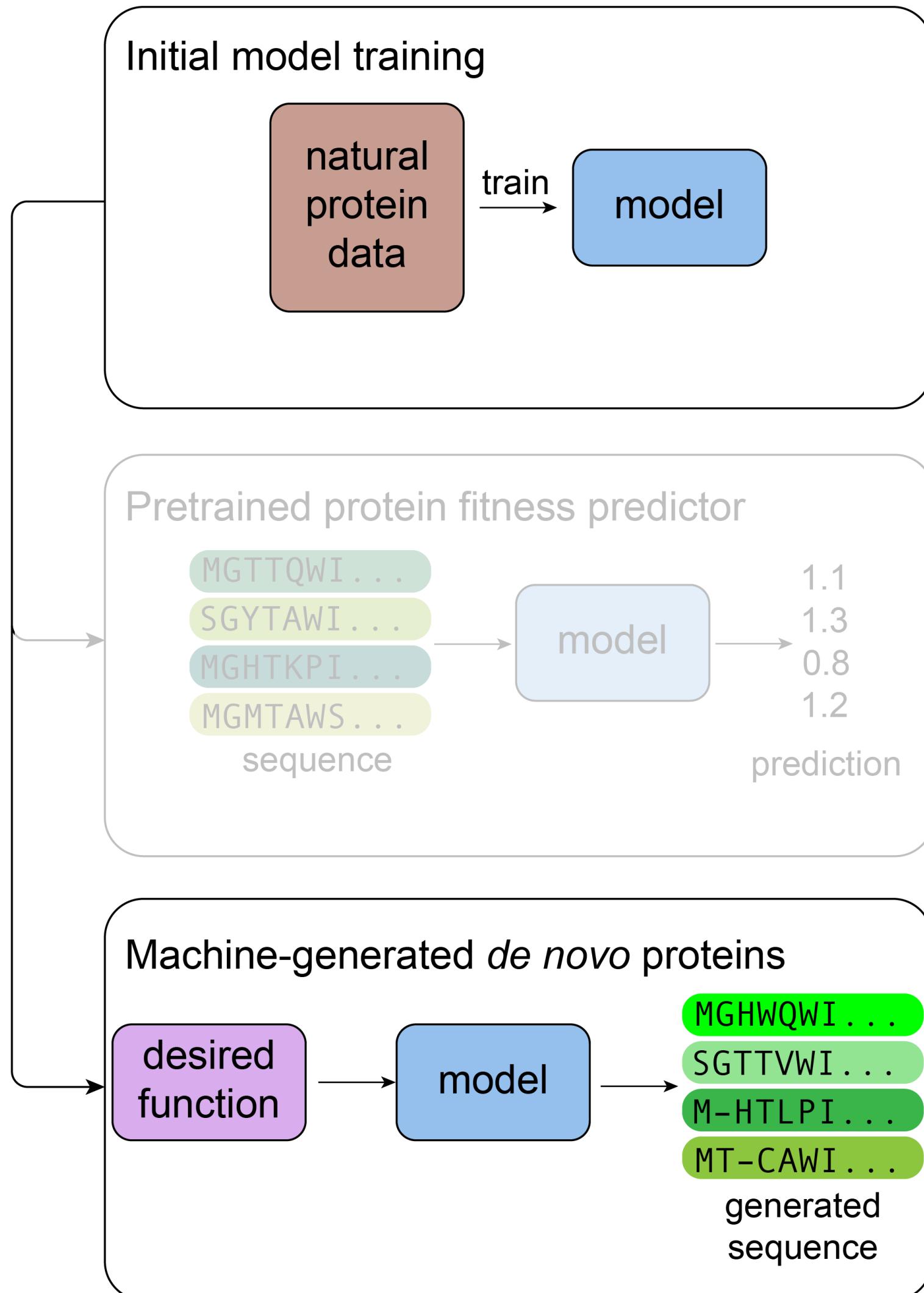


What sequence will give the desired function?

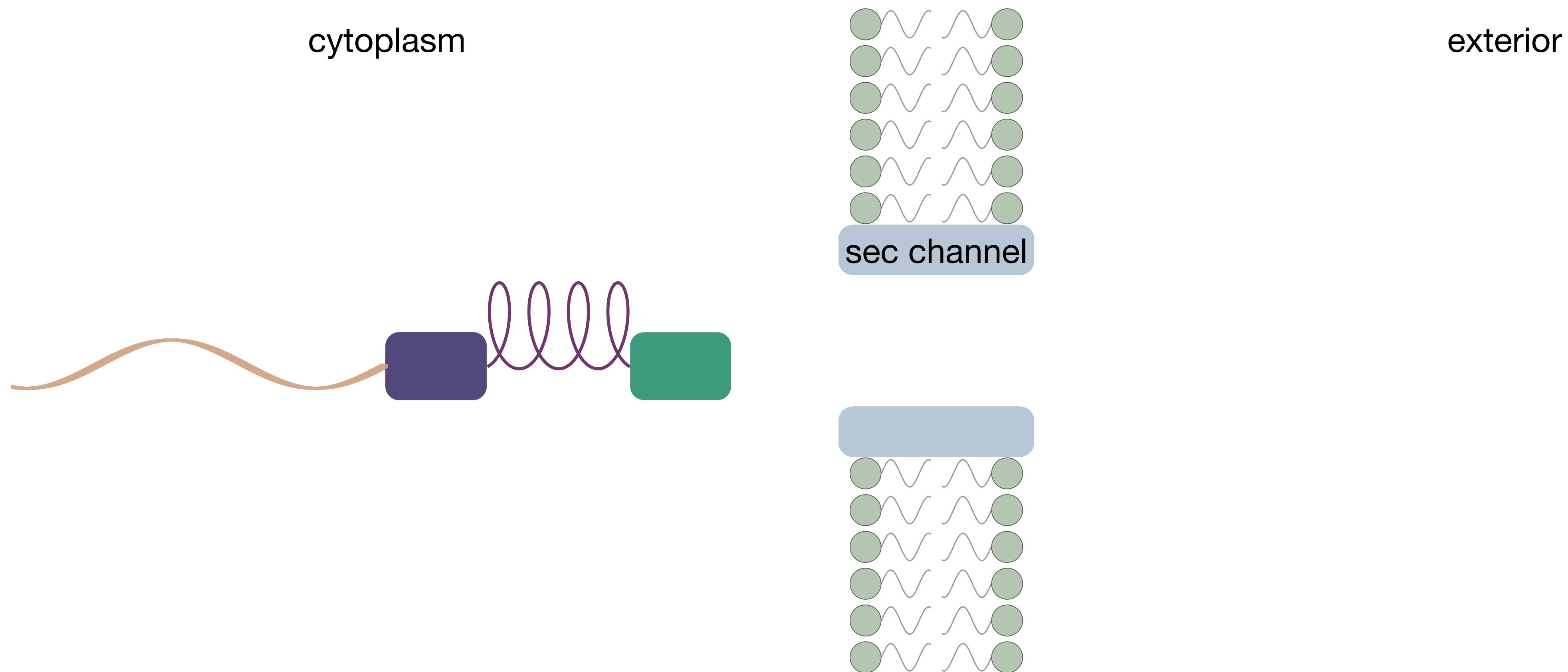
# Use multiple data modalities to discover and design proteins



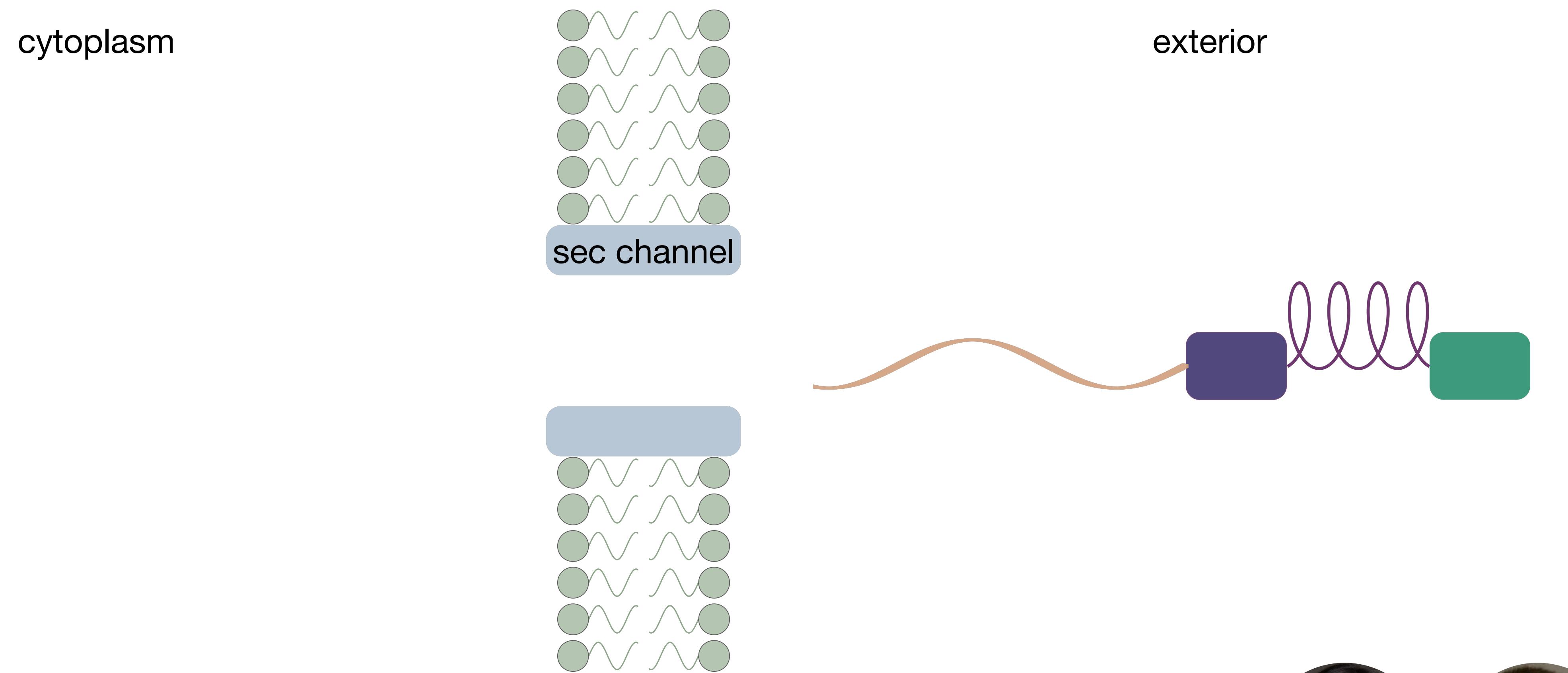
# Use multiple data modalities to discover proteins



# Signal peptides are secretion signals



# Signal peptides are secretion signals

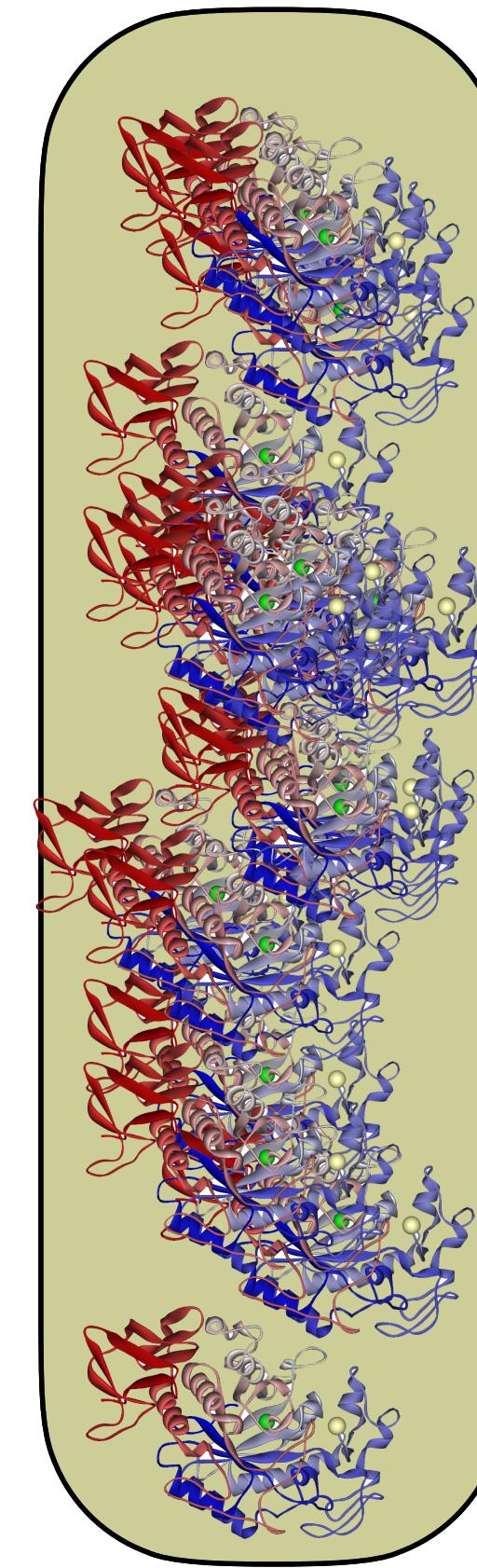


Need: custom SPs for arbitrary proteins

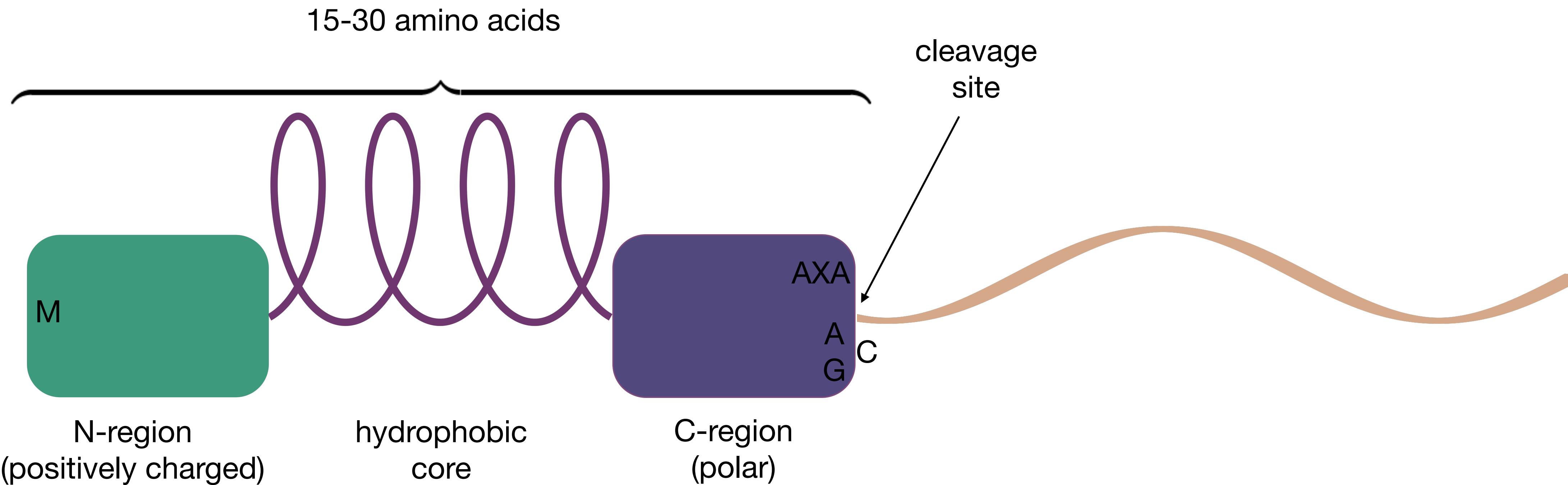


# Signal peptides simplify protein production

- Less burden on cells
- Easier separation

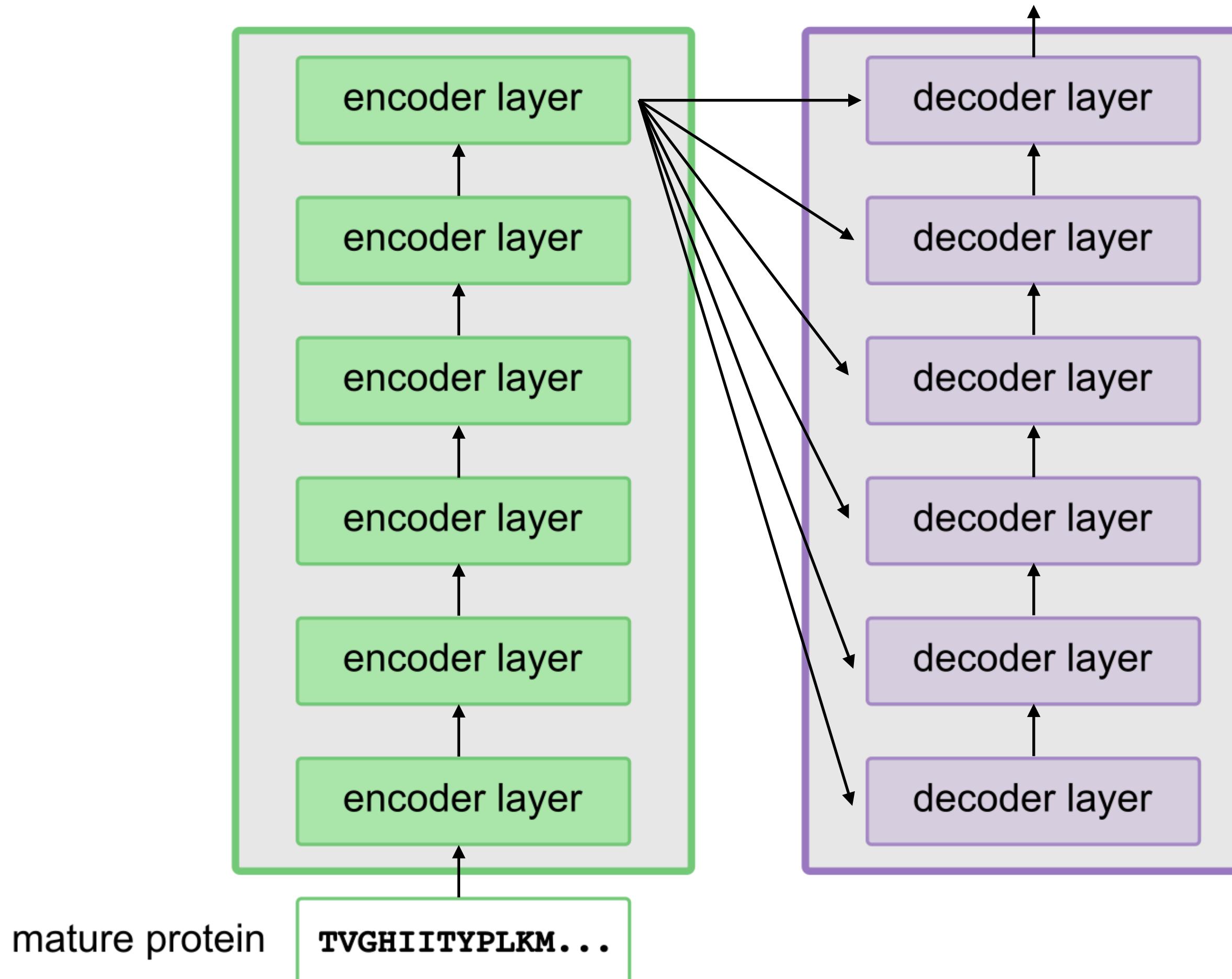


# Signal peptides are structured



- Extra-cellular secretion in prokaryotes
- Secretion from endoplasmic reticulum in eukaryotes
- SP dependent on species and protein
- Rules insufficient for generation

# Machine translation for SP generation



- Transformer architecture
- 40k examples from UniProt
- 40% test accuracy

Actual MNFSKIFIFVVLAVLILCSQTEA  
Predicted MNFKLIFLVALVLMAAFLGQTEG

Actual MLRKSAVLAGVVLLGASAQA  
Predicted MRLSTSALVLGAASSAV

Actual MASRVECLVVCLLVLRAQSSG  
Predicted MNLASVLLLAACHLSVSVNG

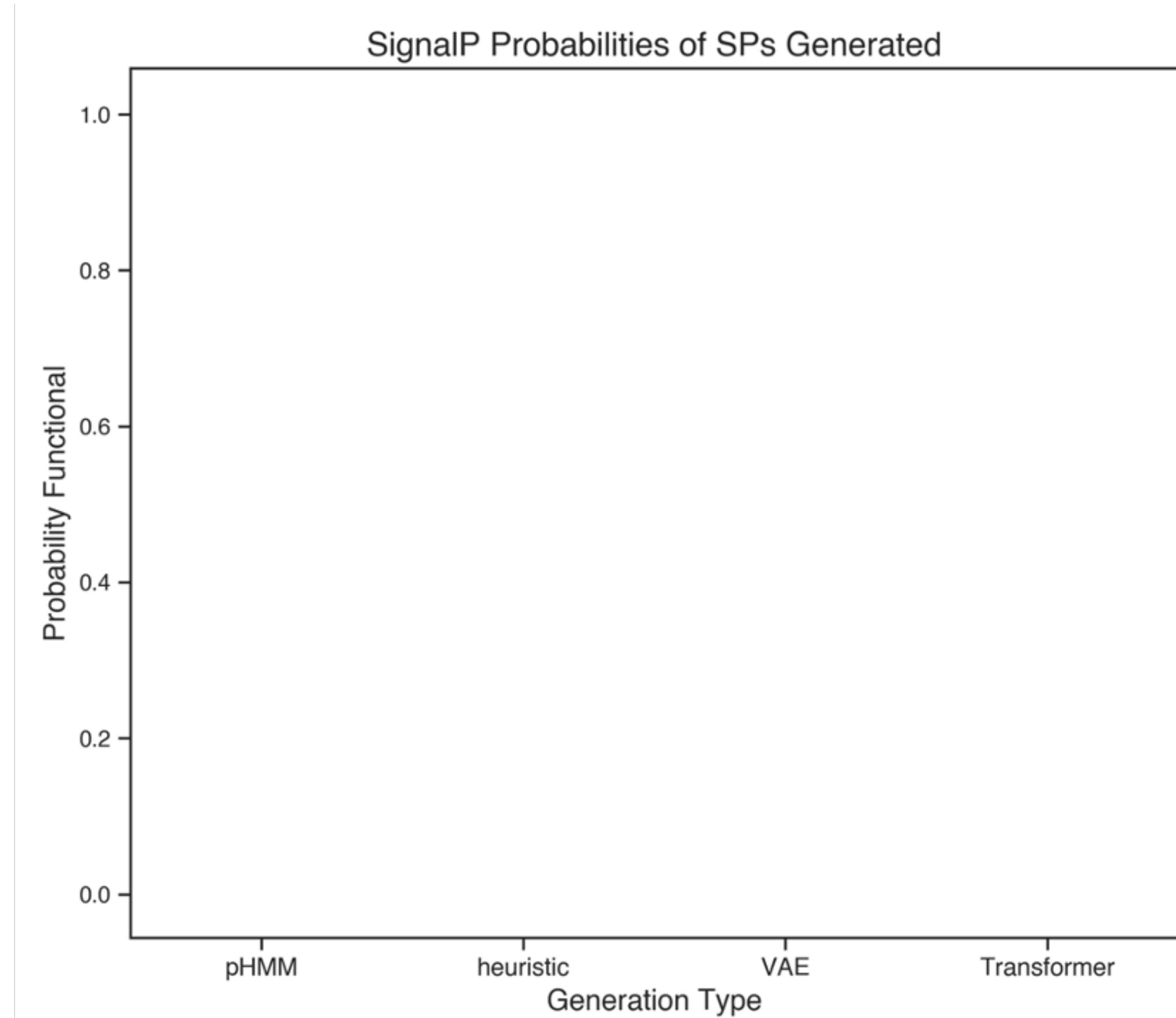
Actual MLWIRAVLPLVAGSDT  
Predicted MLGIWTLLPLVLTYVVRLLSKCVNA

# Generated SPs trick SignalP 5.0

Method	Mean SignalP Probability
Hidden Markov Model	16% $\pm$ 19%
heuristic	71% $\pm$ 25%
<b>Variational Autoencoder (VAE)</b>	<b>92% <math>\pm</math> 15%</b>
Transformer	90% $\pm$ 17%

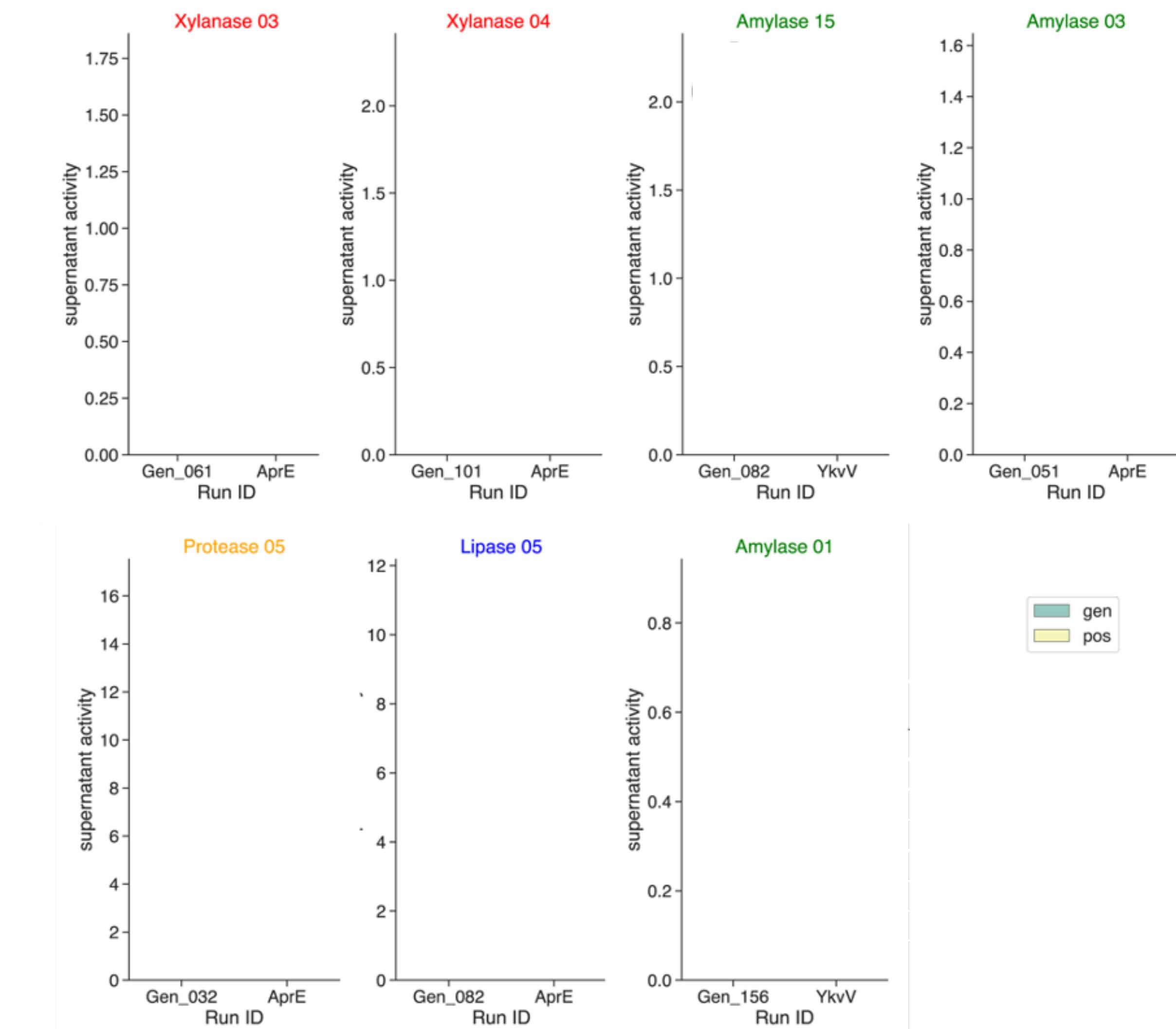
VAE sequences are repetitive

MRKRLALALALAALSLLLLSFGVKALAGSGA  
MRLLLLLLVLVLLAAPAPPGLS  
MKLLLLLVTLTTSVLALQA  
MRLLLLALLAAAAVALASA  
MASSSSSLFVVVLAVLLLLLTLSSA

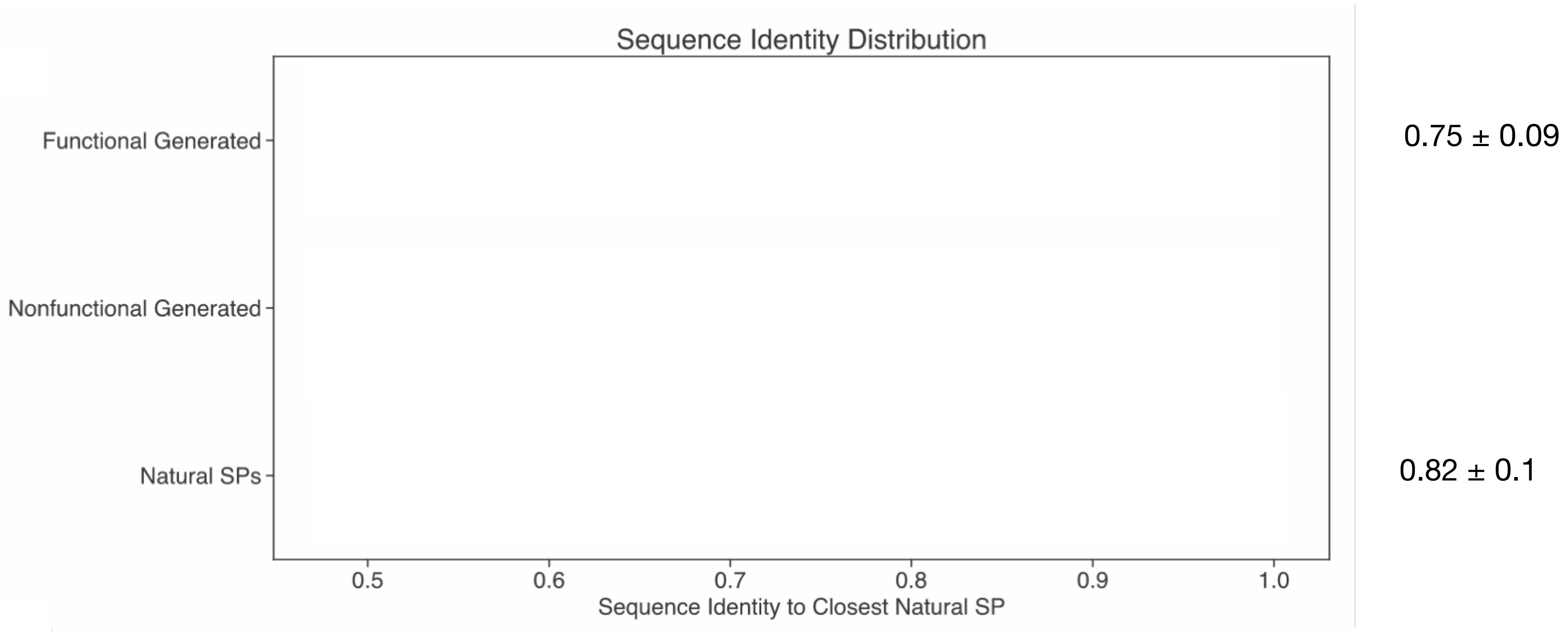


# Generated SPs are functional in *Bacillus subtilis*

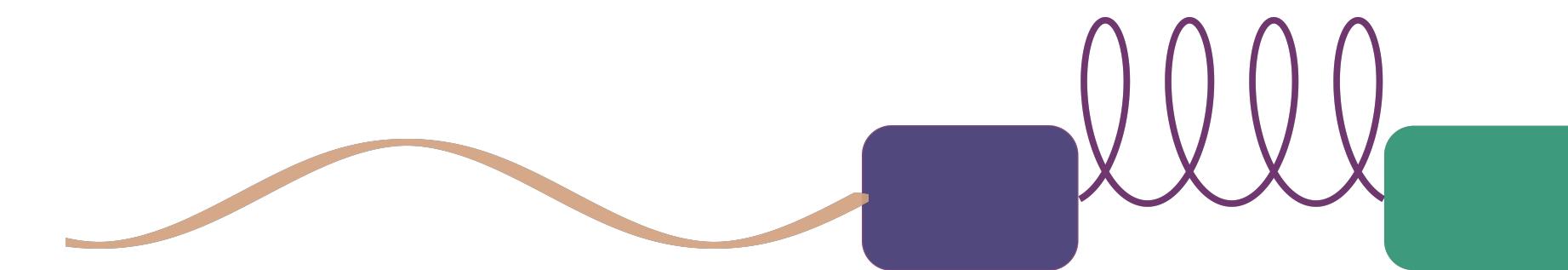
- 79% of natural SPs are functional
- 48% of generated SPs are functional
- Enzymatic activity is comparable



# Generated functional SPs are novel and diverse



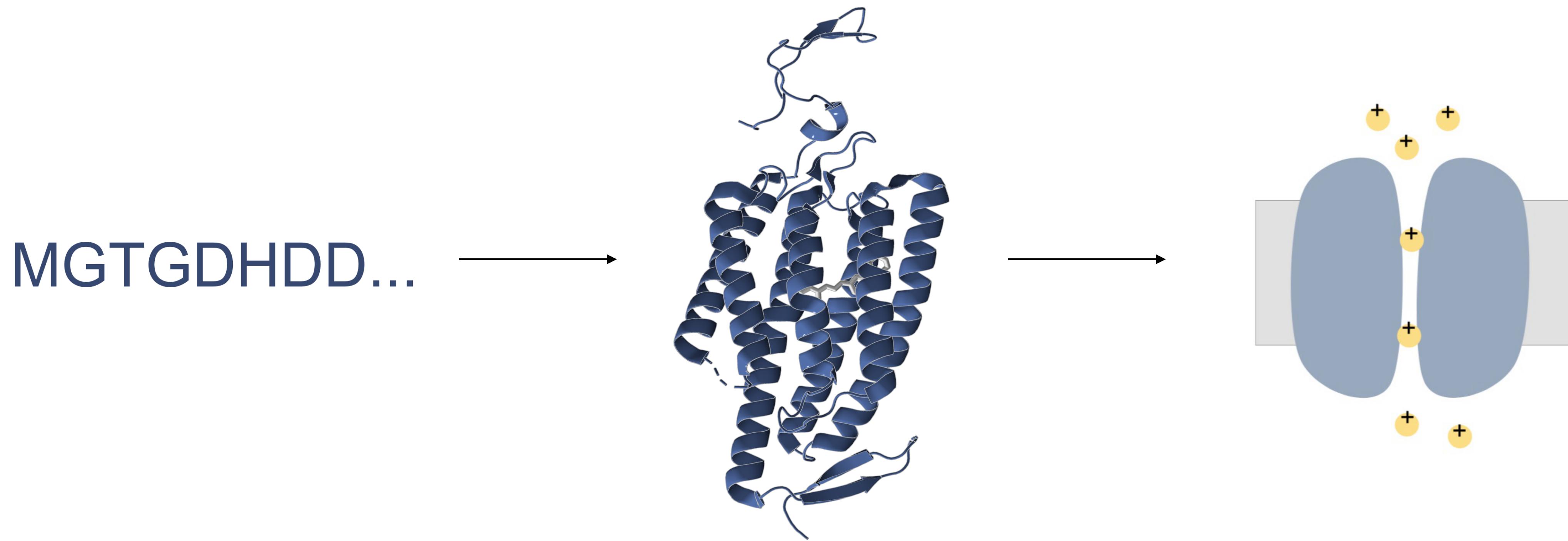
# SP generation conclusions



- Can generate *de novo* signal peptides
- That look like real SPs
- And result in functional secreted enzymes

How about entire proteins?

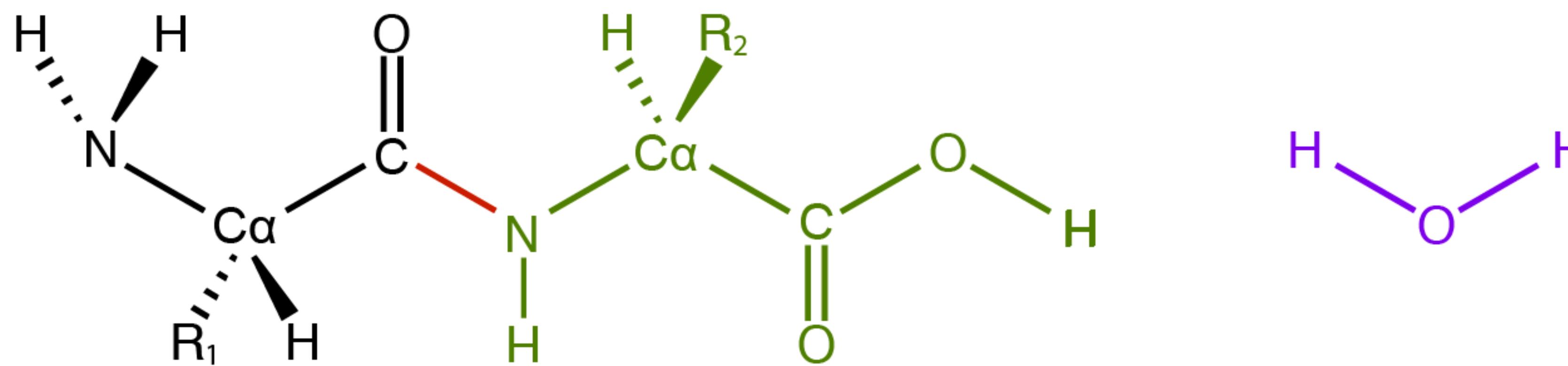
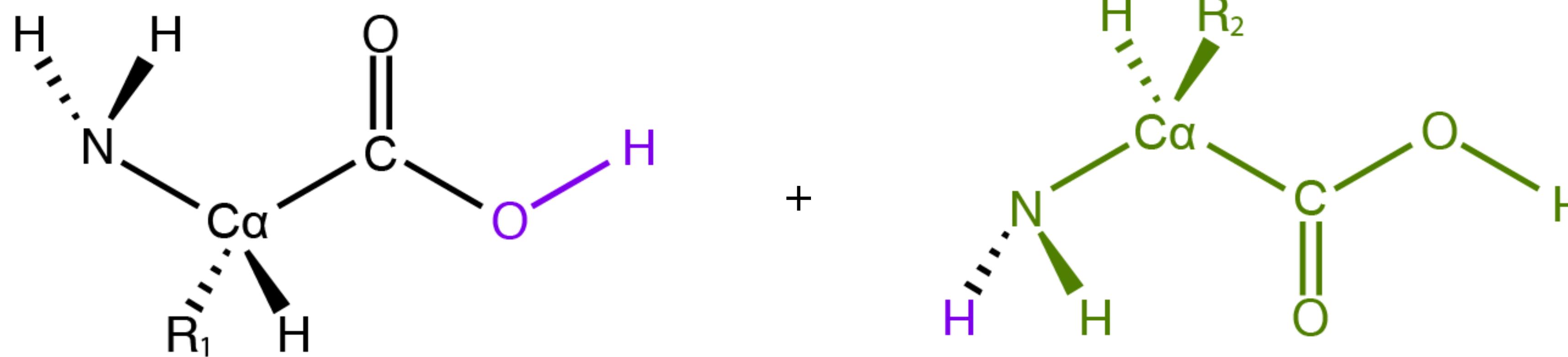
# Generating new, designable structures expands functional space



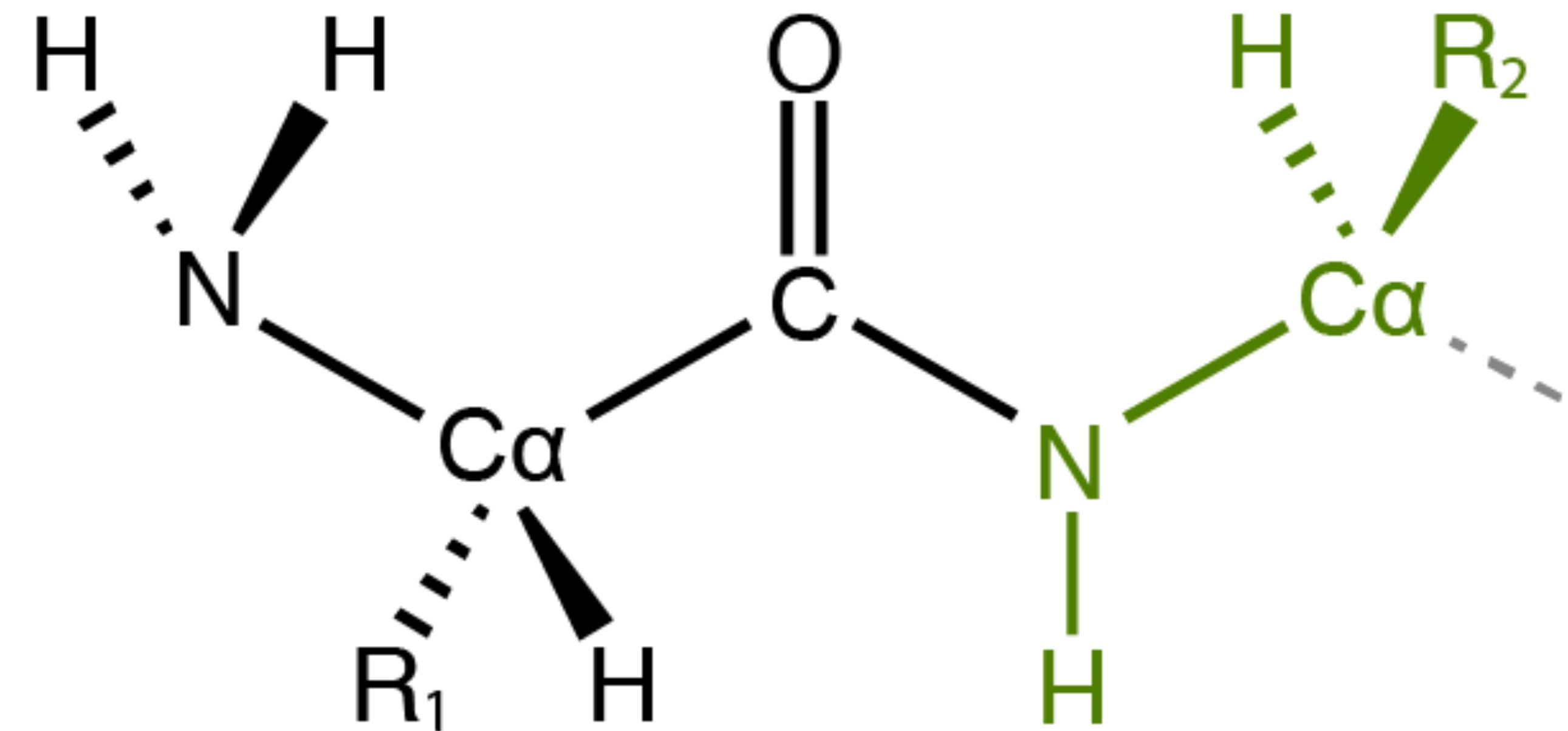
Challenge: Generate diverse and designable structures



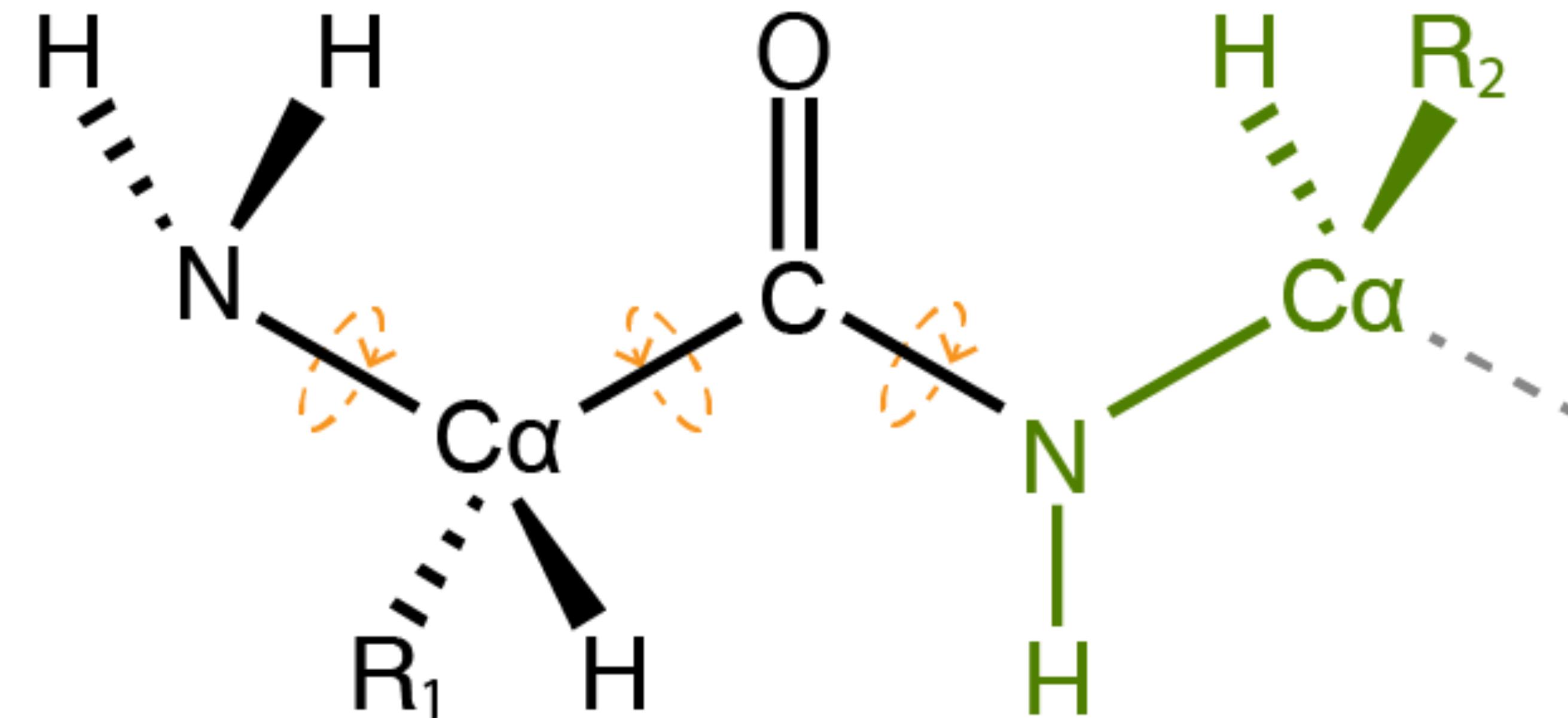
# Proteins are polypeptide chains



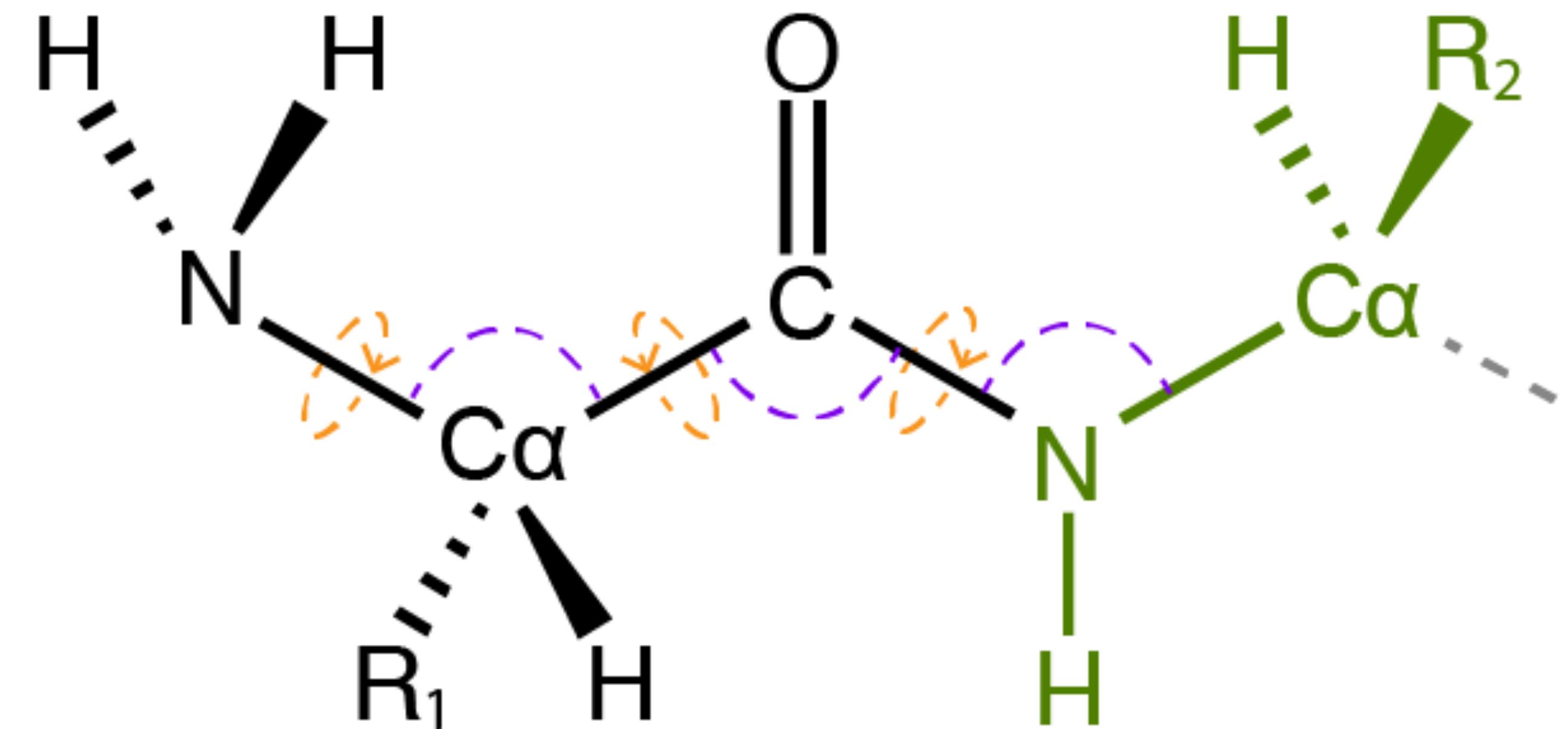
# Protein structure is determined by bond angles



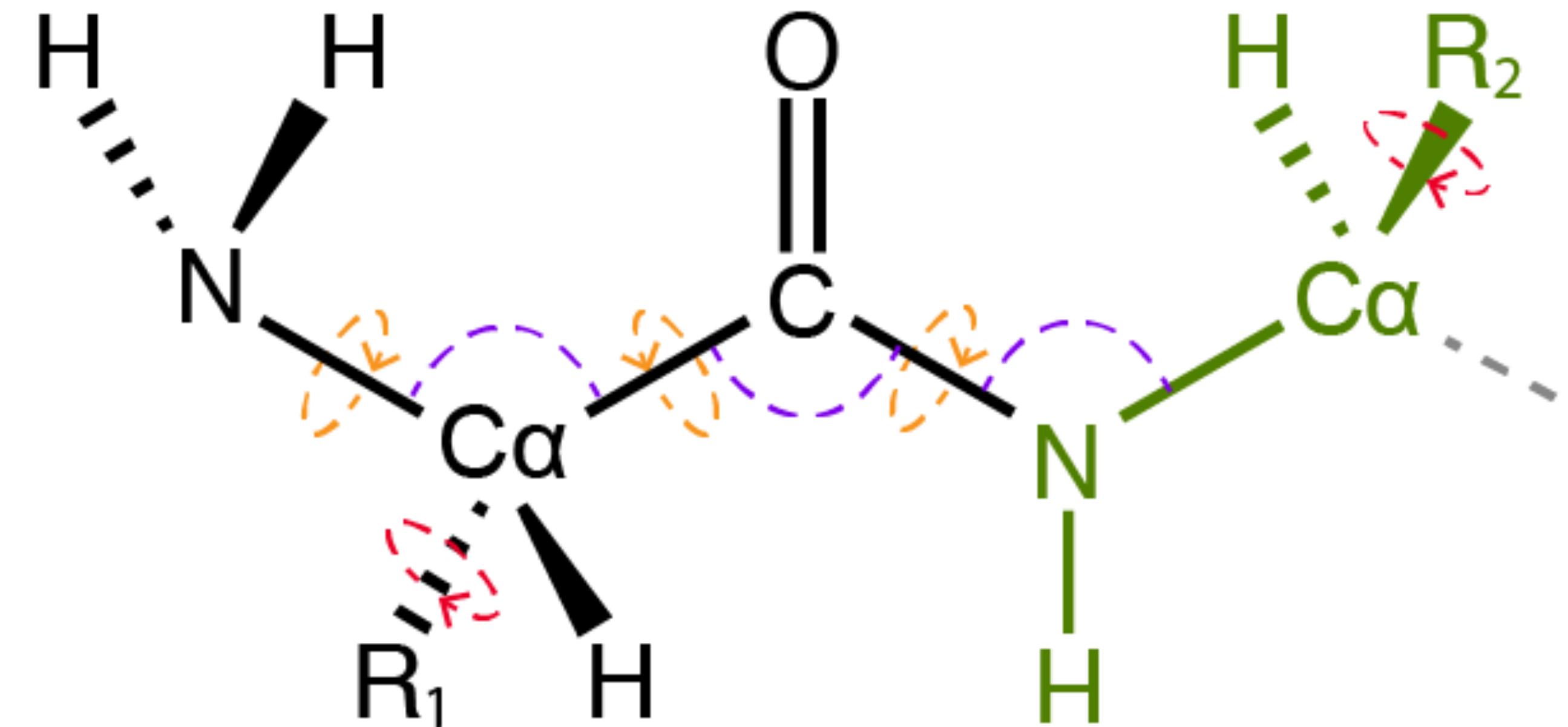
# Protein structure is determined by bond angles



# Protein structure is determined by bond angles

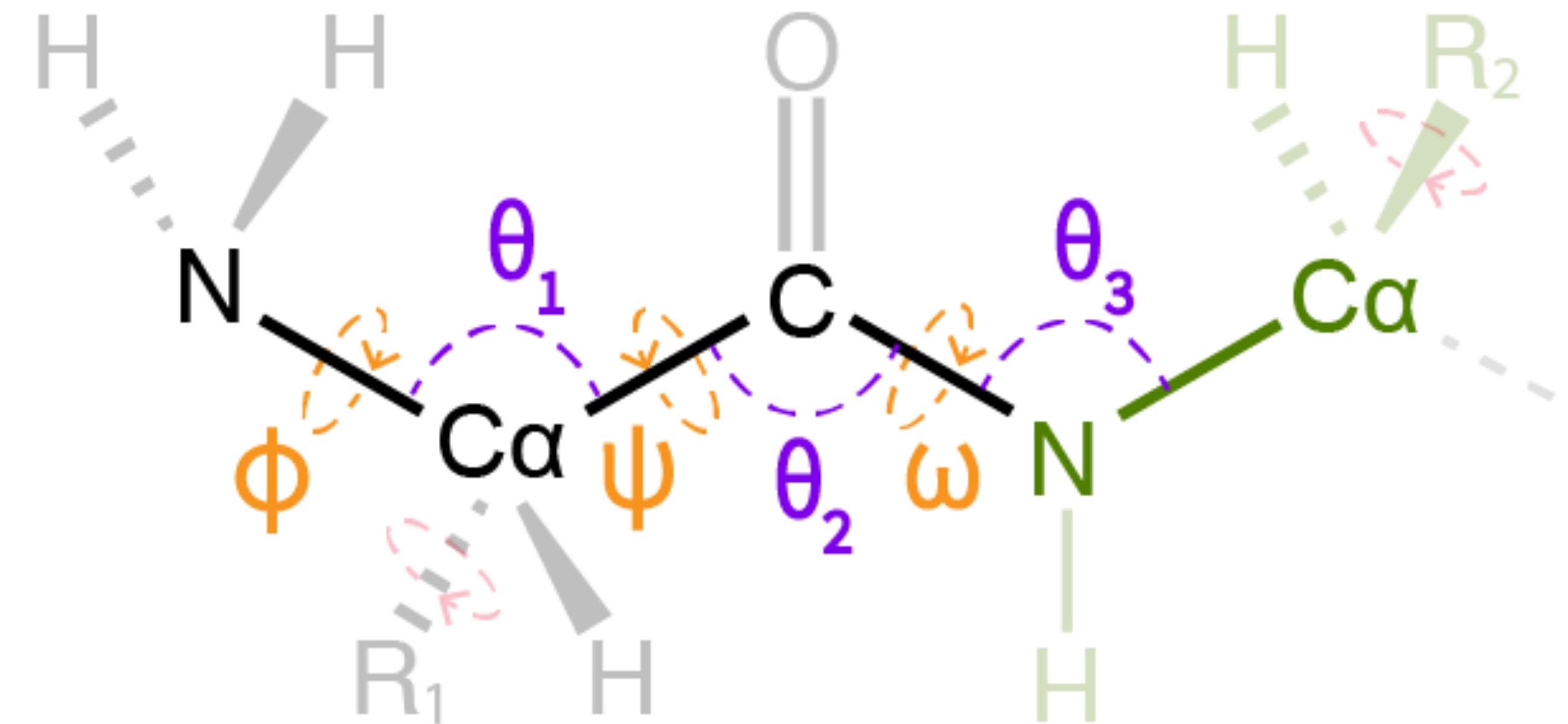


# Protein structure is determined by bond angles



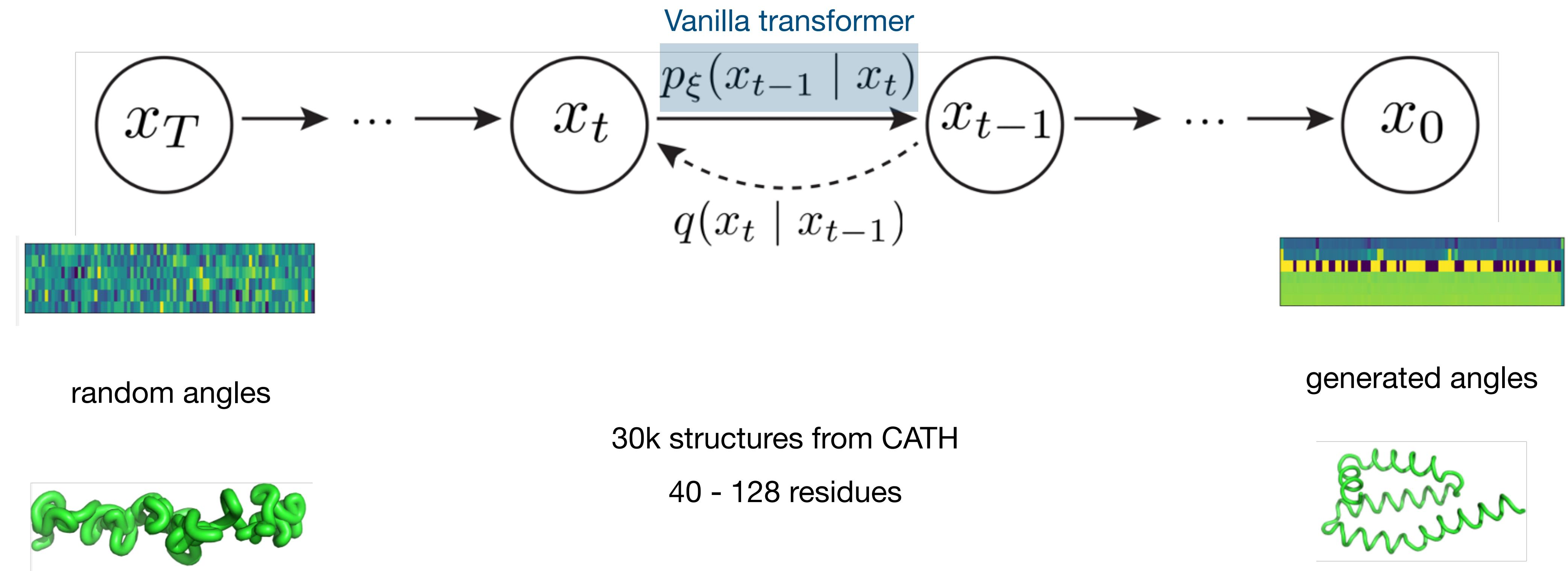
These six angles (for every consecutive pair of amino acids) fully determine the structure

# We generate backbone structures represented by **dihedral** and **bond** angles

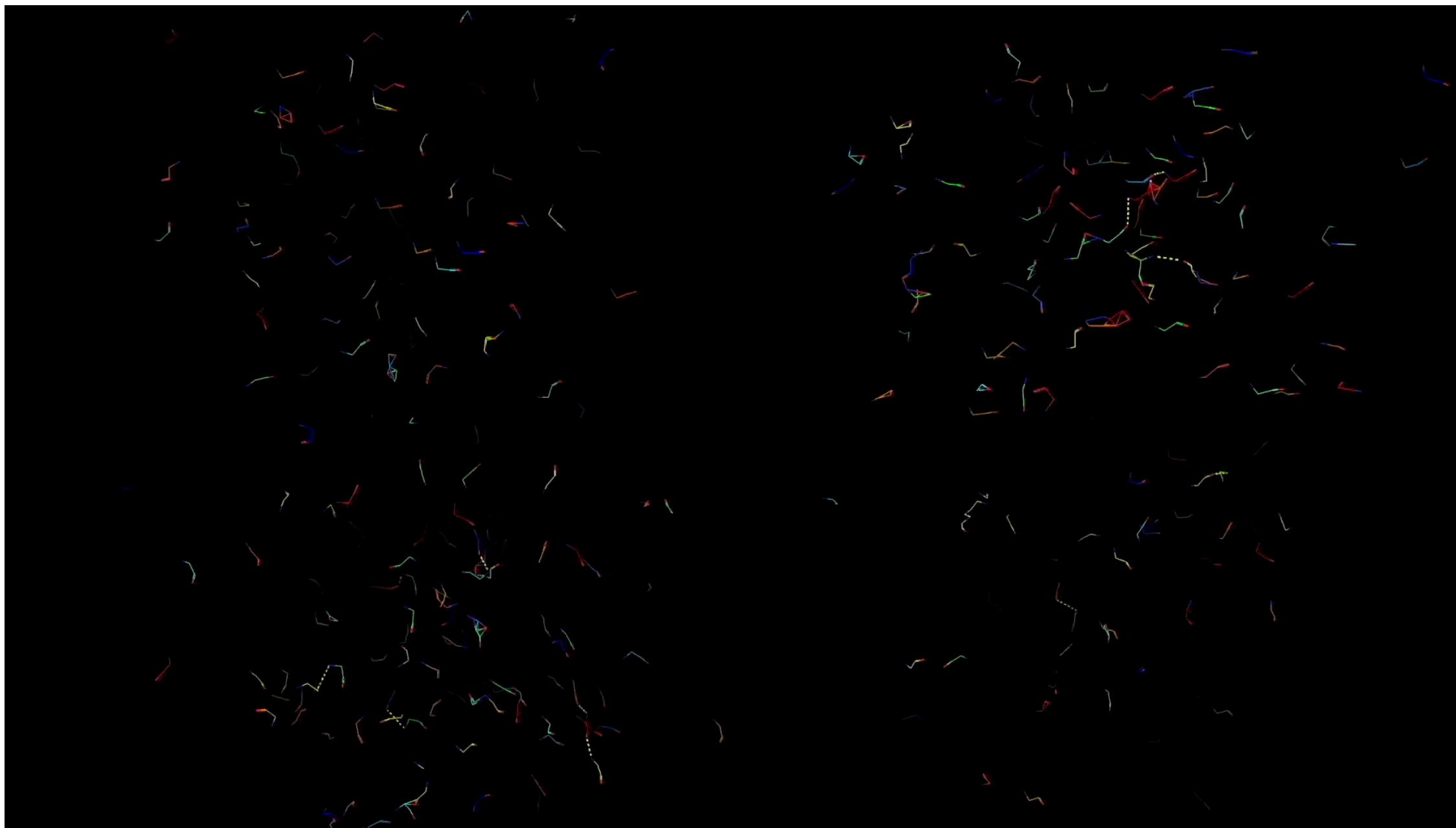


These six angles (for every consecutive pair of amino acids) fully determine the backbone structure

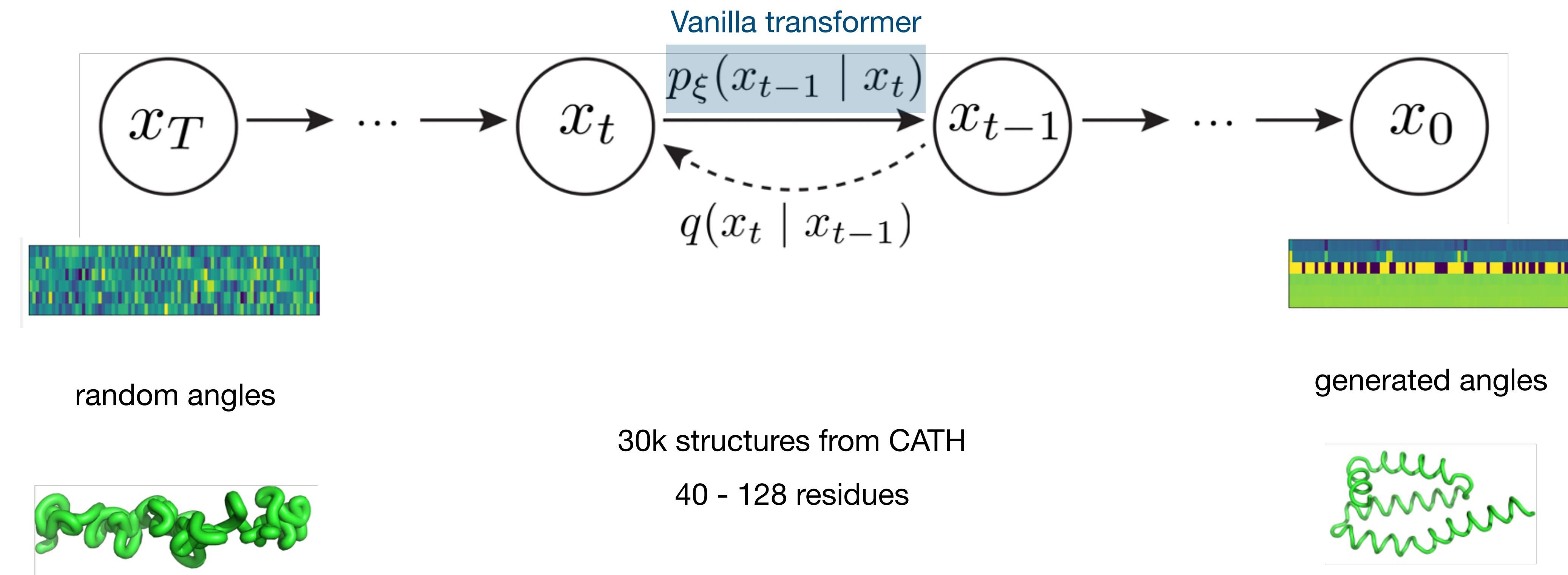
# FoldingDiff uses diffusion to generate angles



# Diffusion on 3D coordinates requires equivariances



# Evaluate generations at 3 levels



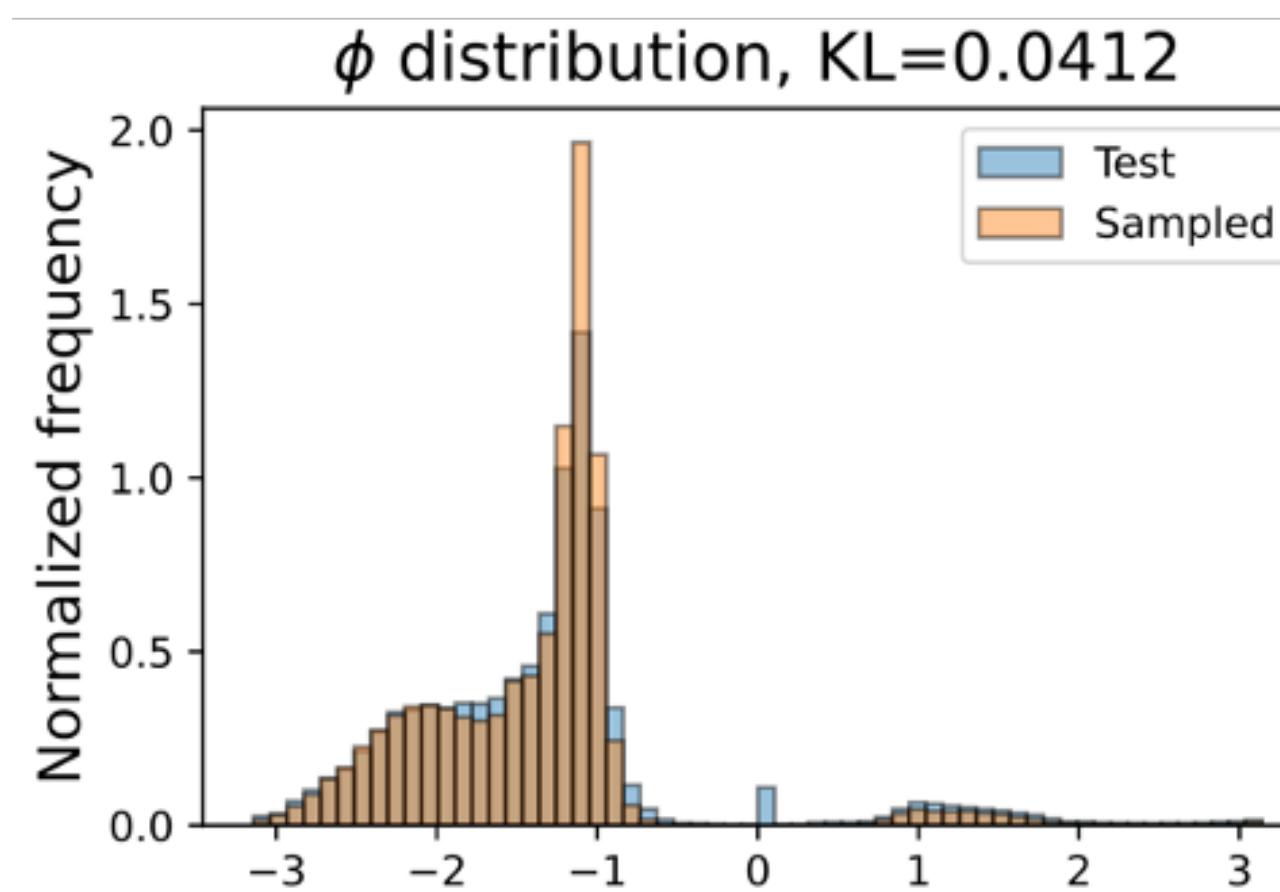
generated angles

structural motifs

overall structures

# Generated angles match test distribution

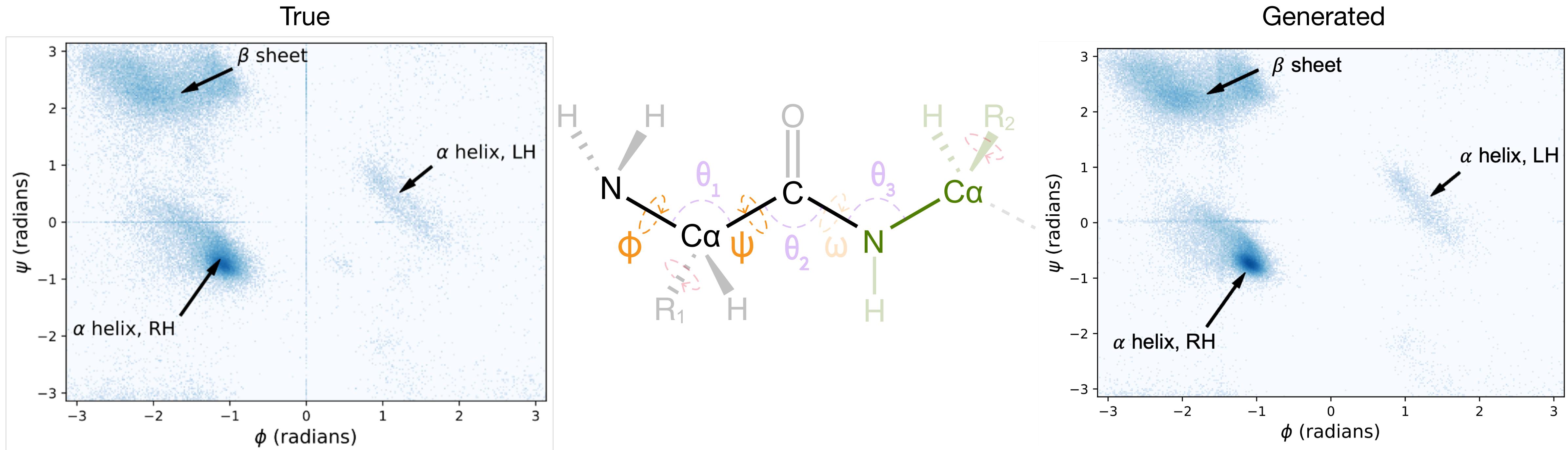
noise  $\longrightarrow$  sample  $\longrightarrow$  compare to true distribution



Generated distributions match natural distribution of individual angles

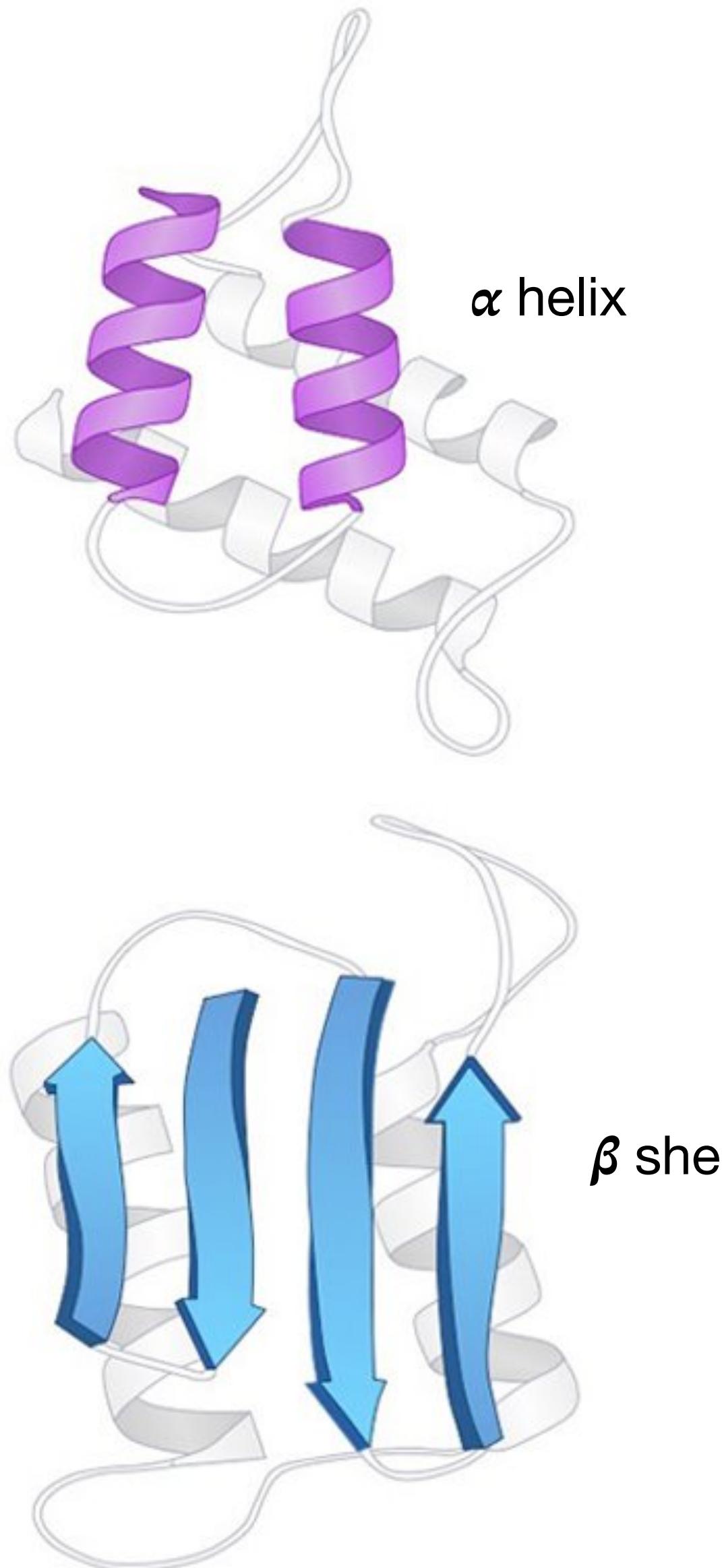
# FoldingDiff captures correlations between angles

noise  $\longrightarrow$  sample  $\longrightarrow$  compare  $(\phi, \psi)$  co-occurrence

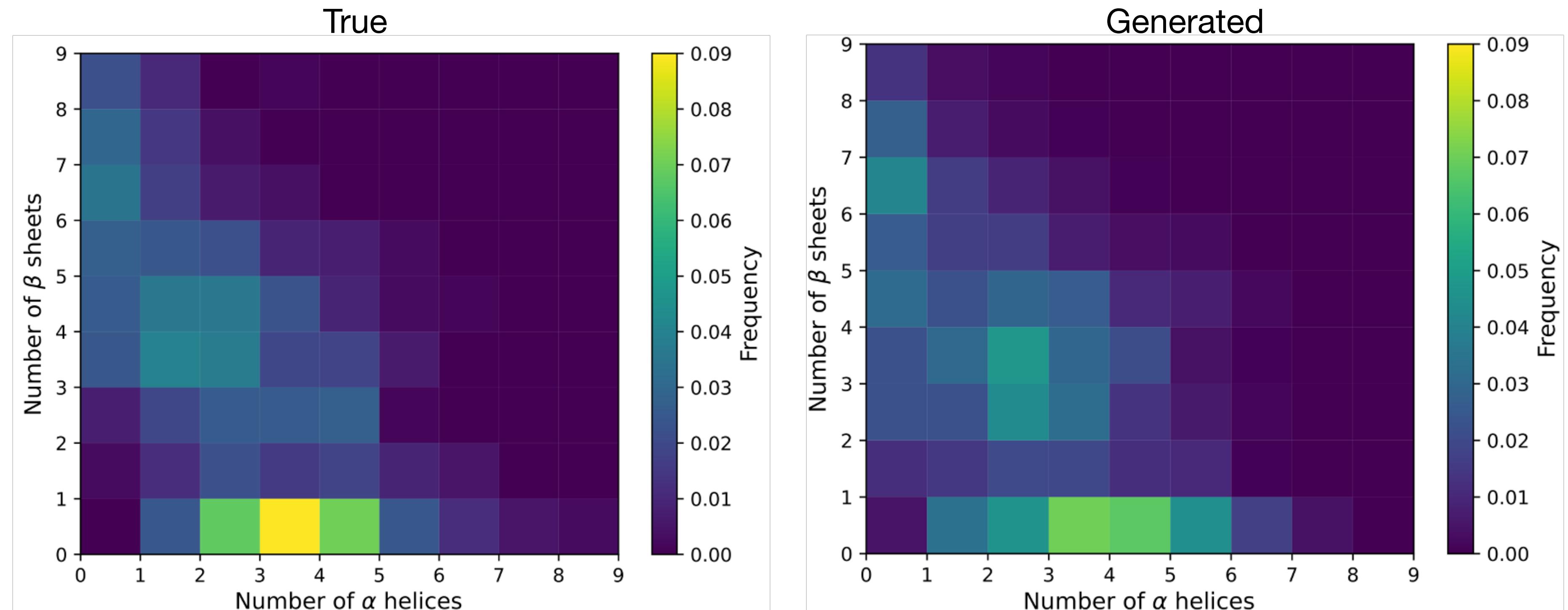


FoldingDiff generates correlations that define common structural motifs

# Generated secondary structures match test structures

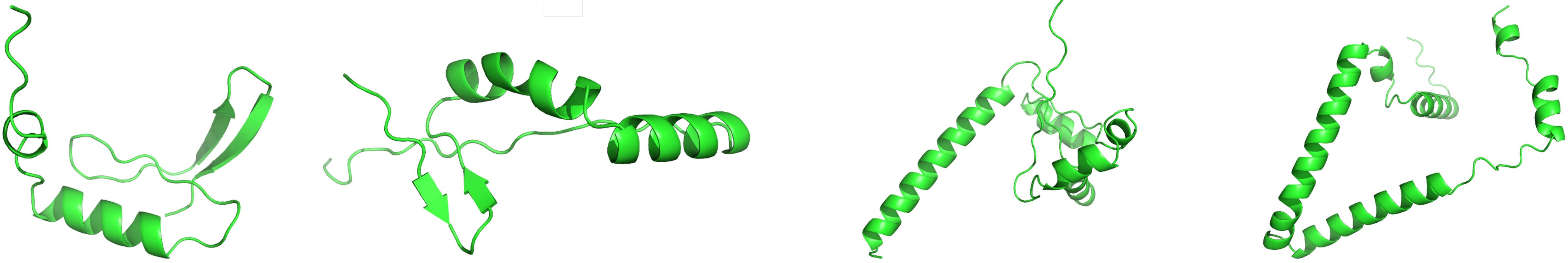


noise  $\longrightarrow$  sample  $\longrightarrow$  measure helix, sheet structures

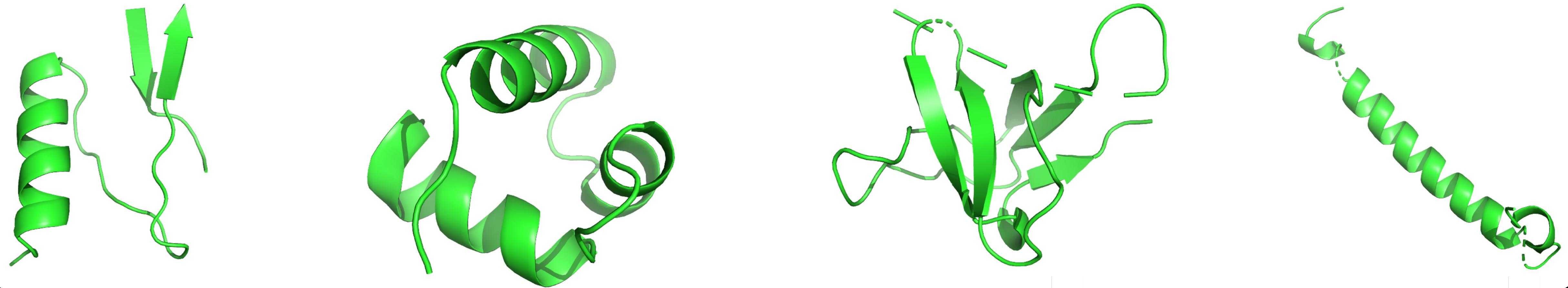


# Generated structures look reasonable

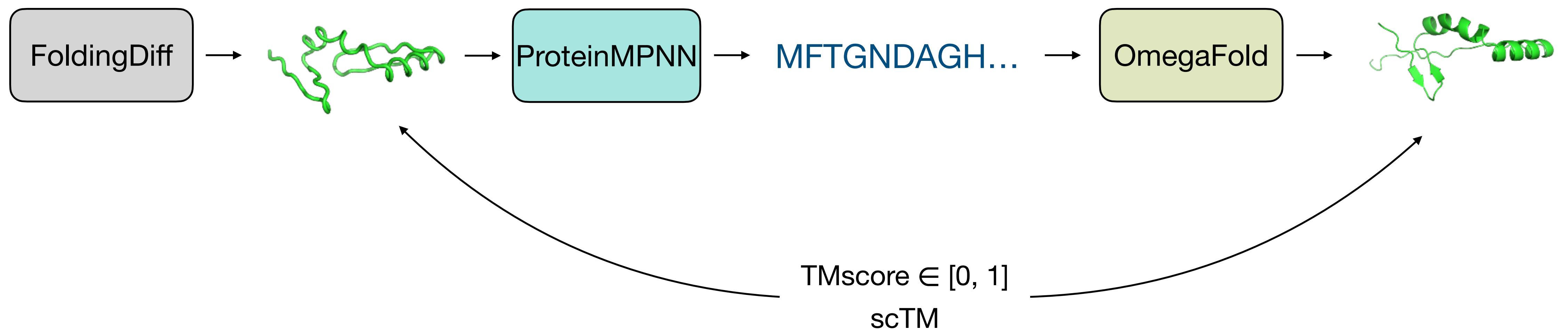
Generated structures



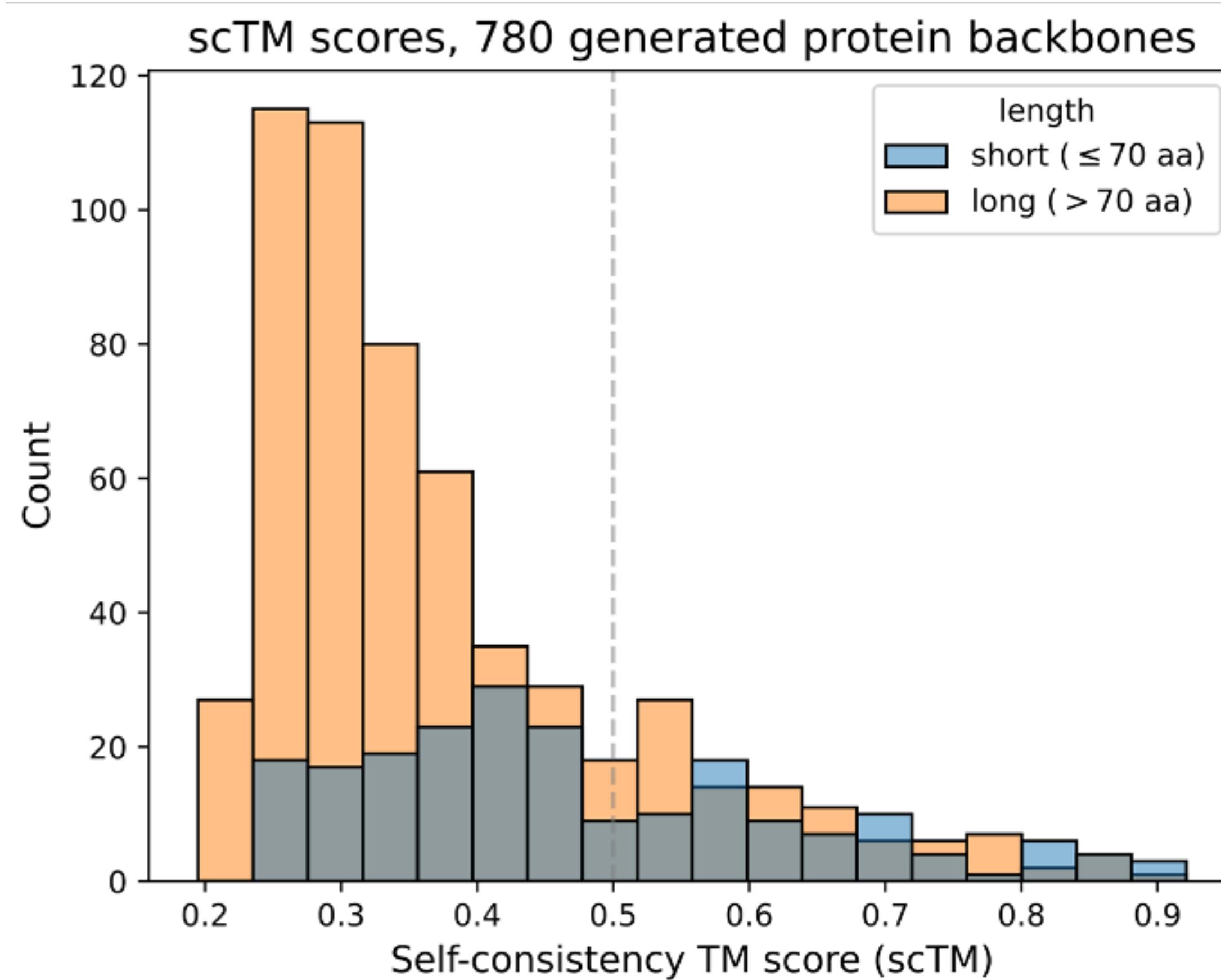
Training set structures



# Measure designability of structures with self-consistency TMscore



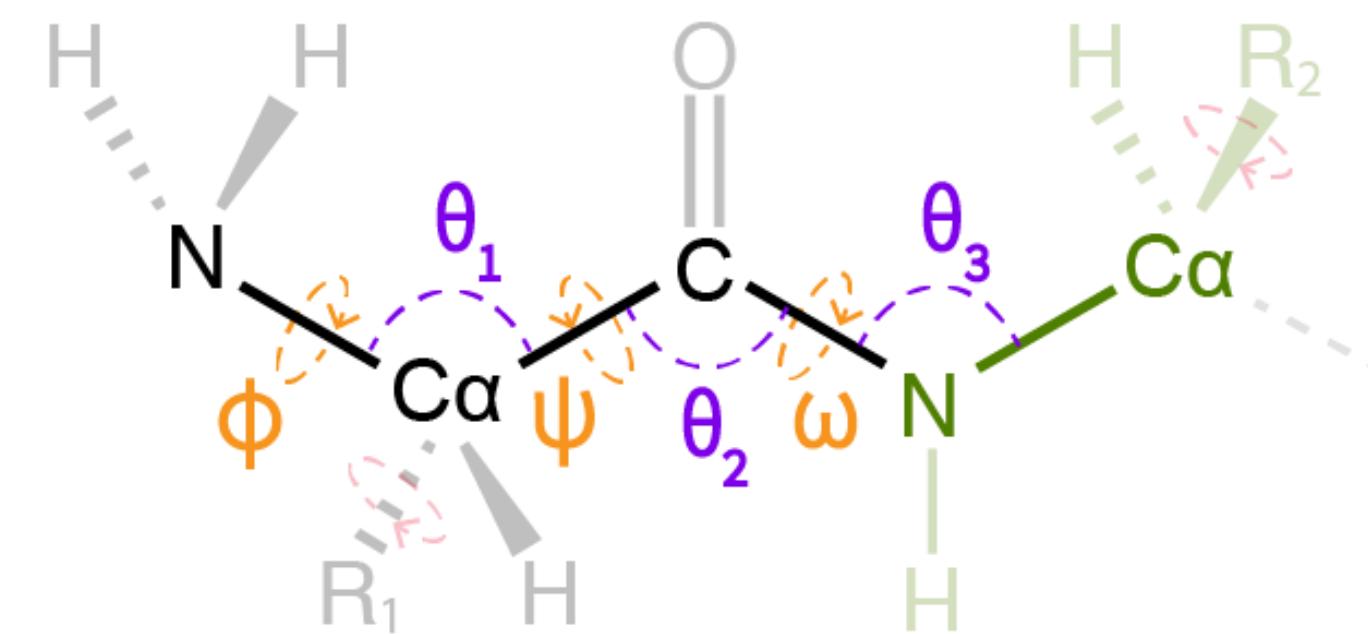
# Many generated structures are designable



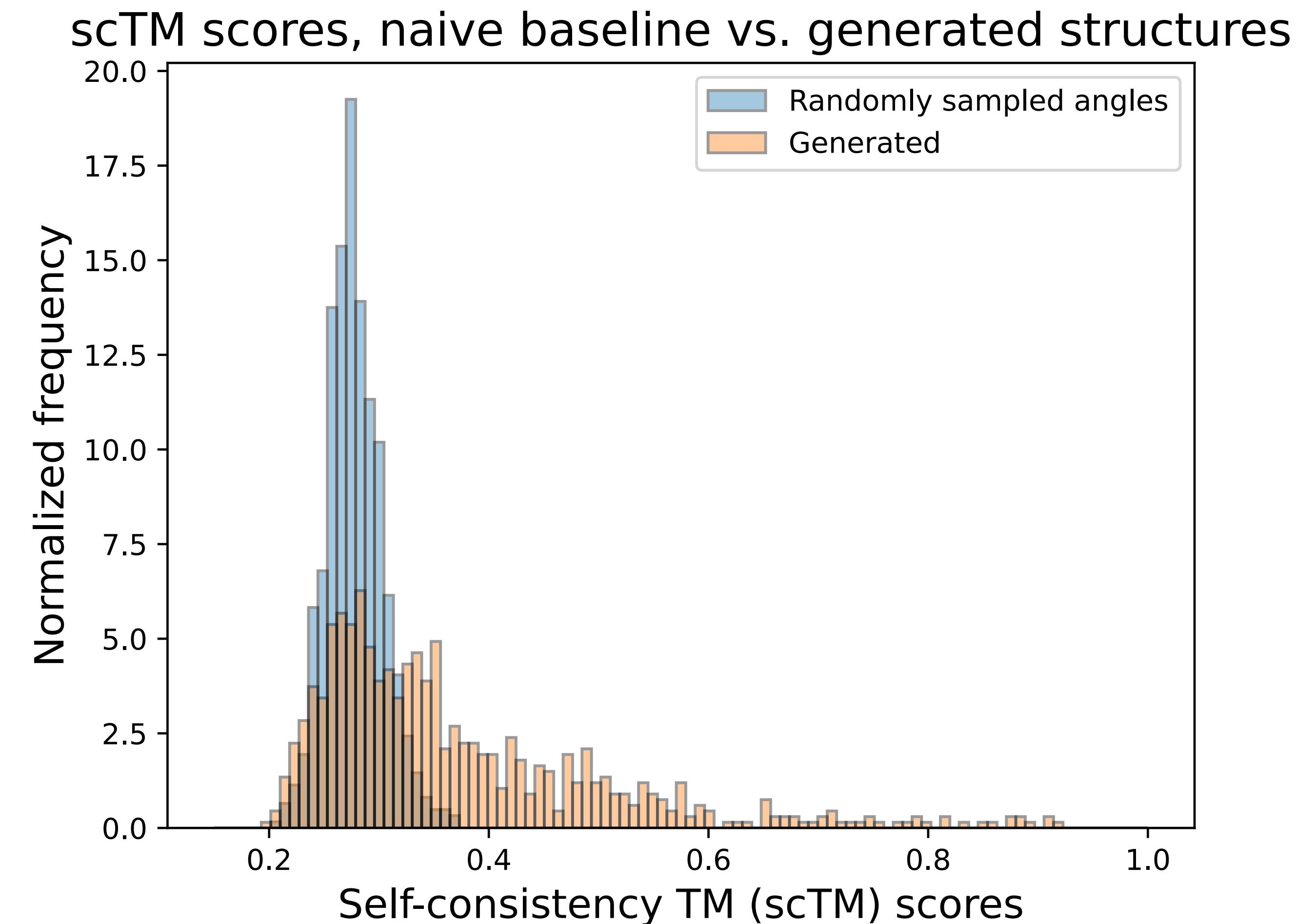
scTM > 0.5	$\leq 70$ aa	$> 70$ aa
FoldingDiff	76/210	87/570
ProtDiff (Trippe et al.)	36/210	56/570

Significant improvements over point cloud diffusion model

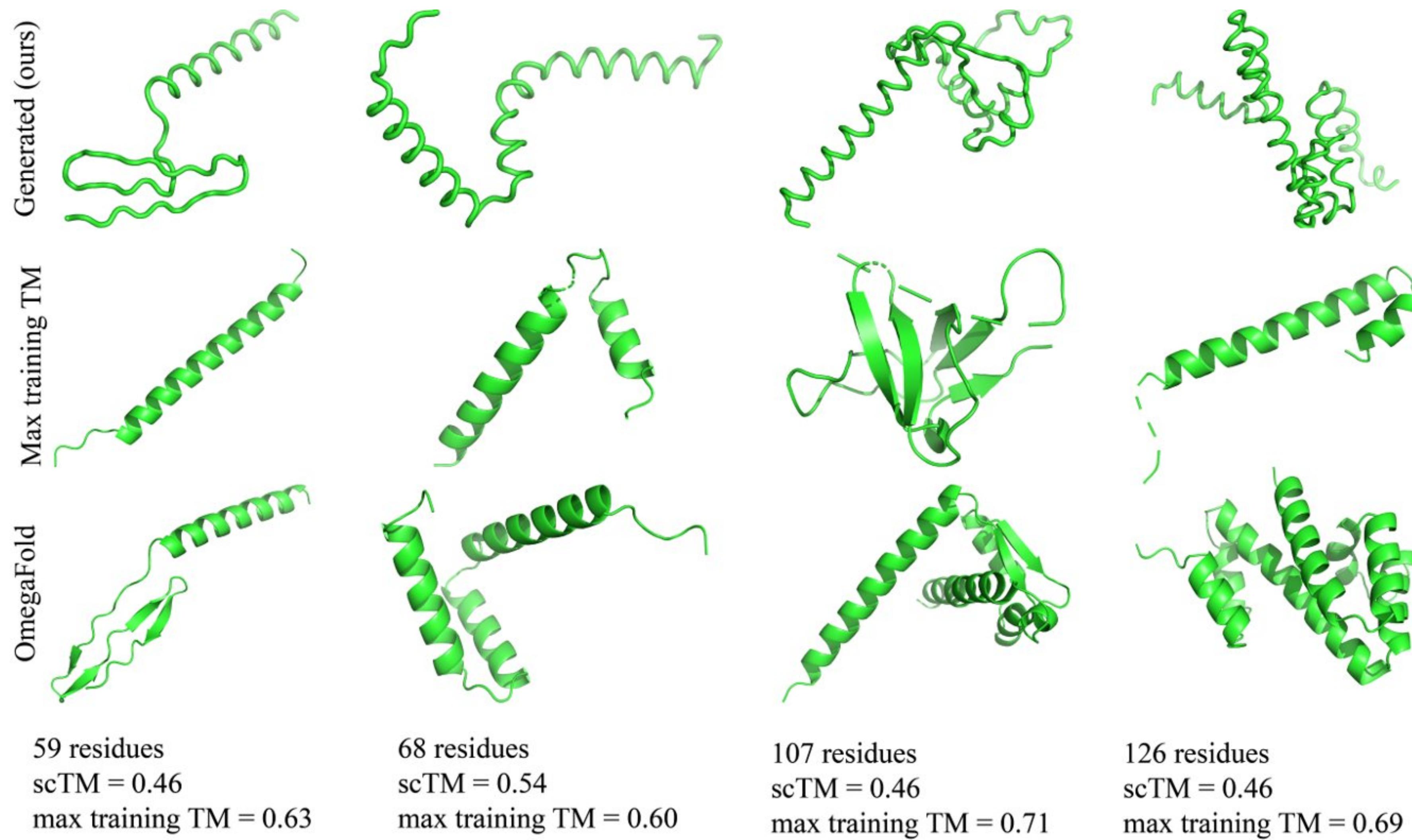
# FoldingDiff structures are better than random baseline



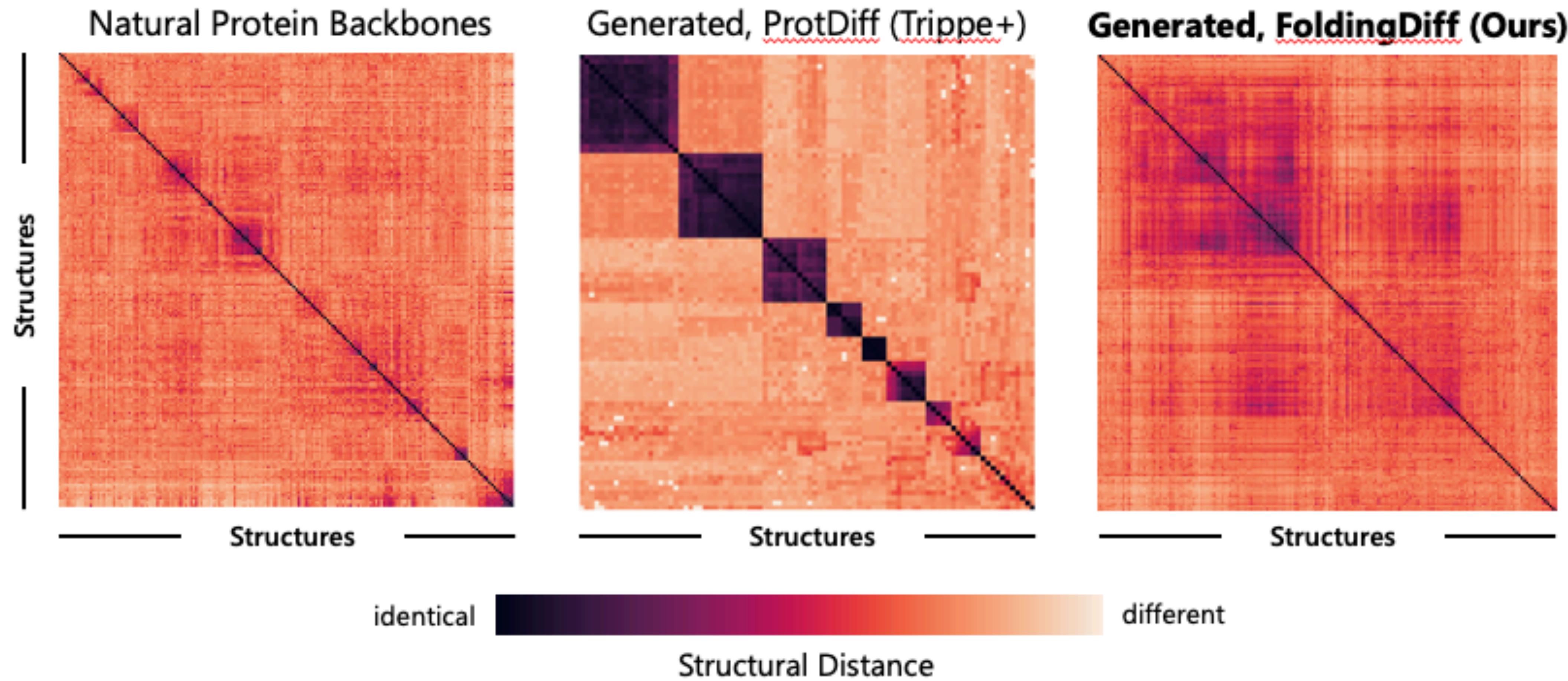
Sample sets of angles  
Preserves Ramachandran plot



# Generated structures are diverse



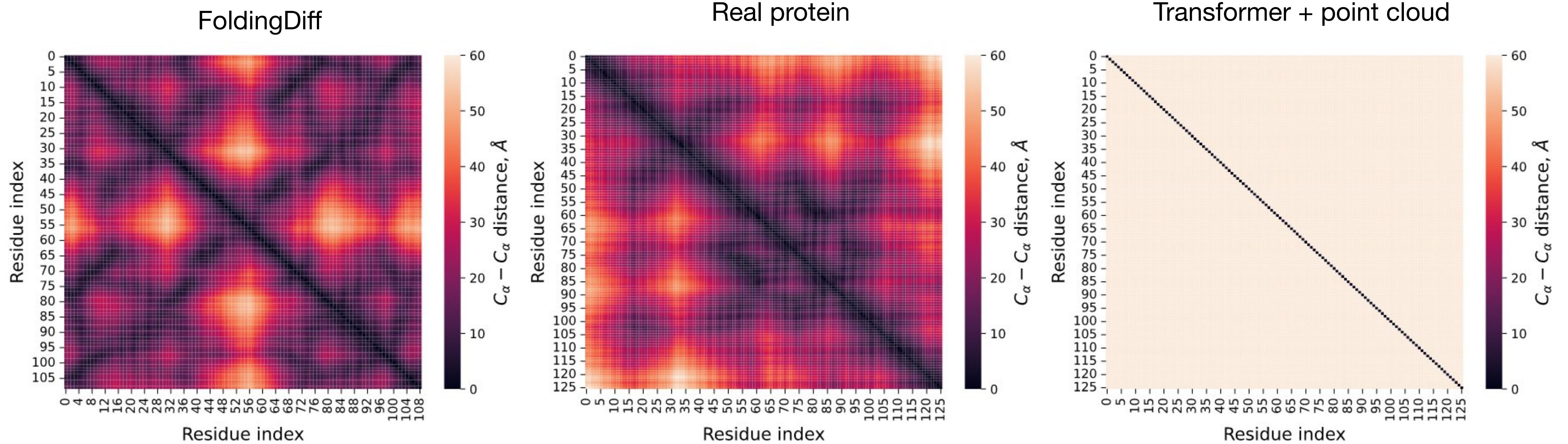
# Generated structures are diverse



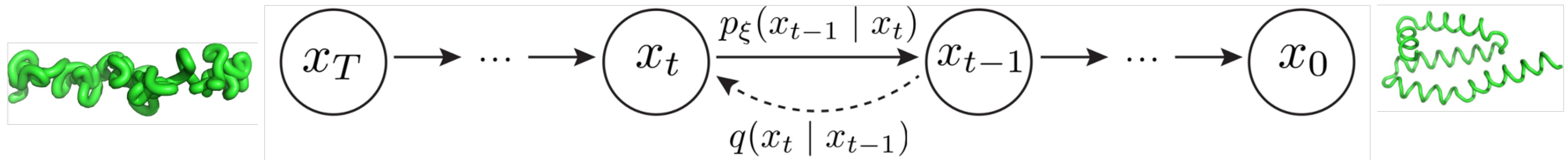
# Generated structures are diverse



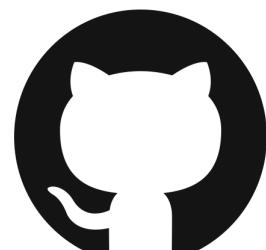
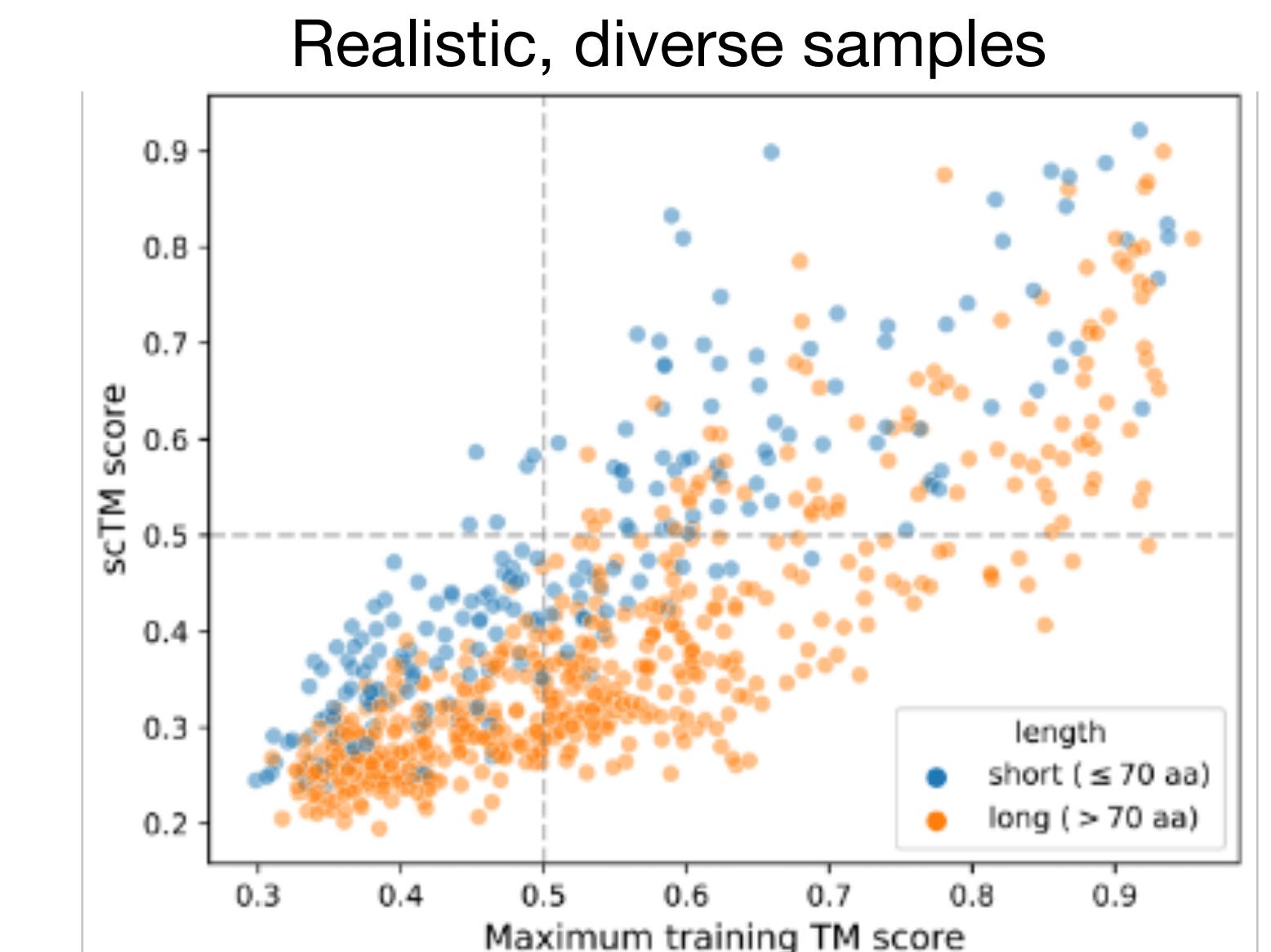
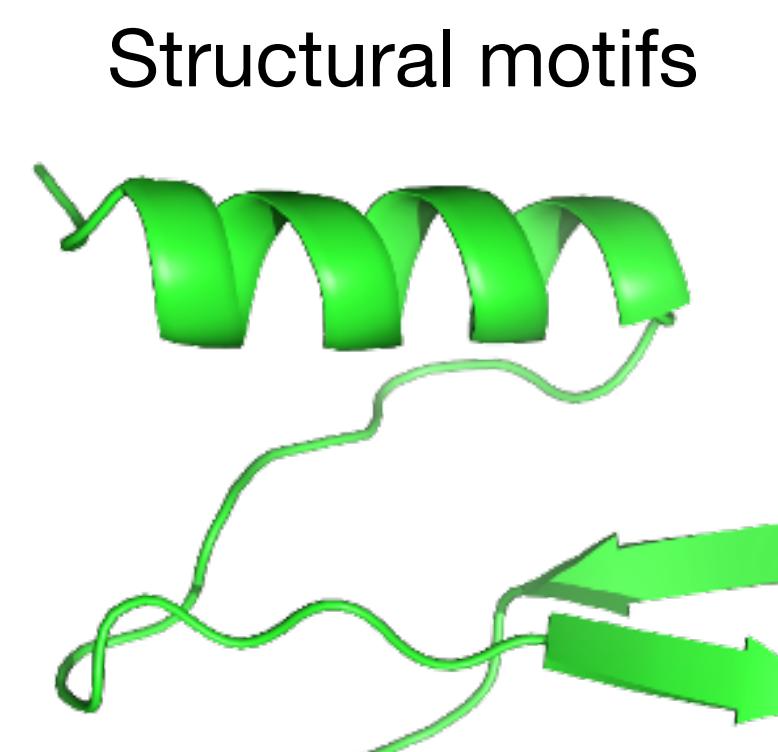
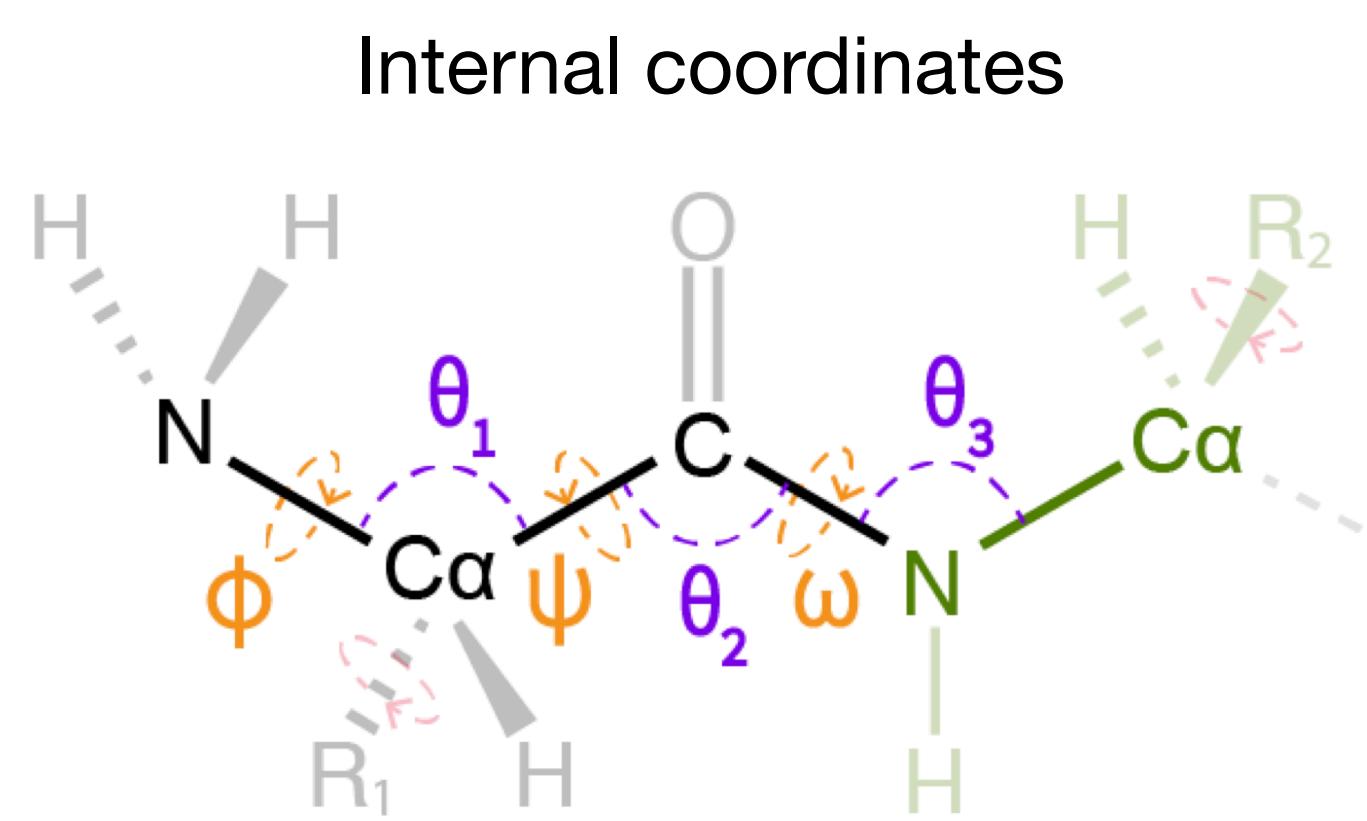
# Angle formulation enables generation from vanilla transformers



# FoldingDiff is first step towards generating new functions



Generate protein backbones by mirroring the folding process



[github.com/microsoft/foldingdiff](https://github.com/microsoft/foldingdiff)

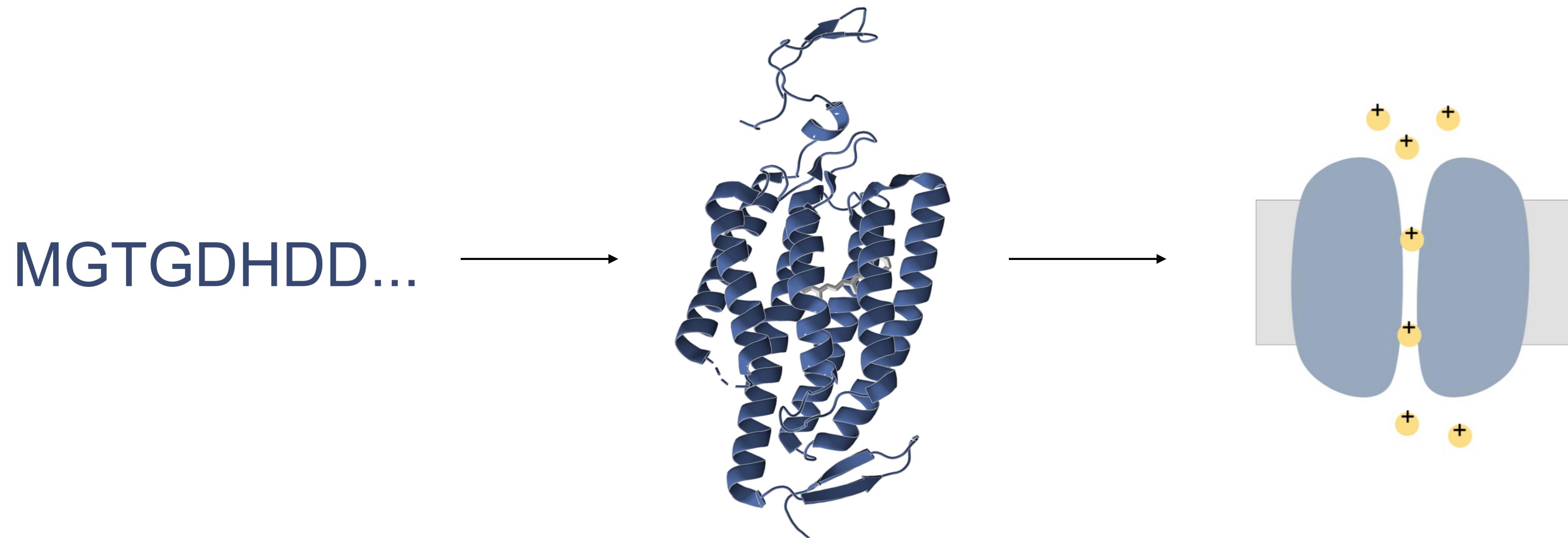


[huggingface.co/spaces/wukevin/foldingdiff](https://huggingface.co/spaces/wukevin/foldingdiff)



BioNeMo?

# Generating new, valid sequences expands functional space



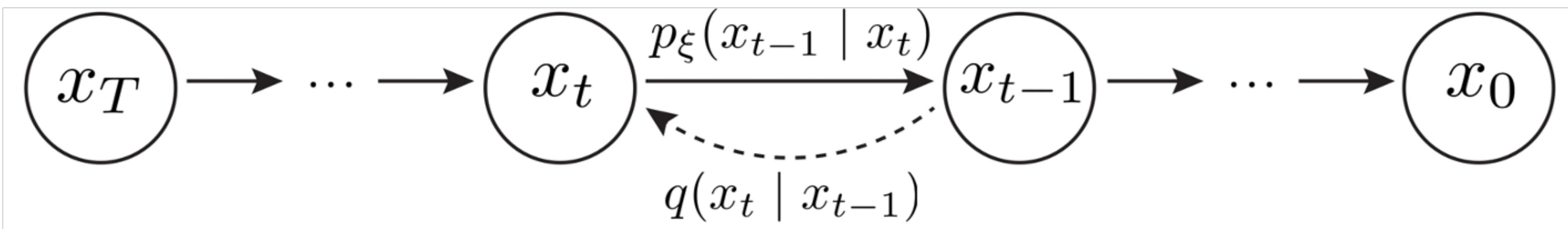
# Generating new, valid sequences expands functional space



Challenge: Generate valid and diverse sequences

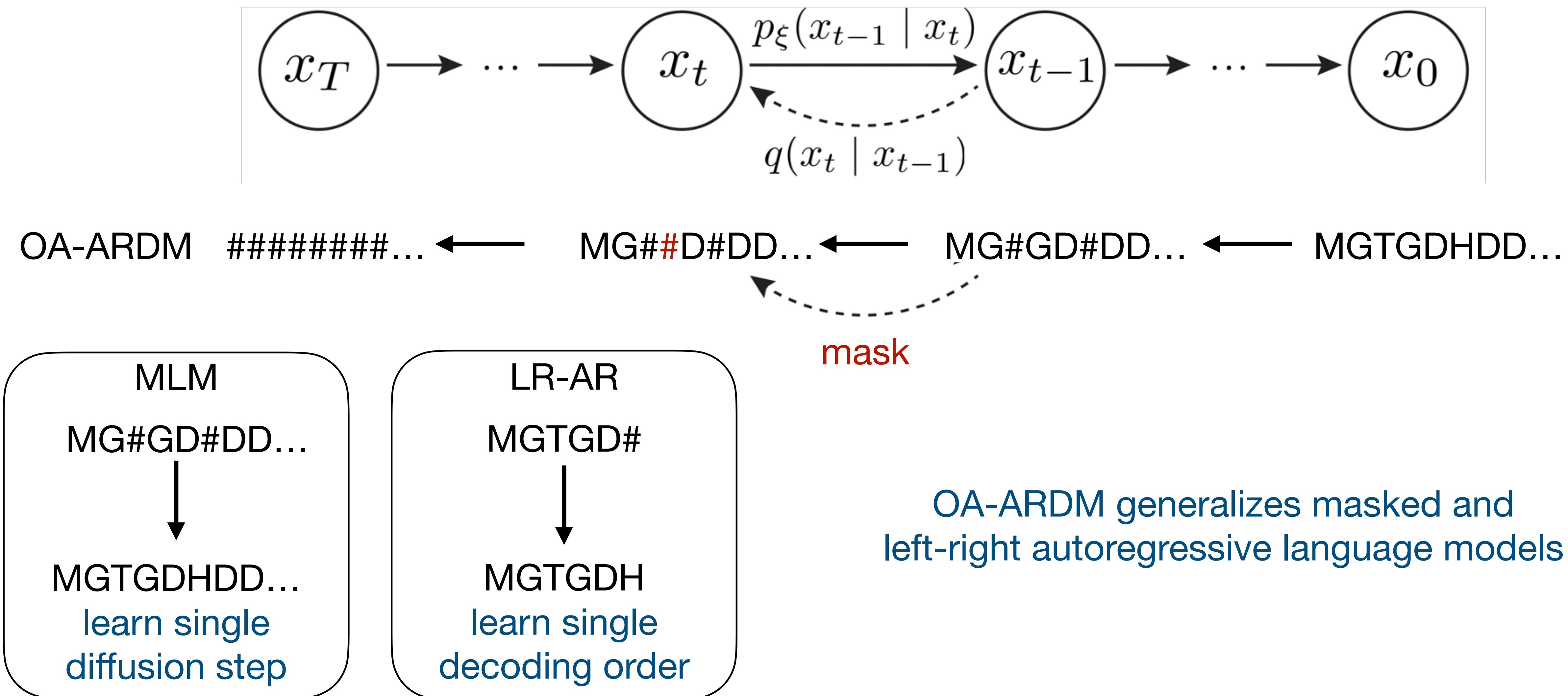


# Use discrete diffusion to generate sequences

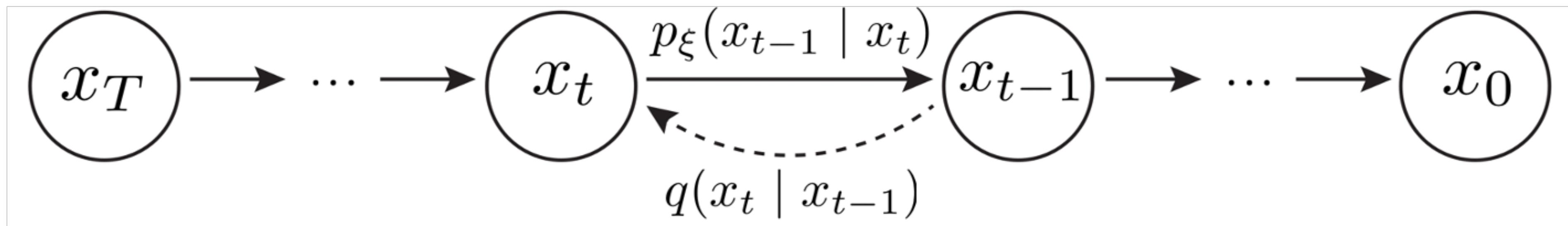


MGTGDHDD...

# Order-agnostic autoregressive diffusion links diffusion, masked language models, and autoregressive language models



# Discrete diffusion mutates towards a random string

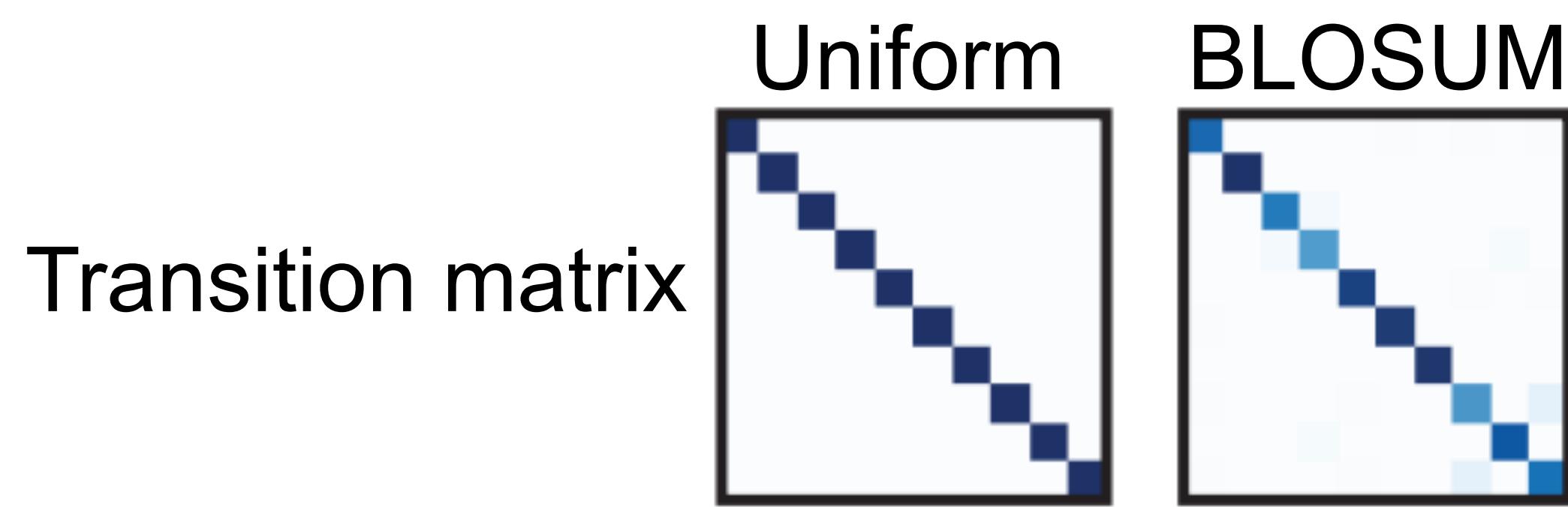


OA-ARDM    #####... ← MG##D#DD... ← MG#GD#DD... ← MGTGDHDD...

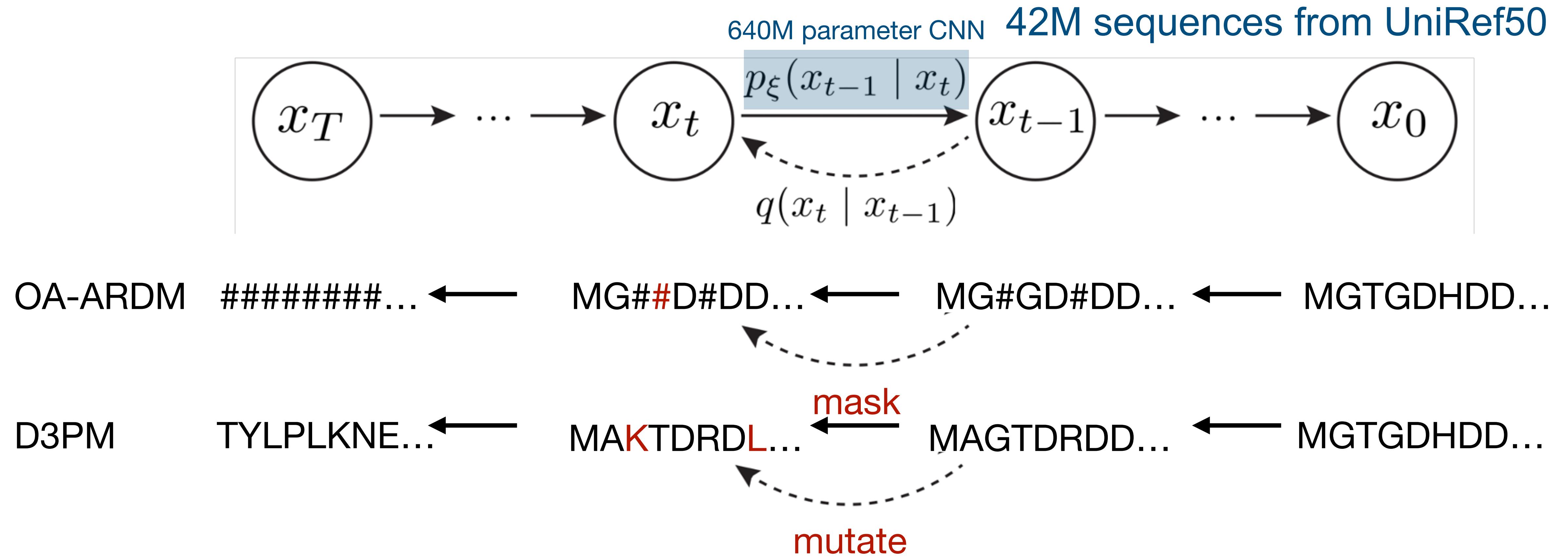
D3PM    TYLPLKNE... ← MAK**T**DRDL... ← MAGTDRDD... ← MGTGDHDD...

mask

mutate



# Train with same dataset and architecture for fair comparisons



# Scale and diffusion improve generated amino acid distribution

Model

MFTGNDAGH...  
MHGAPOKLO...  
MIEASWQNI...  
MTYVVNMAD...



Natural amino acid usage

Model	Params	KL
<b>Valid</b>	-	0.00059

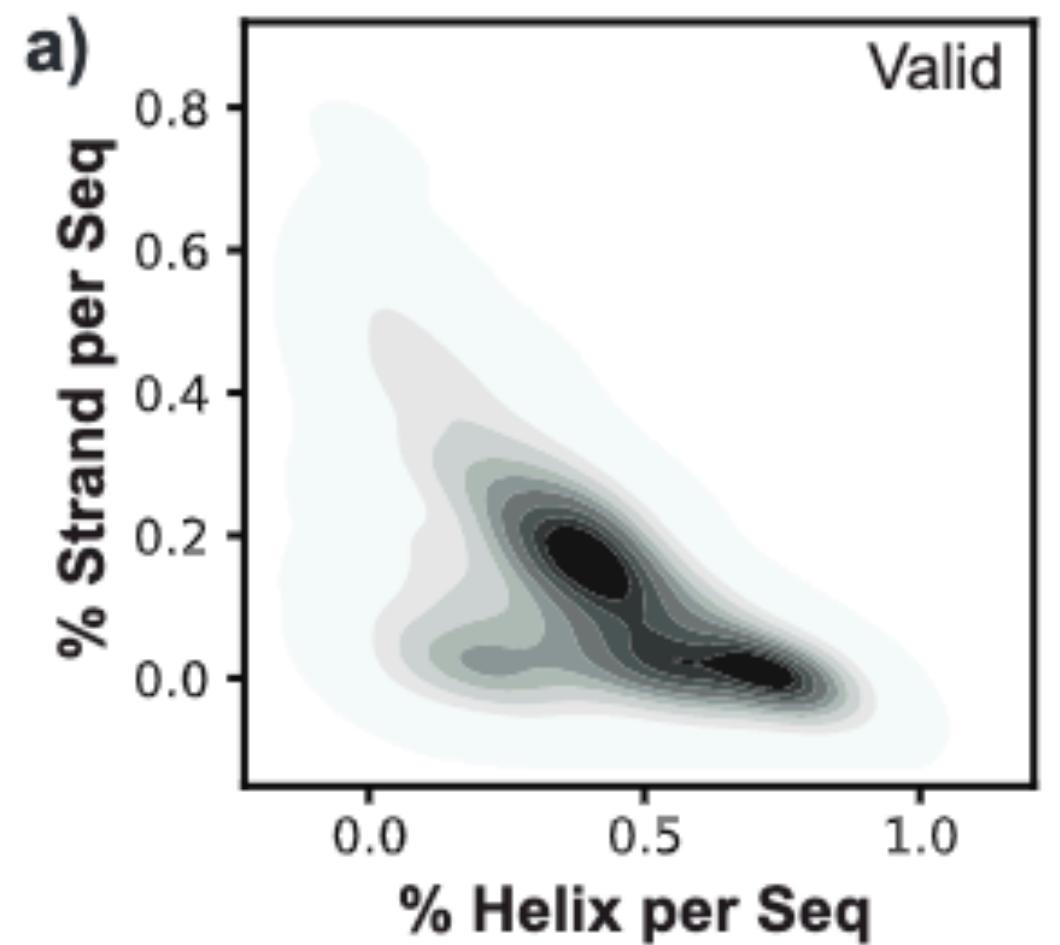
Unlike NLP, OA-ARDM is better than left-right

As expected, masked language models are not good at unconditional generation

<b>Random</b>	-	0.165
---------------	---	-------

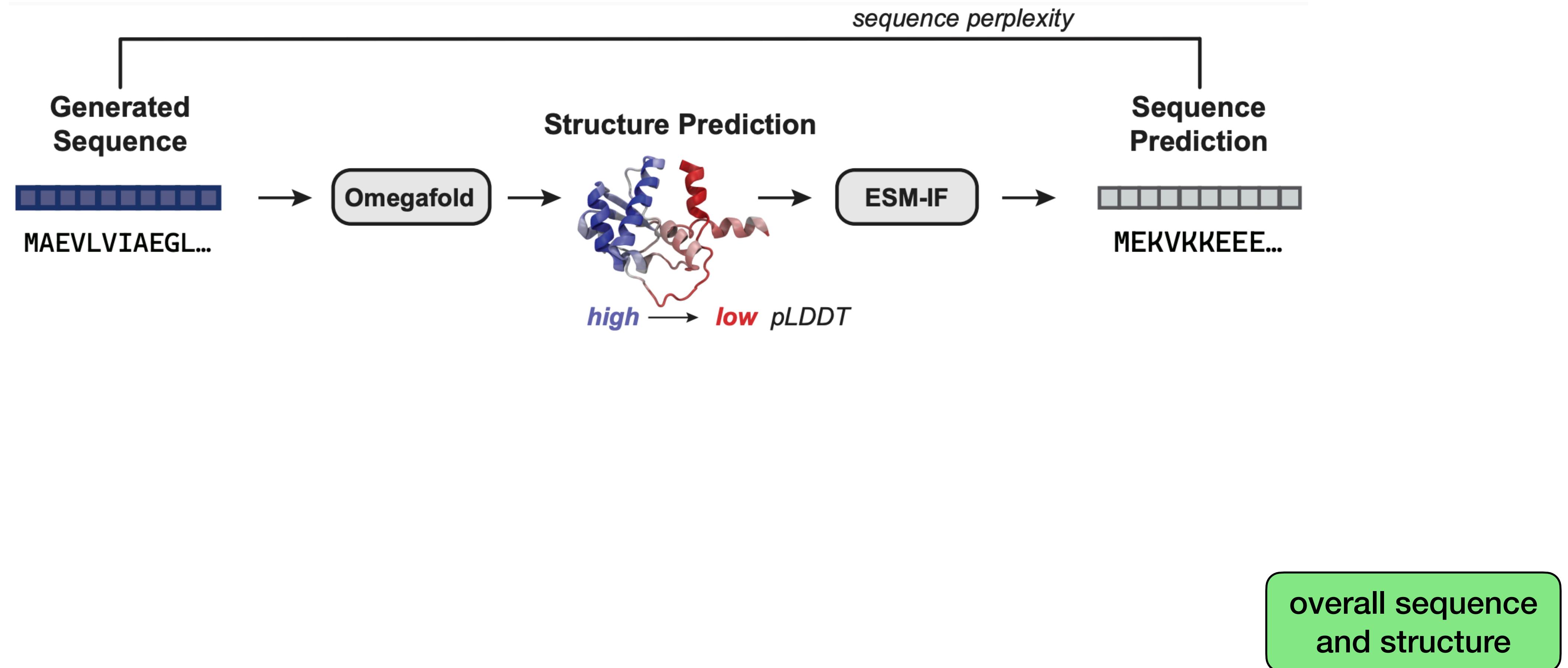
generated residues

# Diffusion models recapitulate secondary structure distributions

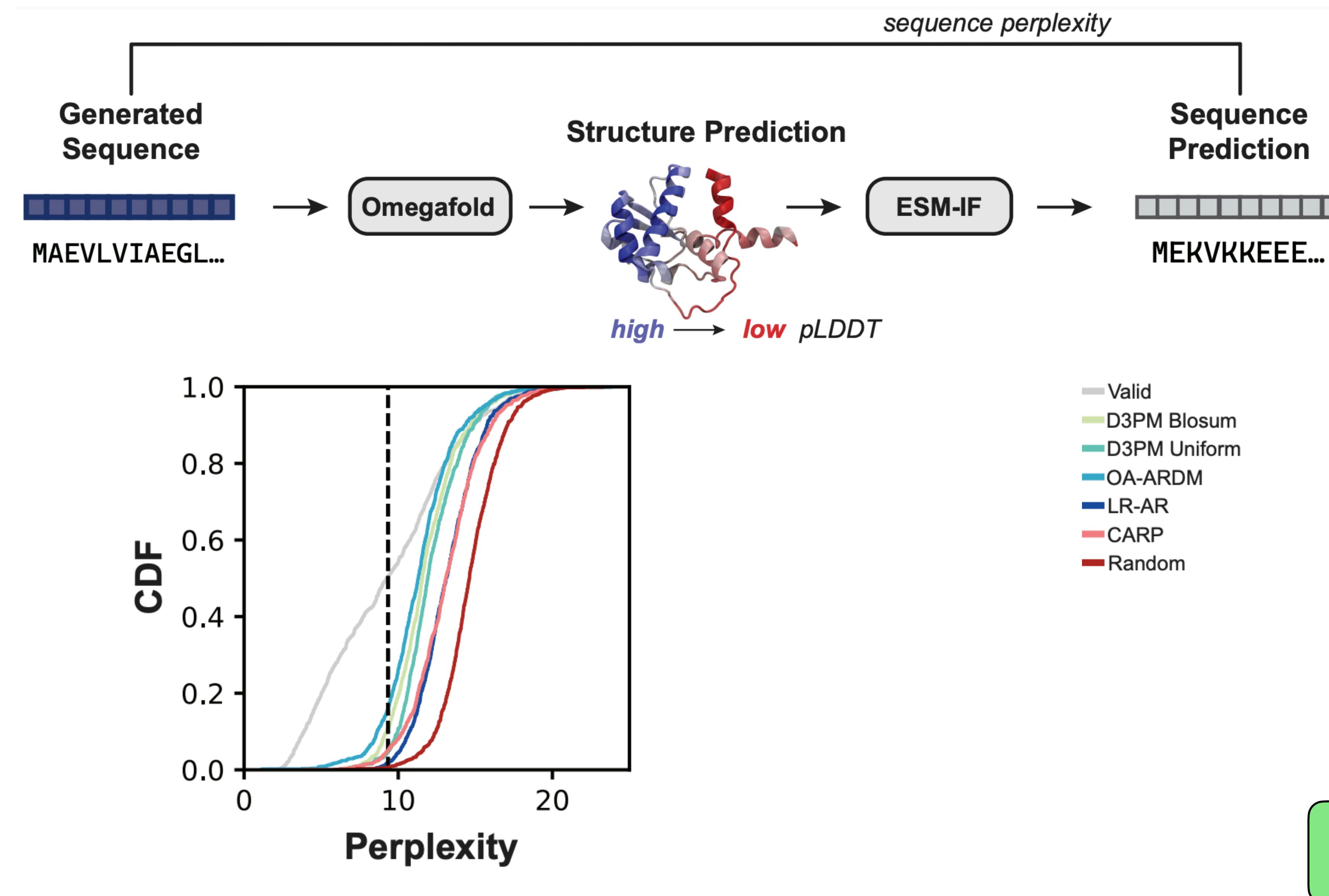


structural motifs

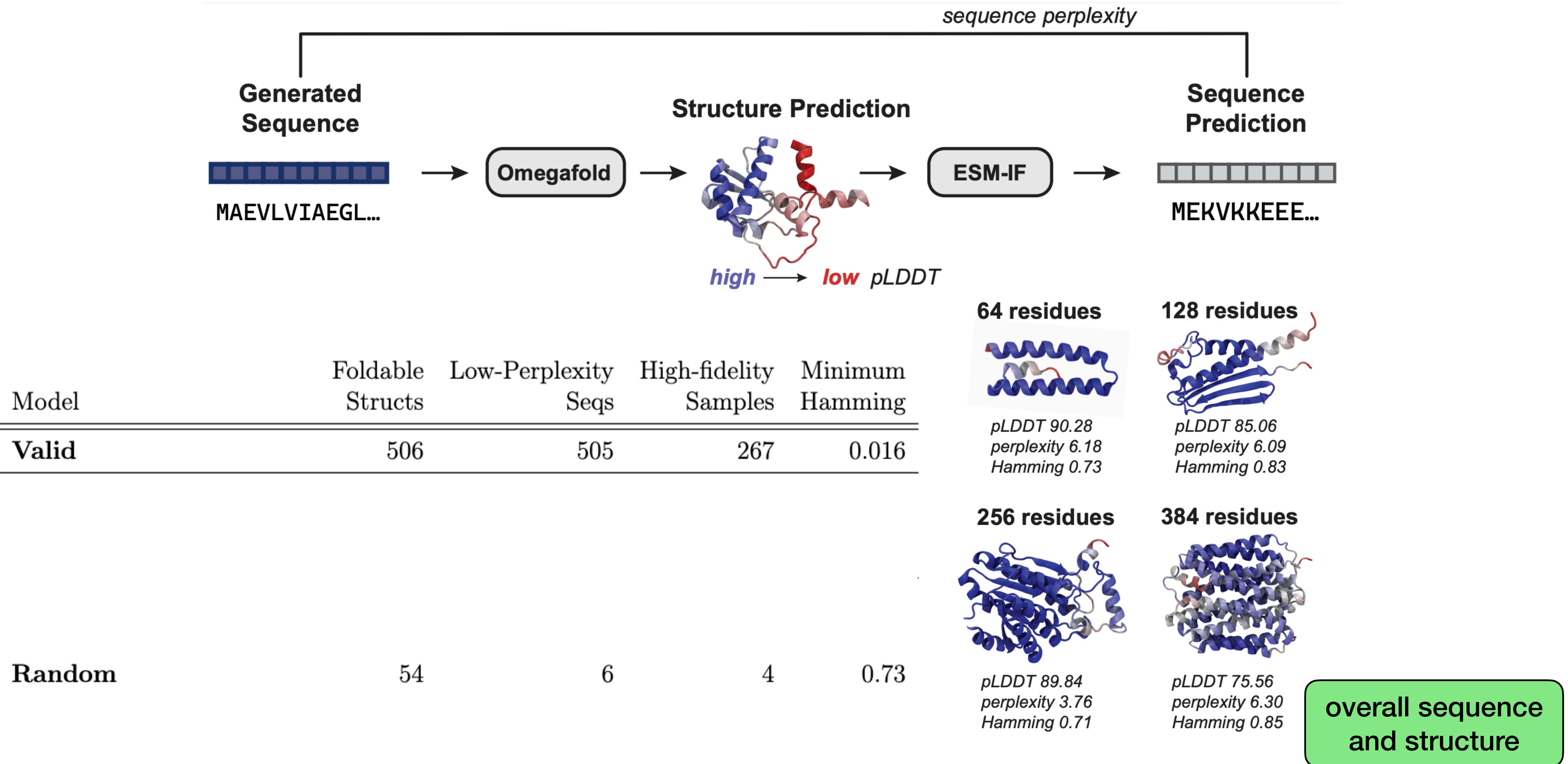
# Evaluate structural confidence and self-consistency



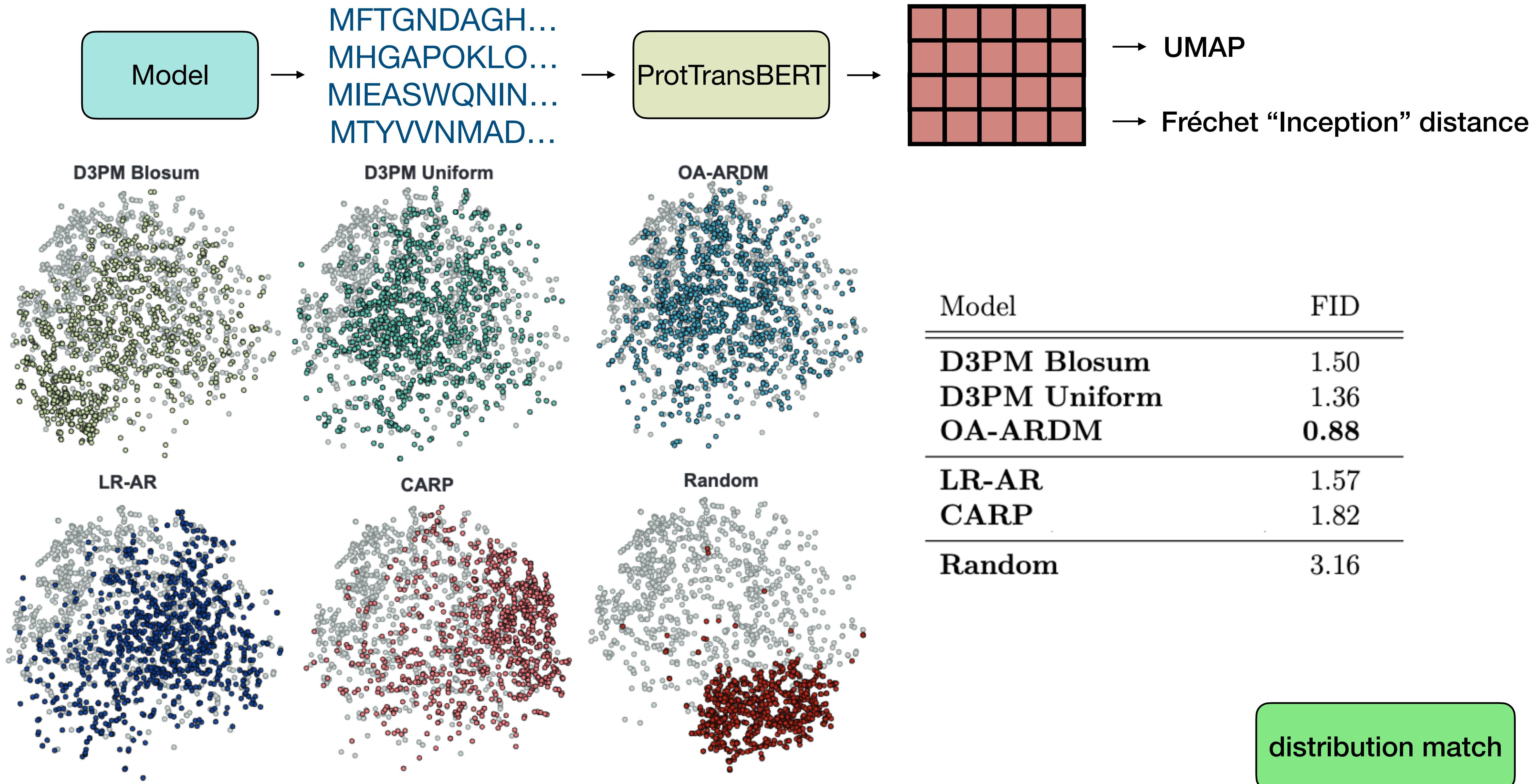
# OA-ARDM produces more structured and self-consistent sequences



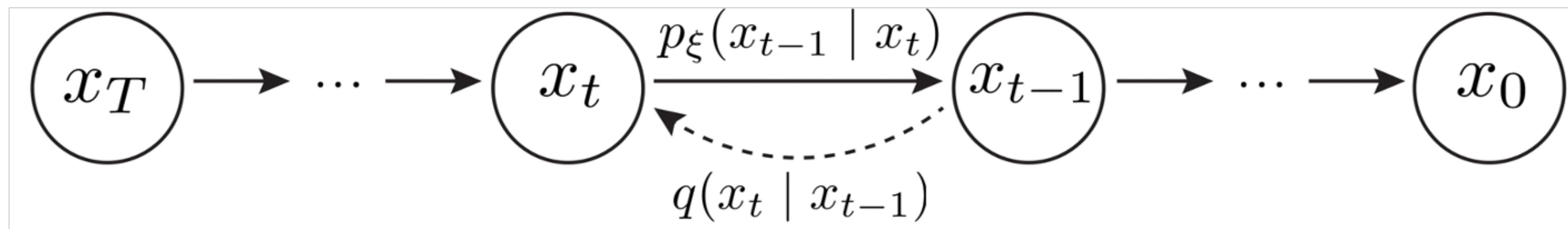
# OA-ARDM produces more structured and self-consistent sequences



# OA-ARDM better recapitulates the natural sequence distribution



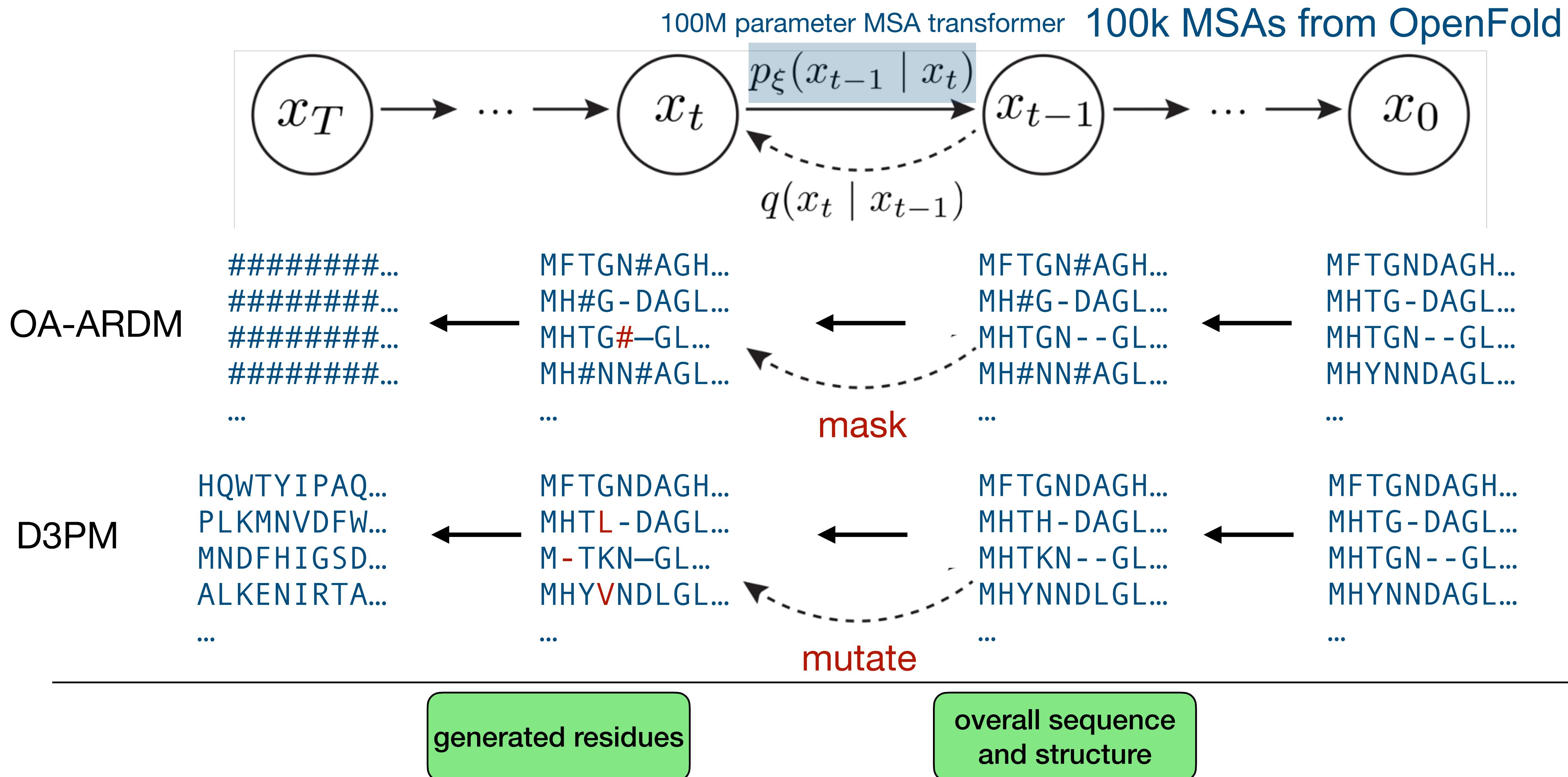
# Generate multiple sequence alignments to leverage evolutionary information



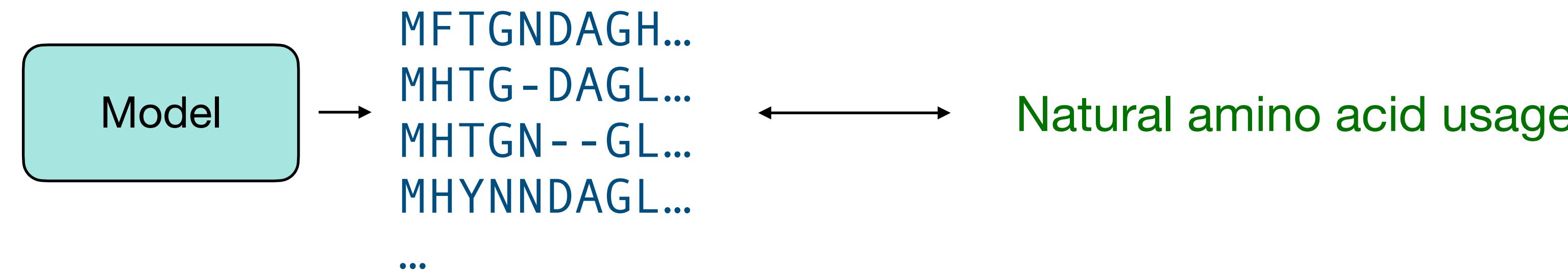
query MFTGNDAGH...

homology  
search

# Generate multiple sequence alignments to leverage evolutionary information



# OA-ARDM better recapitulates the natural amino acid distribution

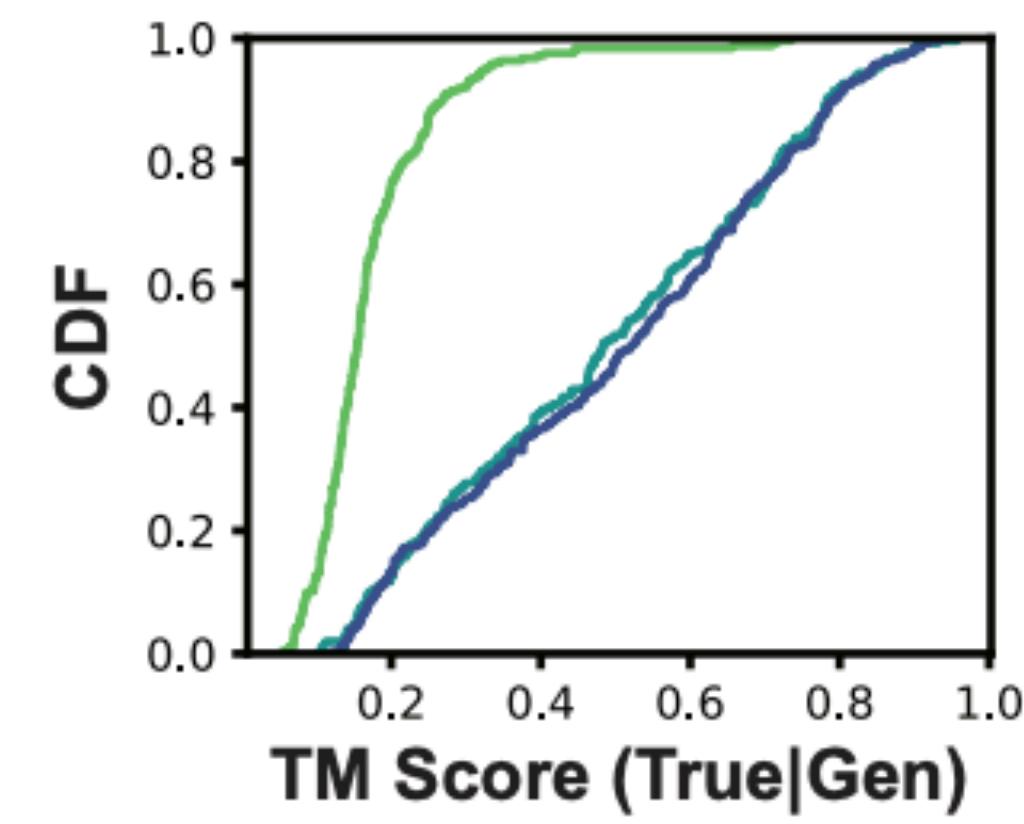
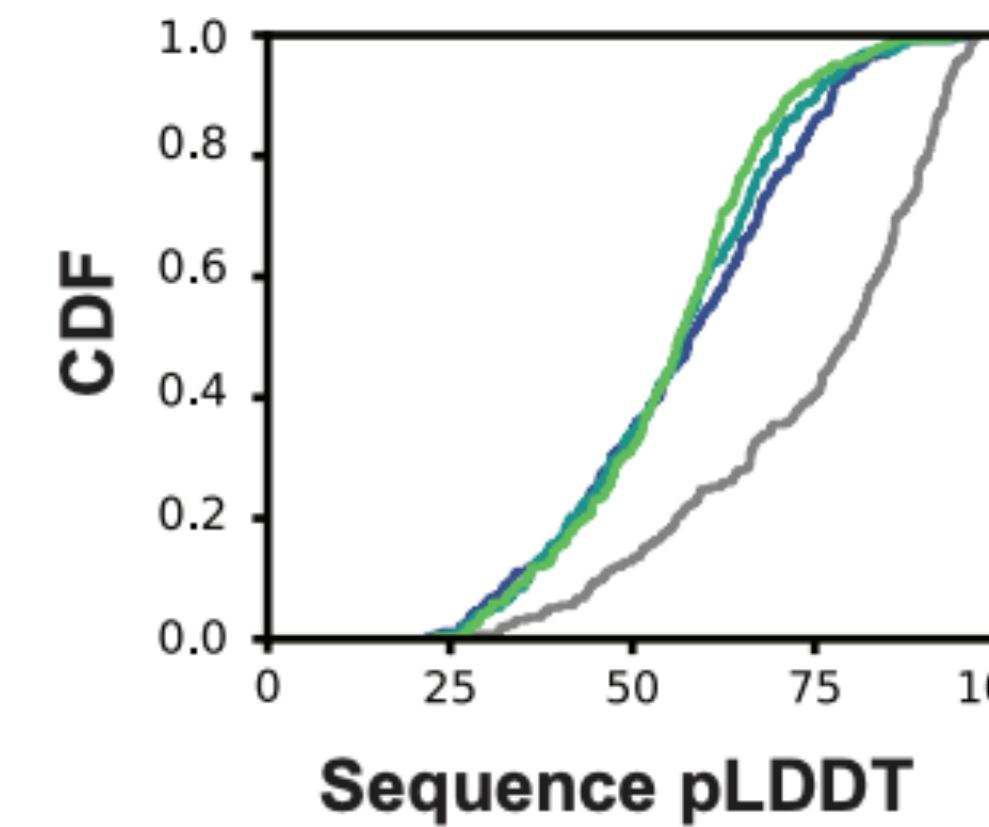
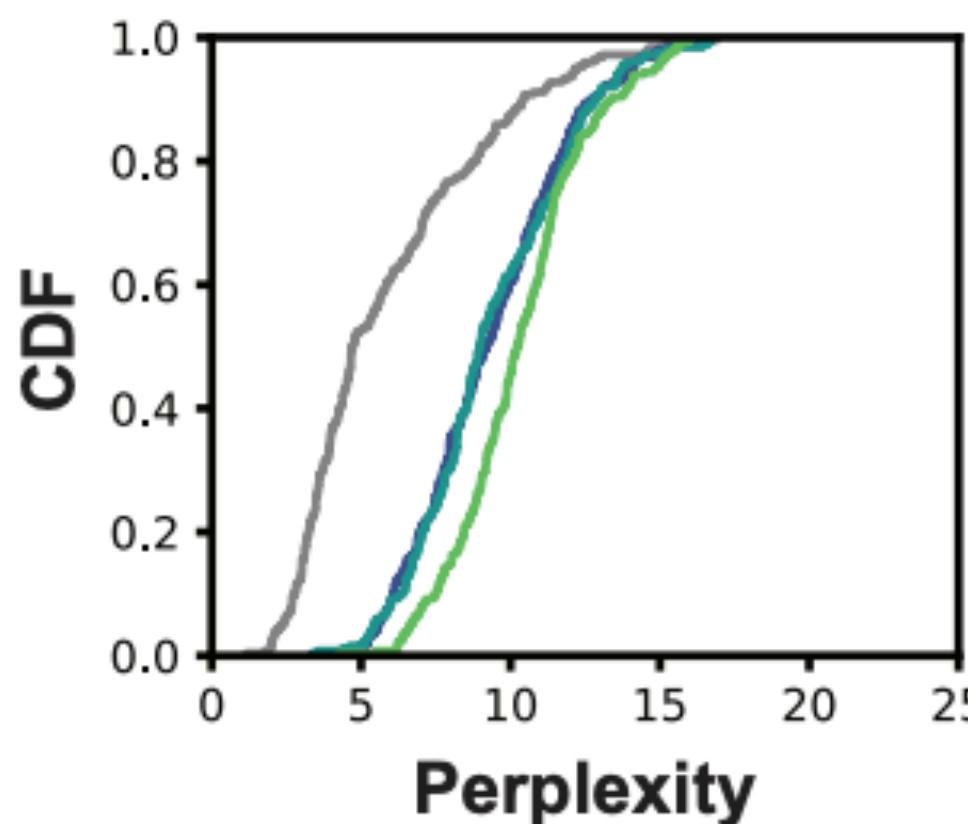
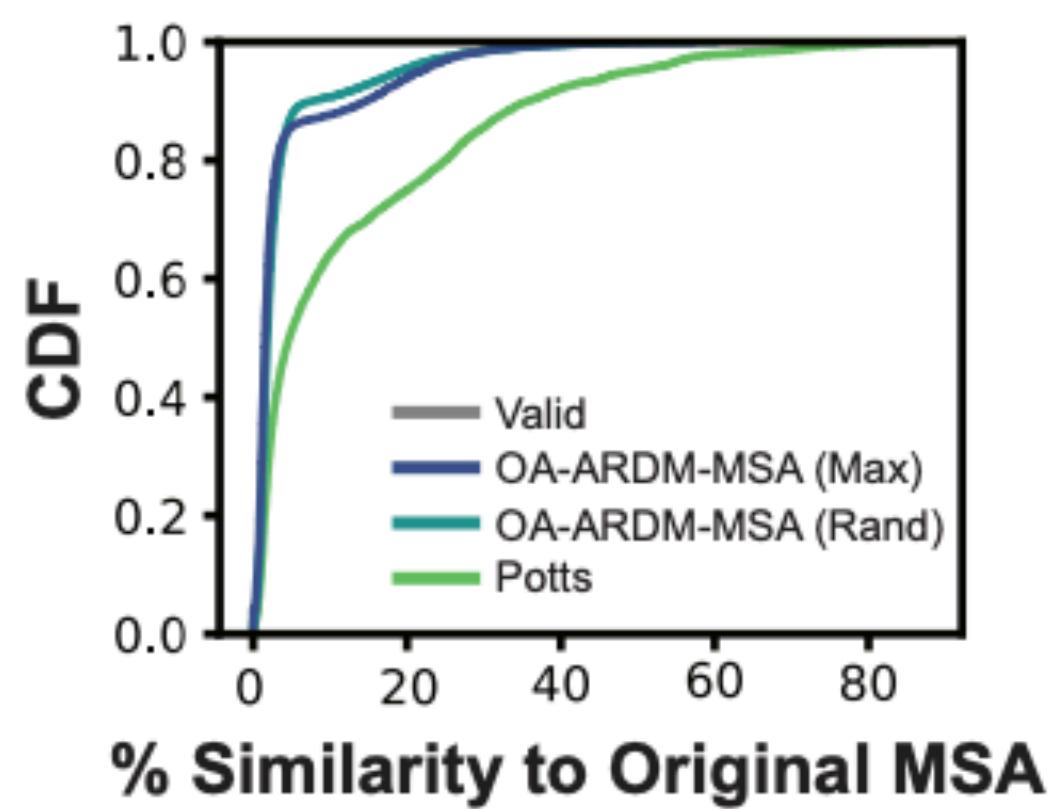
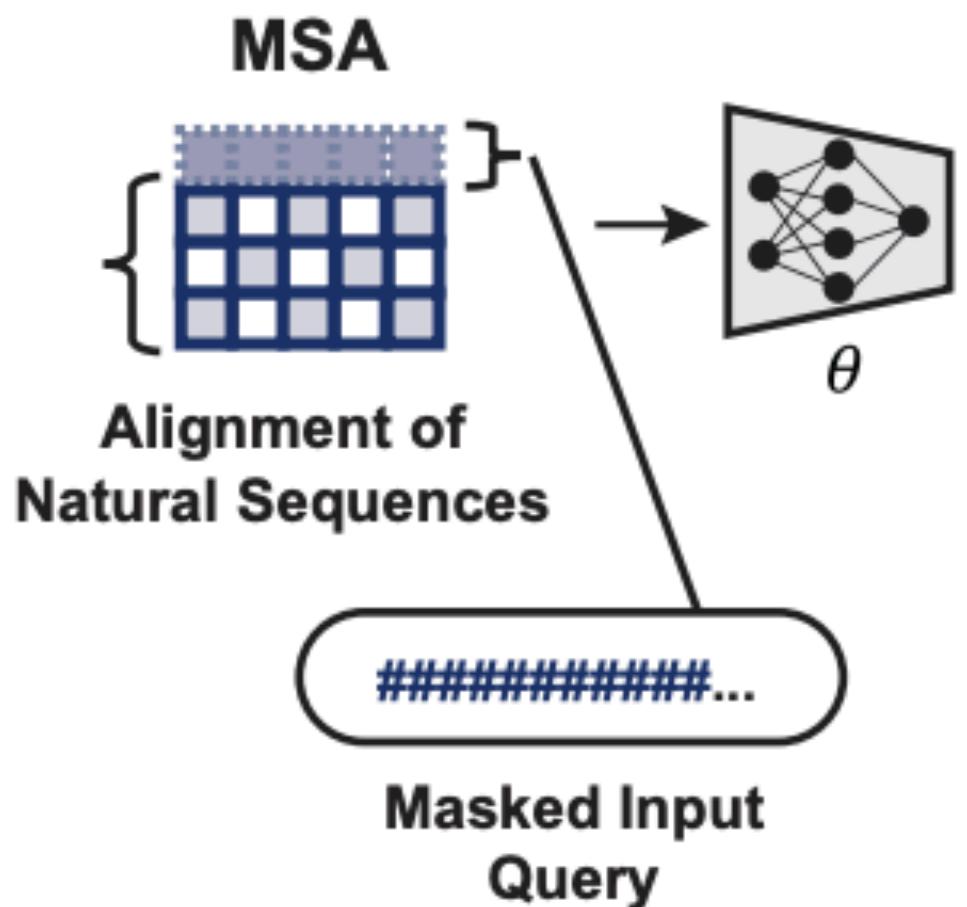


Model	Params	KL (AA+GAP)	KL (AA)
<b>Valid</b>	-	0.056	0.00093
ESM-MSA-1b (Rao et al., 2021)	100M	>1	> 1

generated residues

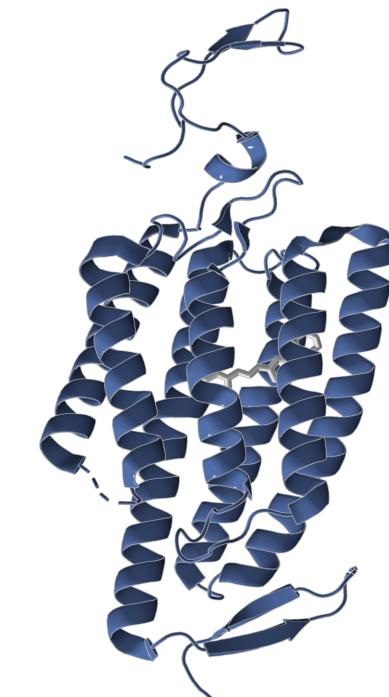
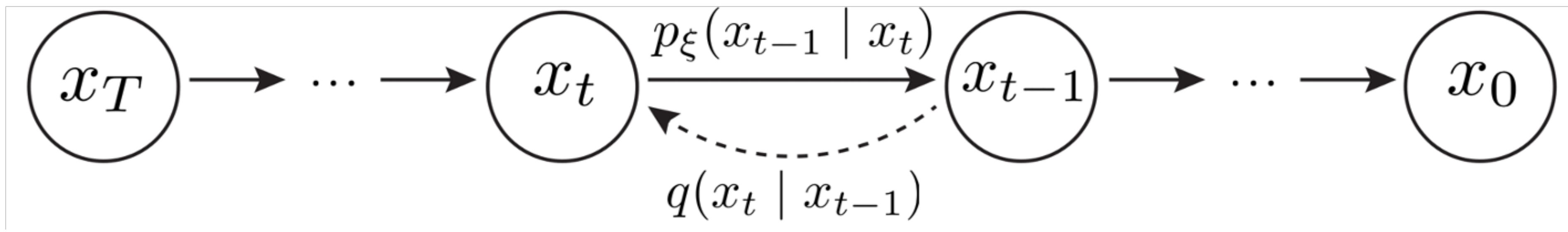
# OA-ARDM can generate realistic queries conditioned on the MSA

## Conditional Query Generation



overall sequence  
and structure

# We use diffusion models to generate structure, sequence, and MSA

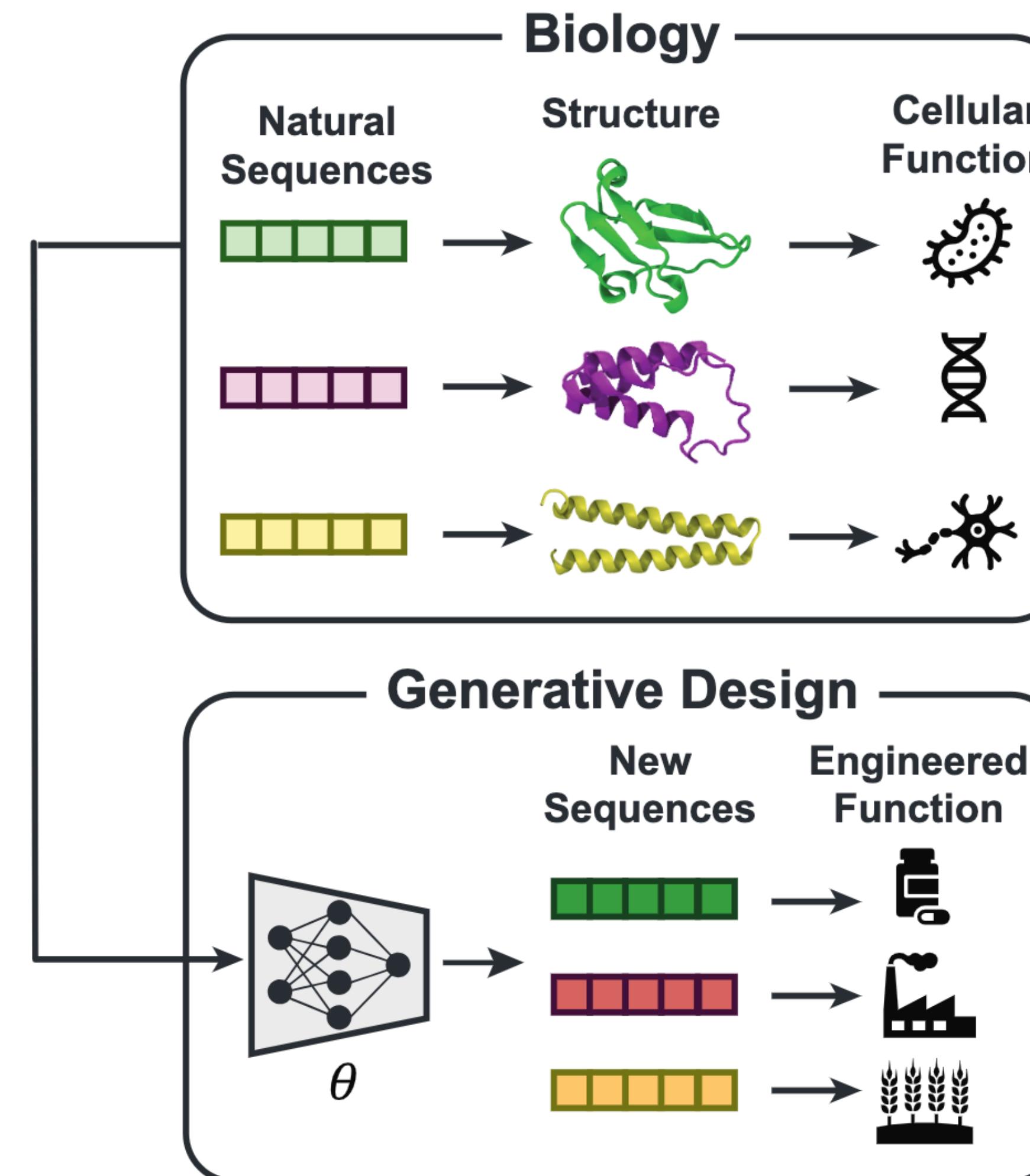
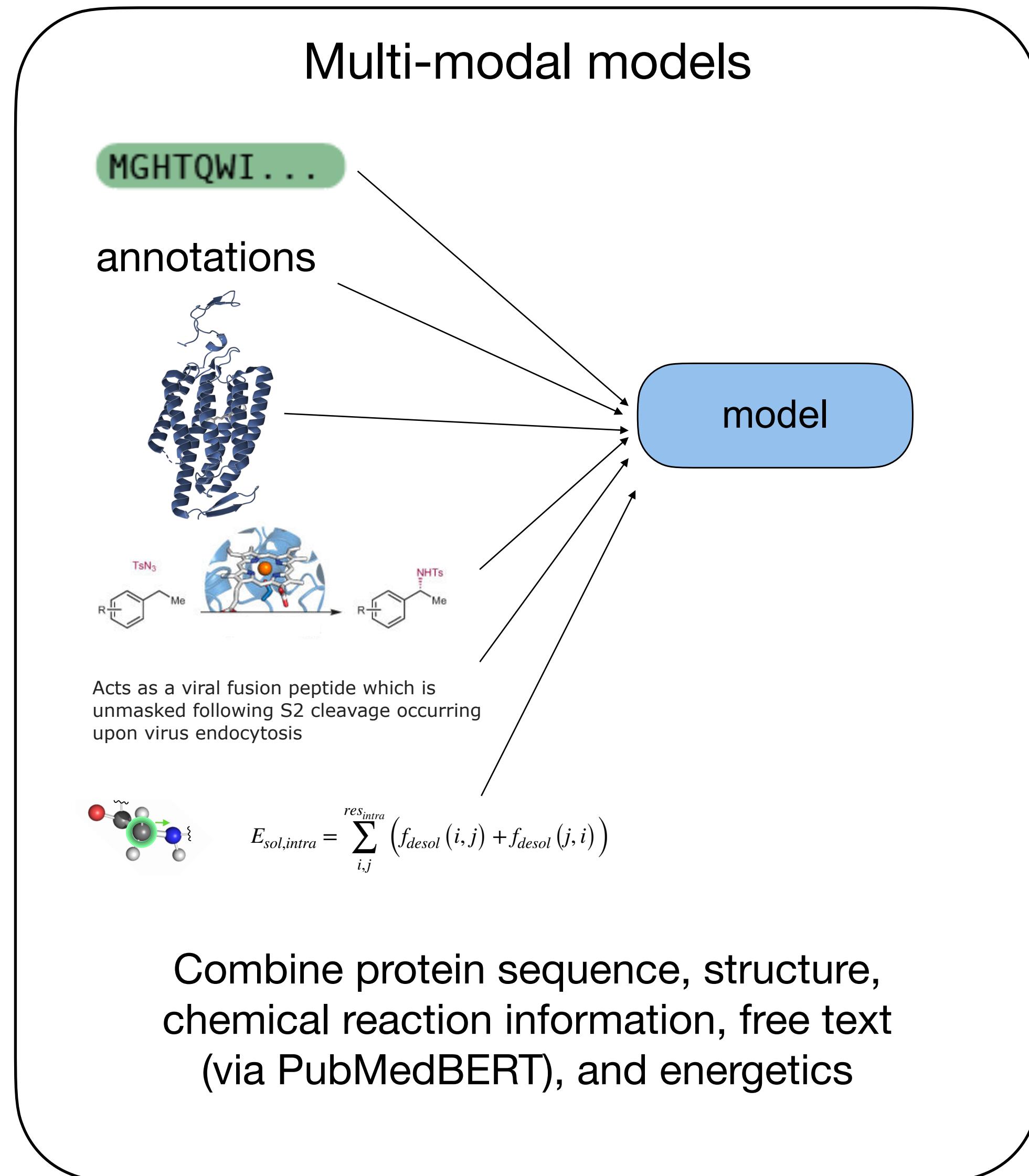


MFTGNDAGH...  
MHGAPOKLO...  
MIEASWQNI...  
MTYVVNMAD...

MFTGNDAGH...  
MHTG-DAGL...  
MHTGN--GL...  
MHYNNDAGL...  
...

Eventually always need sequences!

# Can we generate proteins with new functions?



# Acknowledgments



BioML at MSR New England