

# Pretraining and Representation Learning on Proteins

Kevin Kaichuang Yang  
Microsoft Research New England  
 @KevinKaichuang



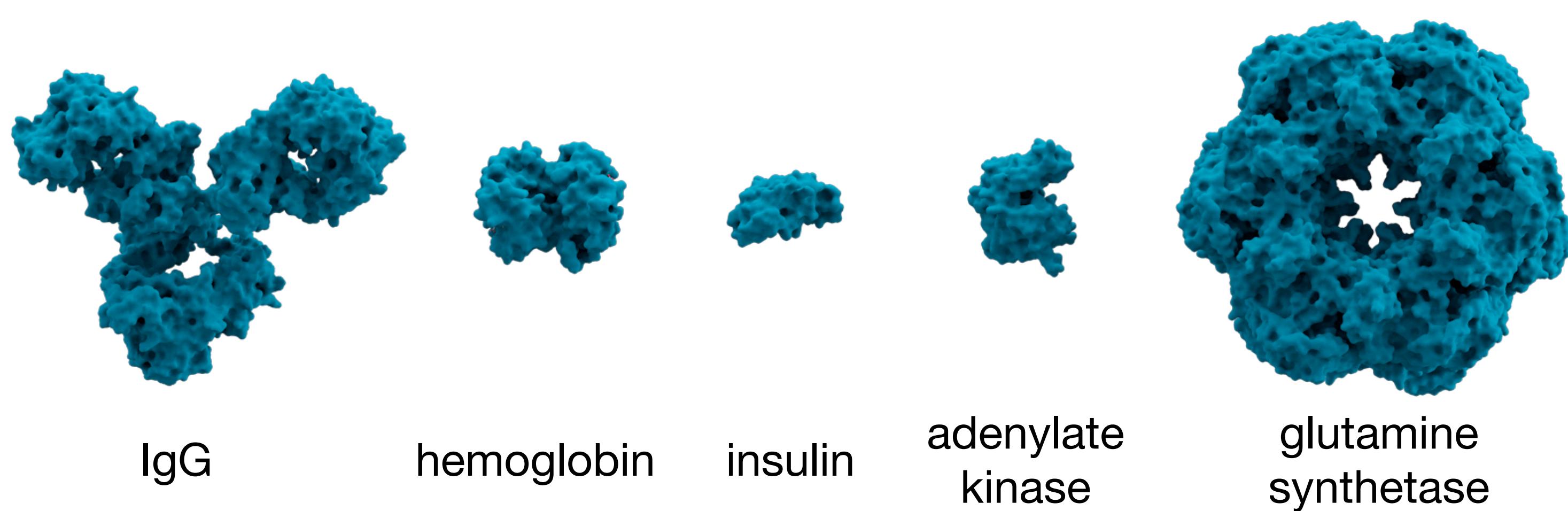
# Proteins are biology's actuators

# Proteins are biology's actuators

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

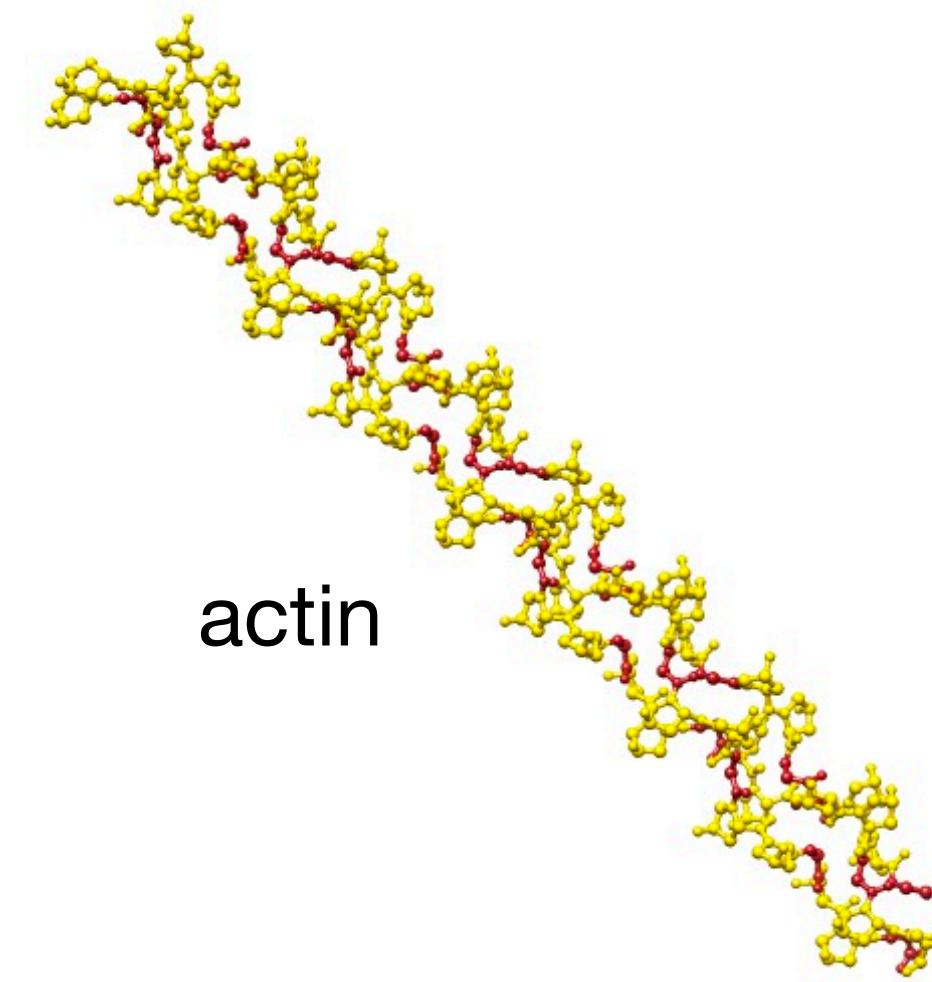
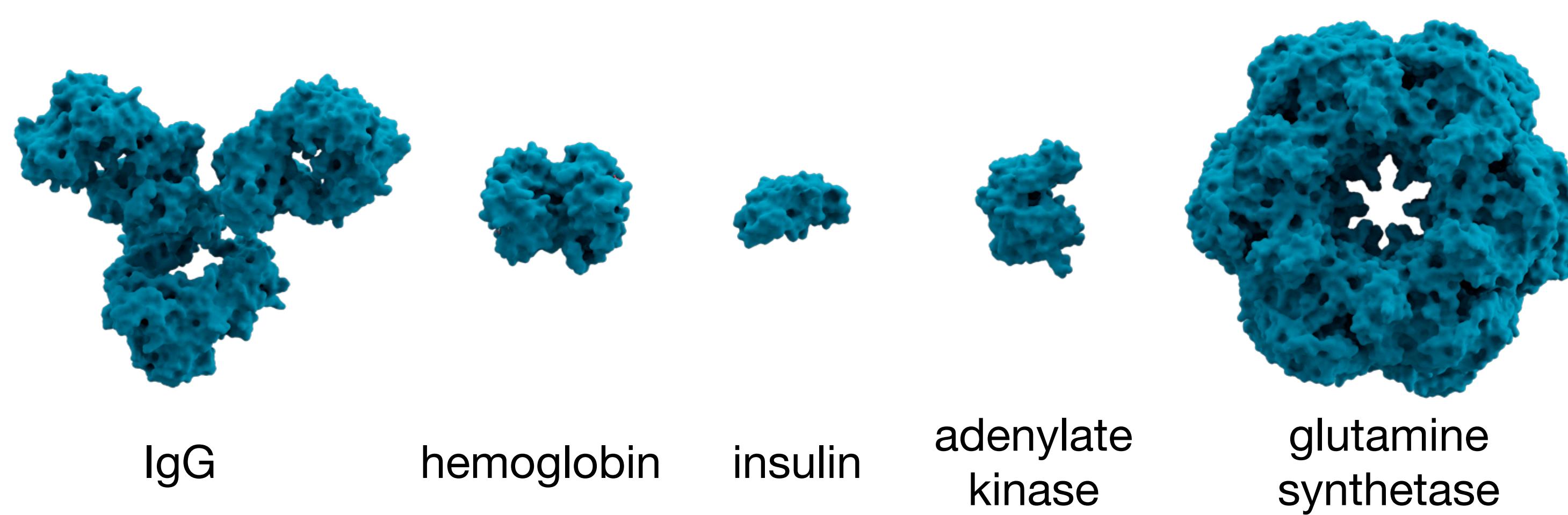
# Proteins are biology's actuators

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



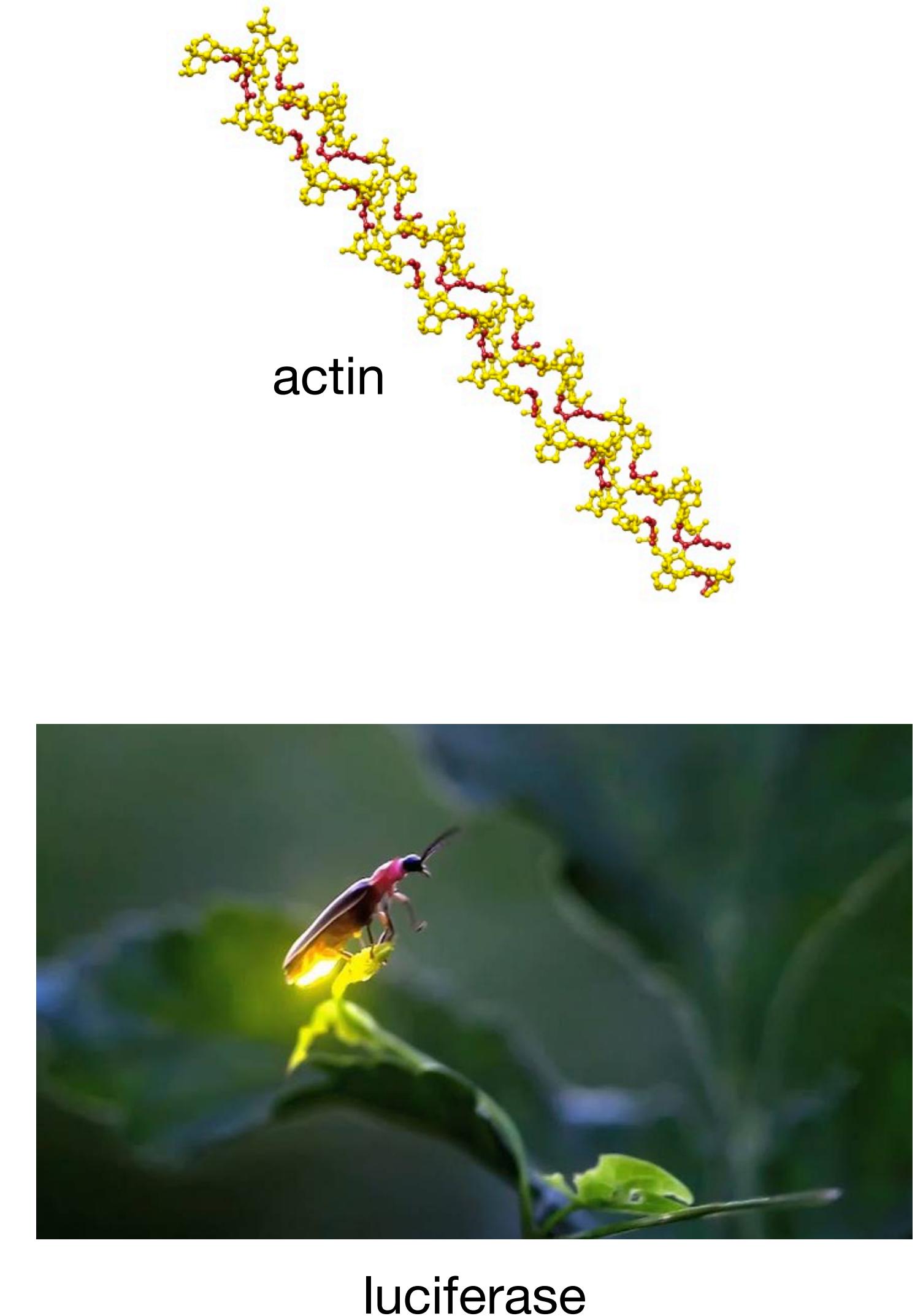
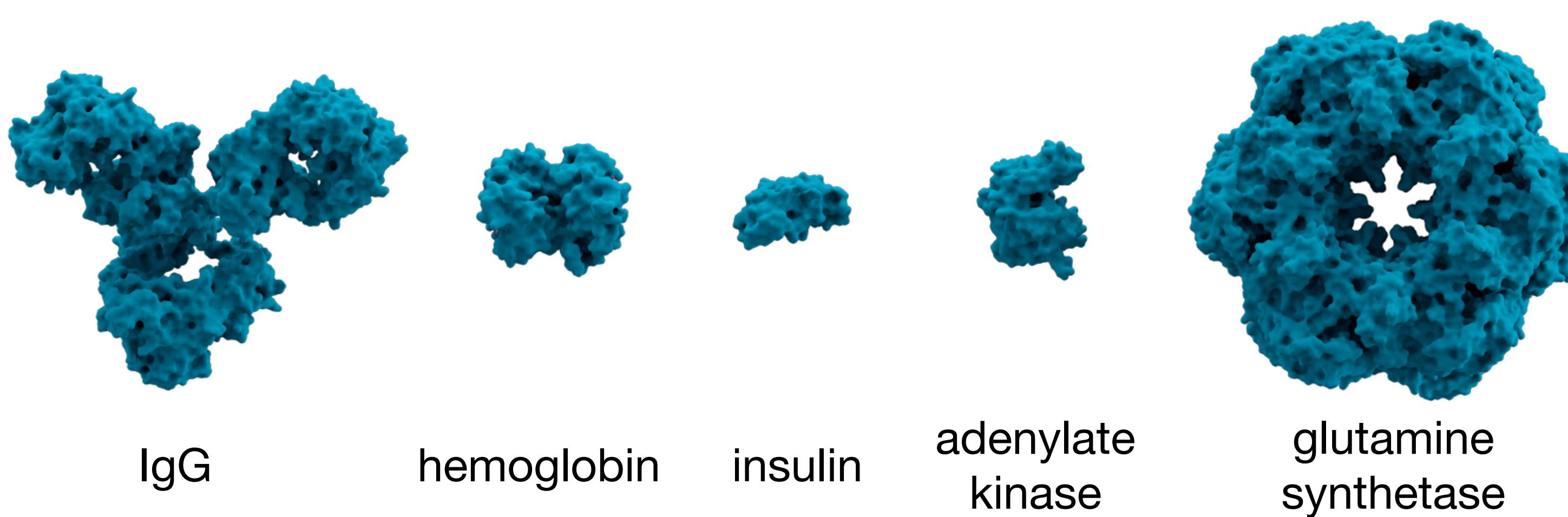
# Proteins are biology's actuators

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

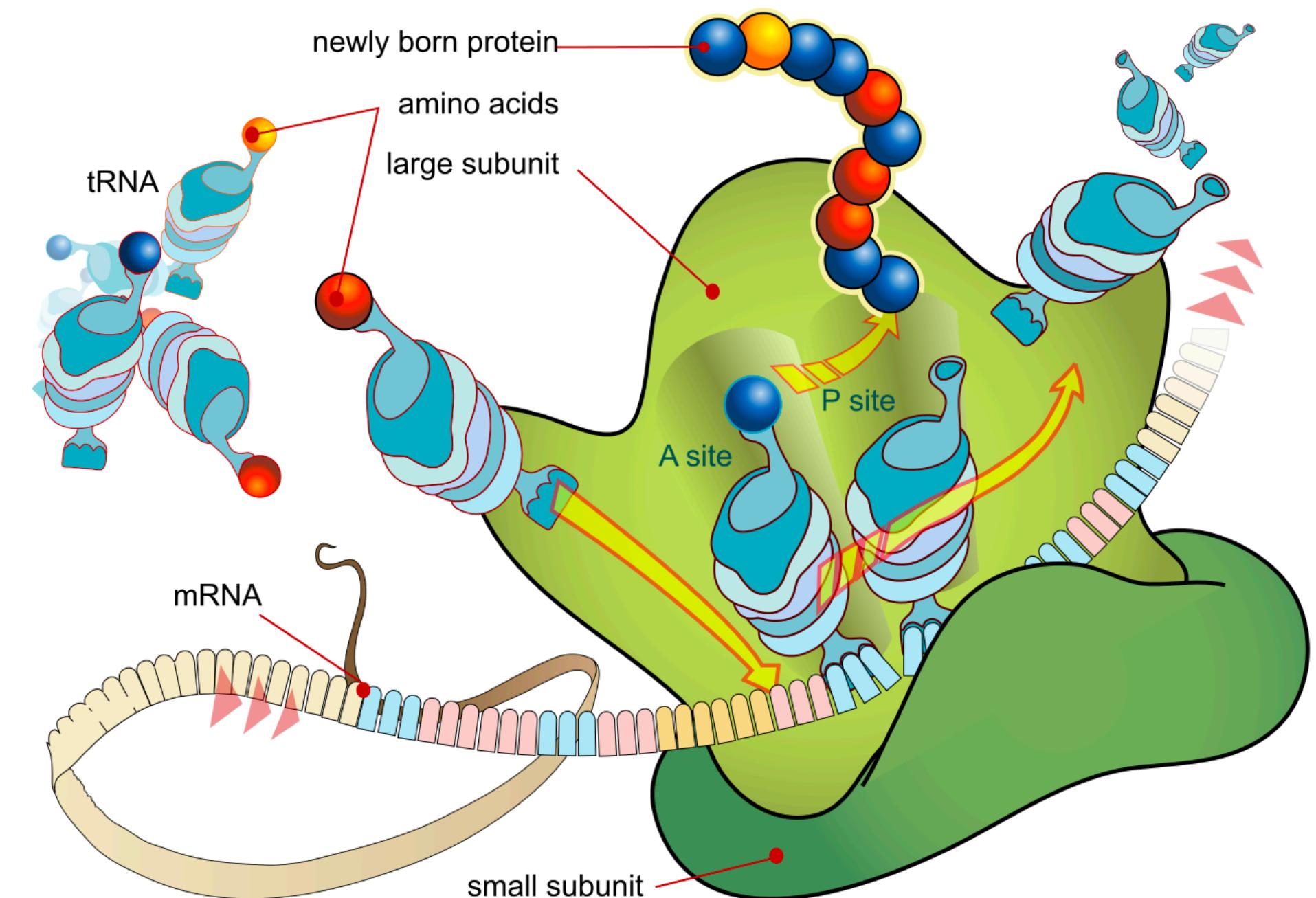
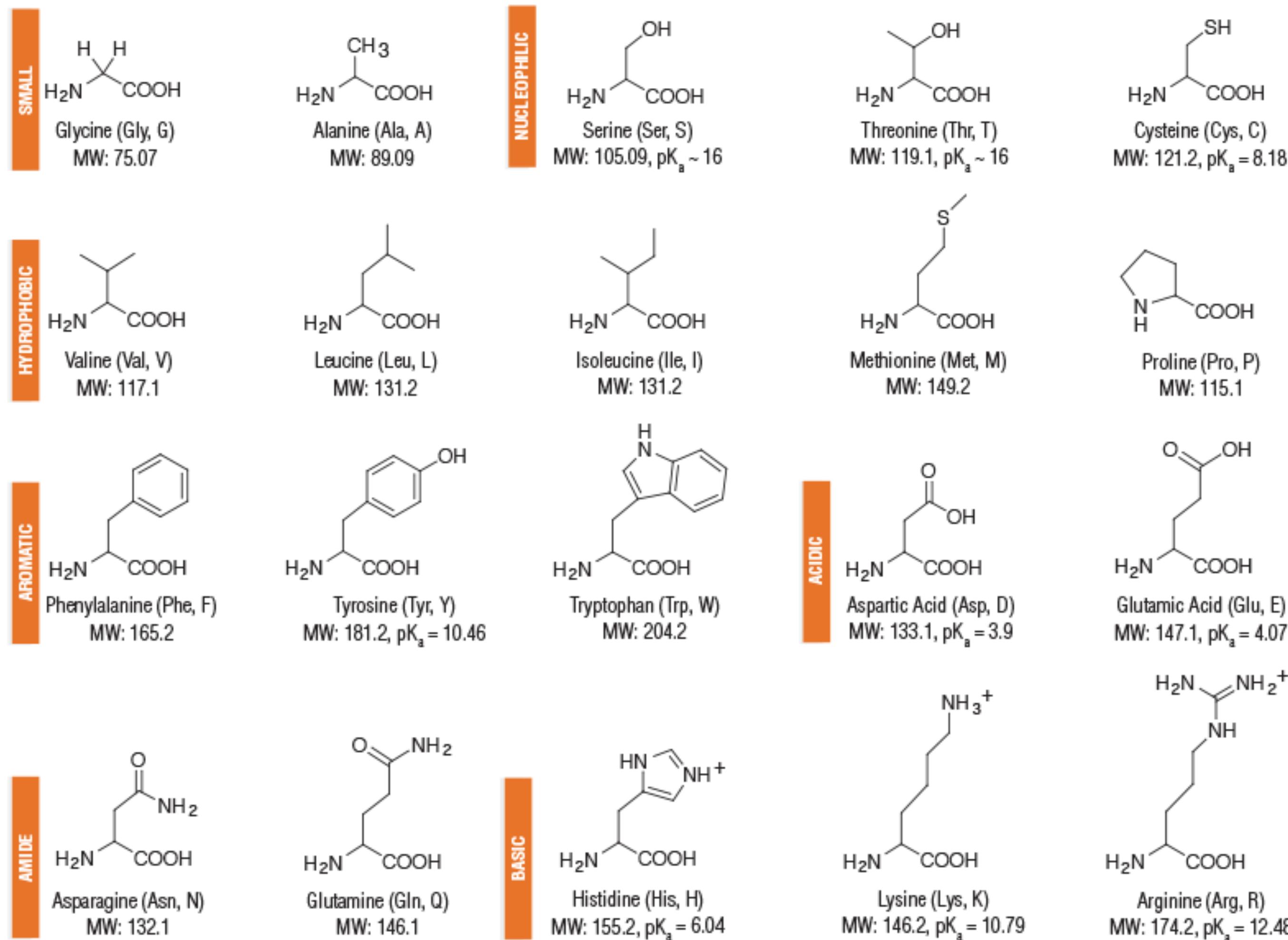


# Proteins are biology's actuators

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



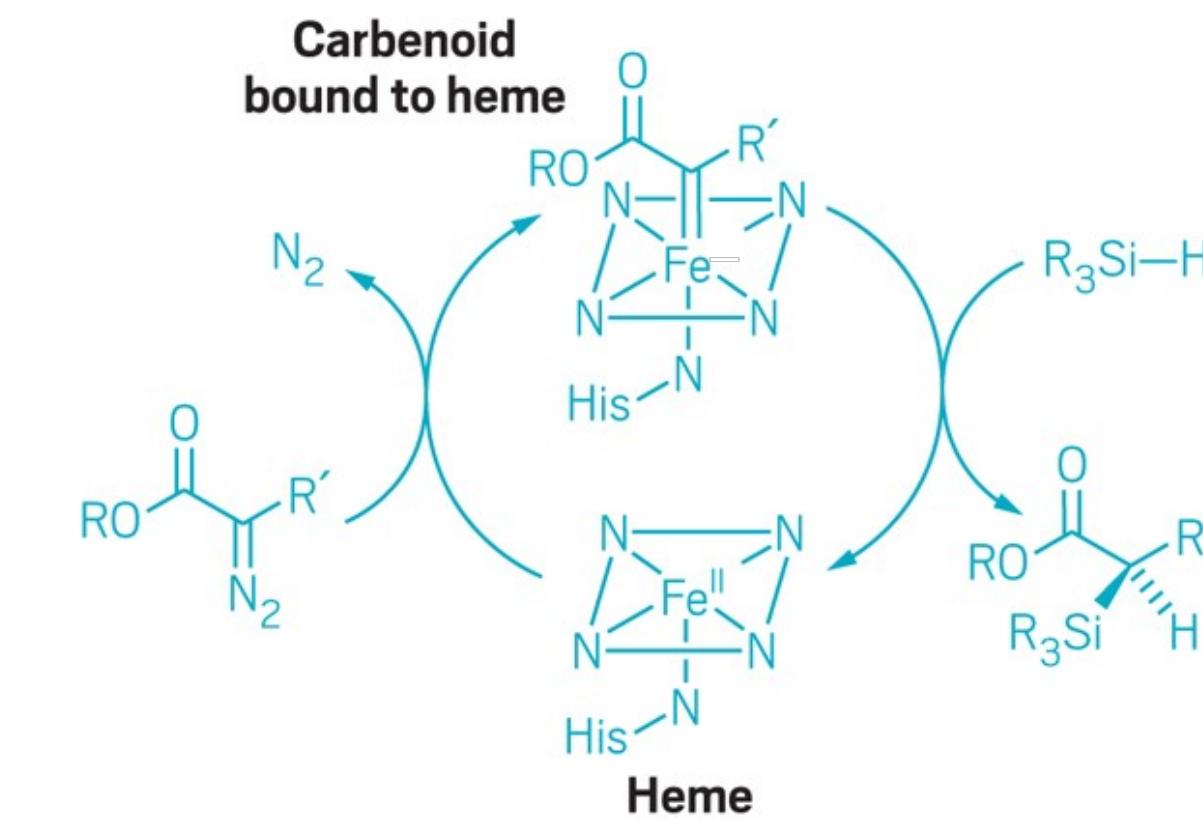
# Diversity arises from 20 building blocks



# Why design proteins?

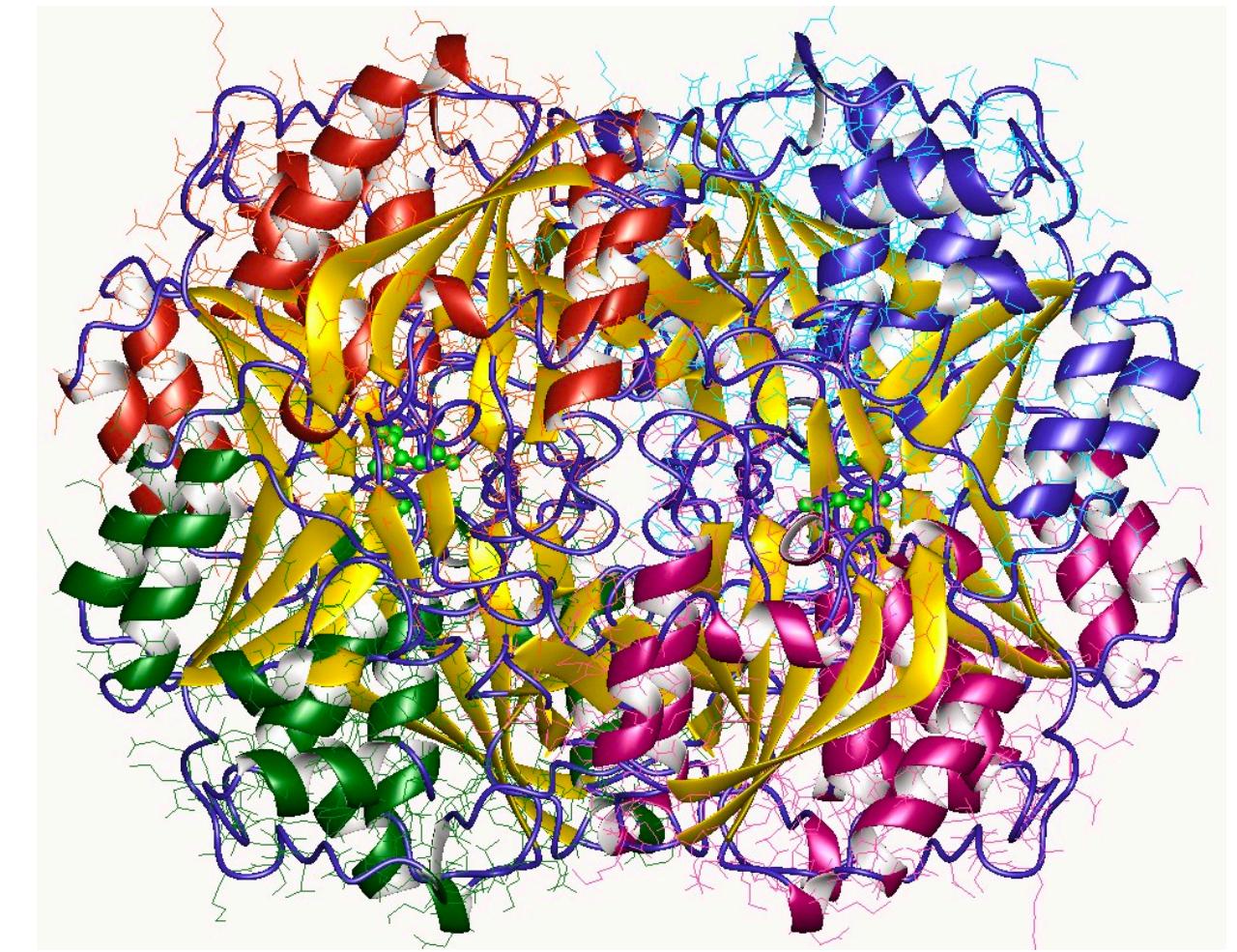
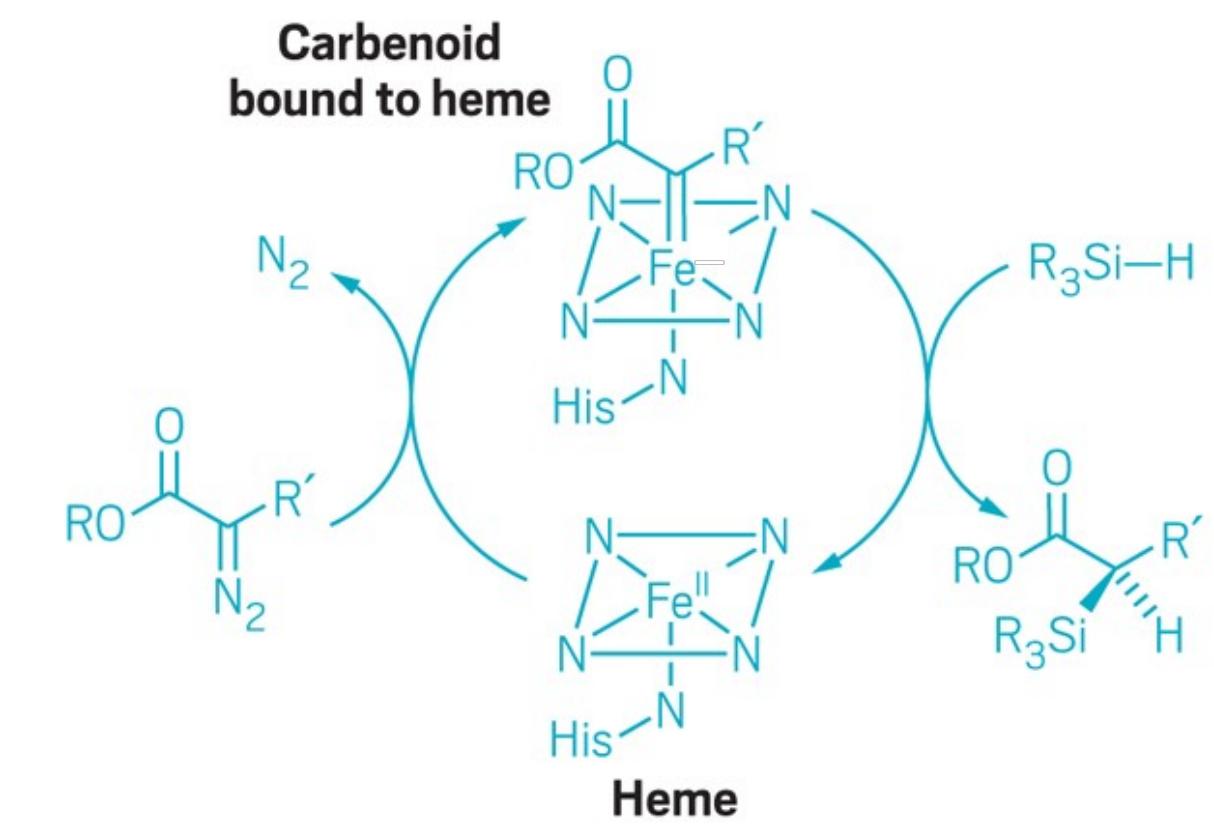
# Why design proteins?

- New chemistry



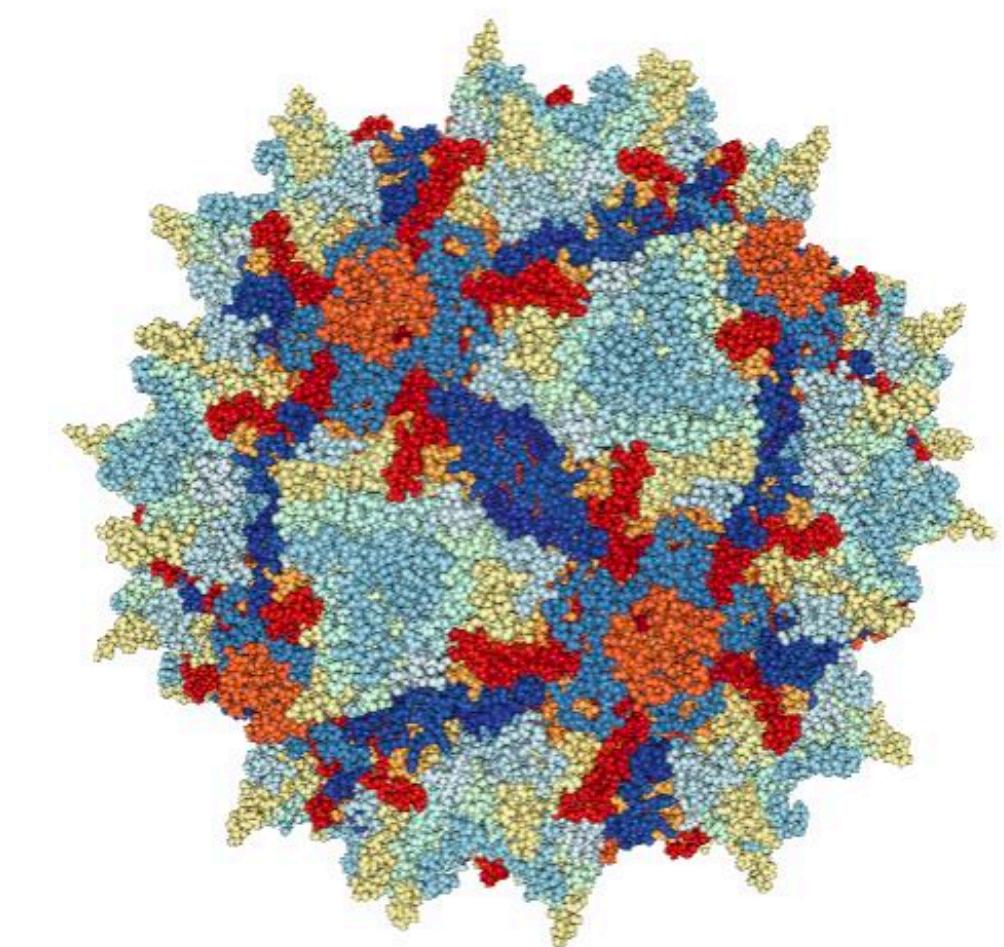
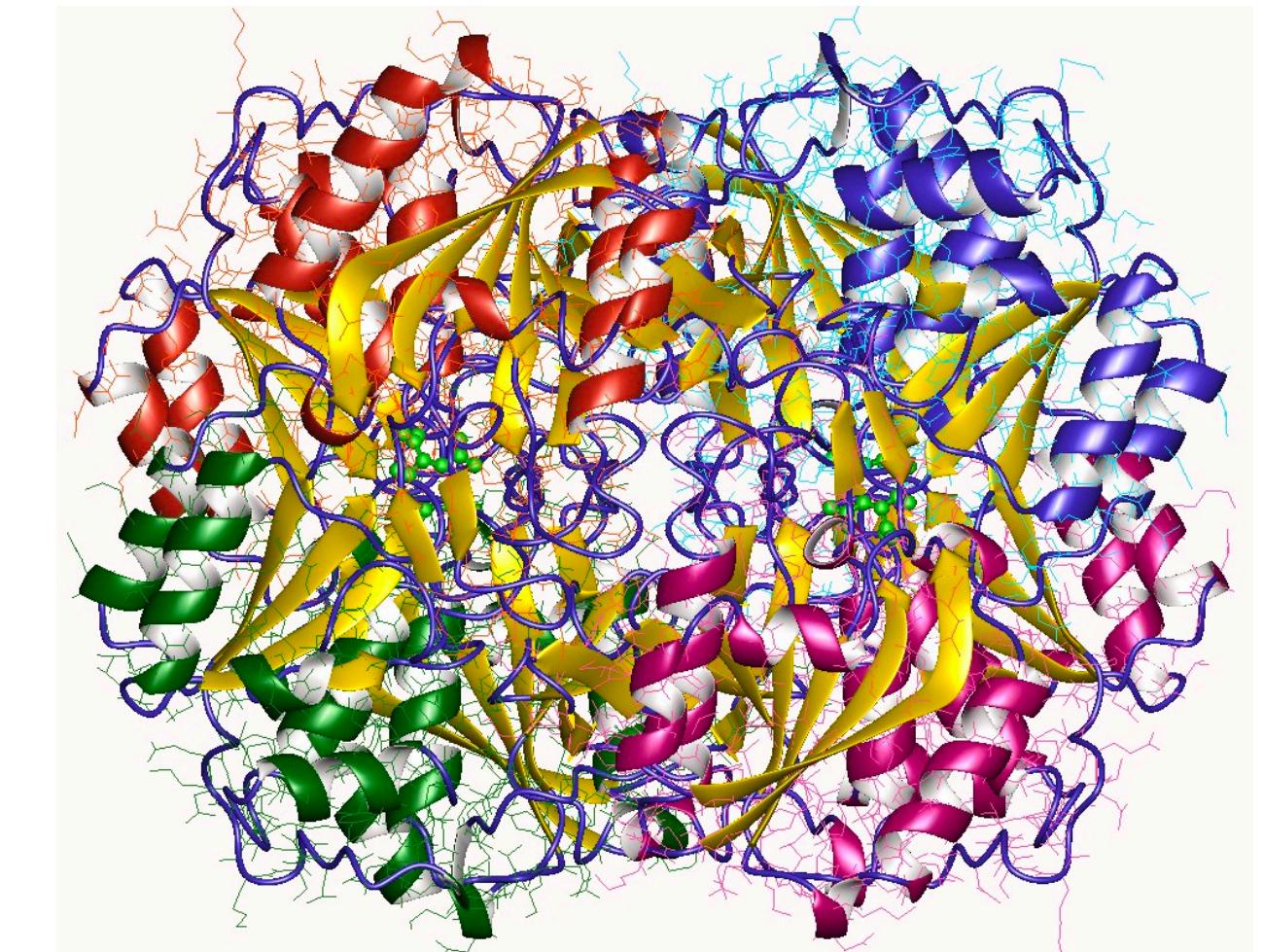
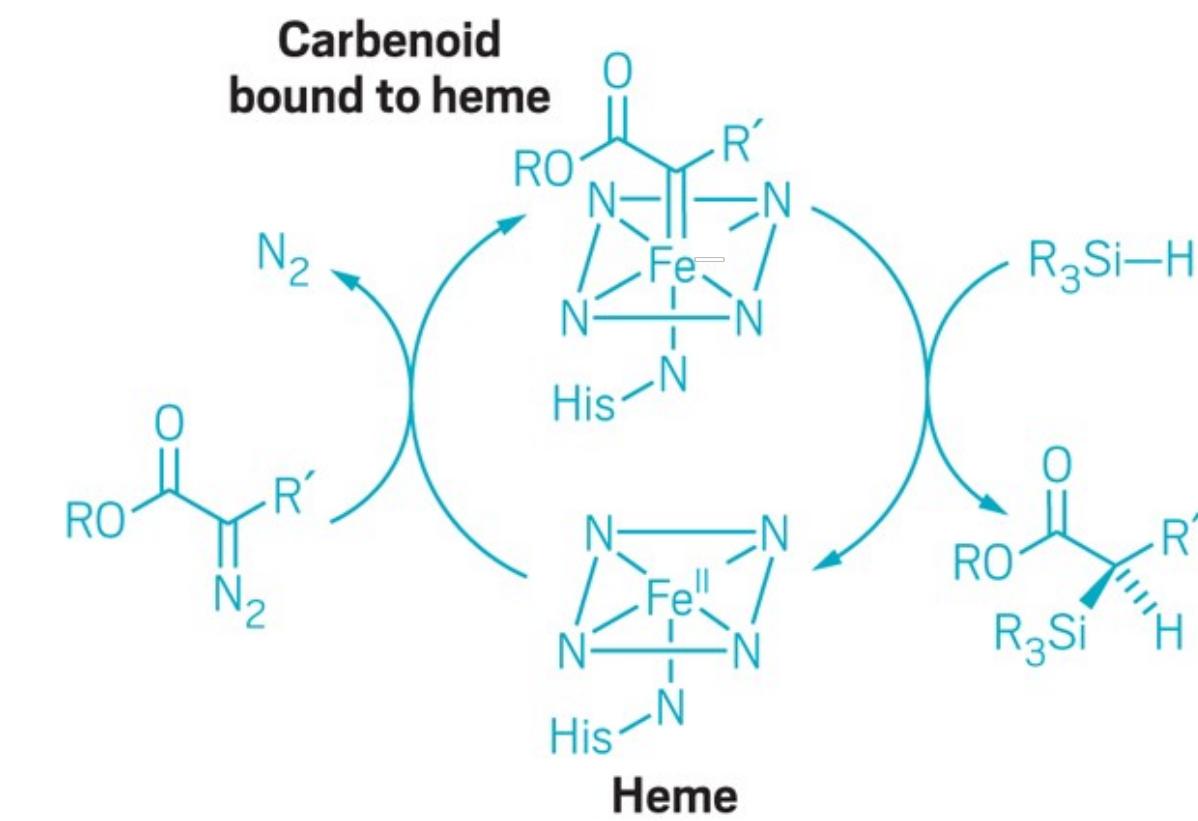
# Why design proteins?

- New chemistry
- Therapeutics



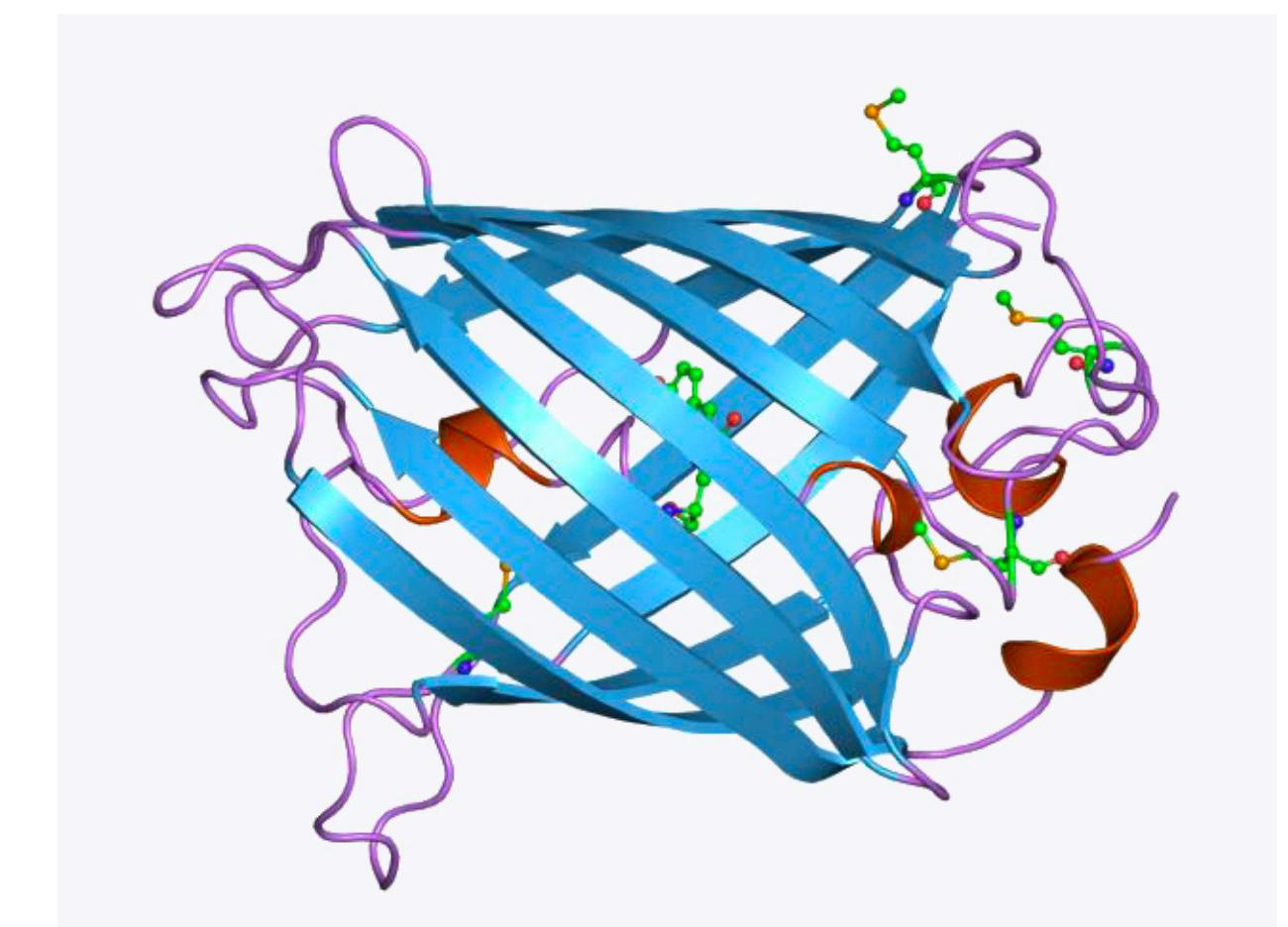
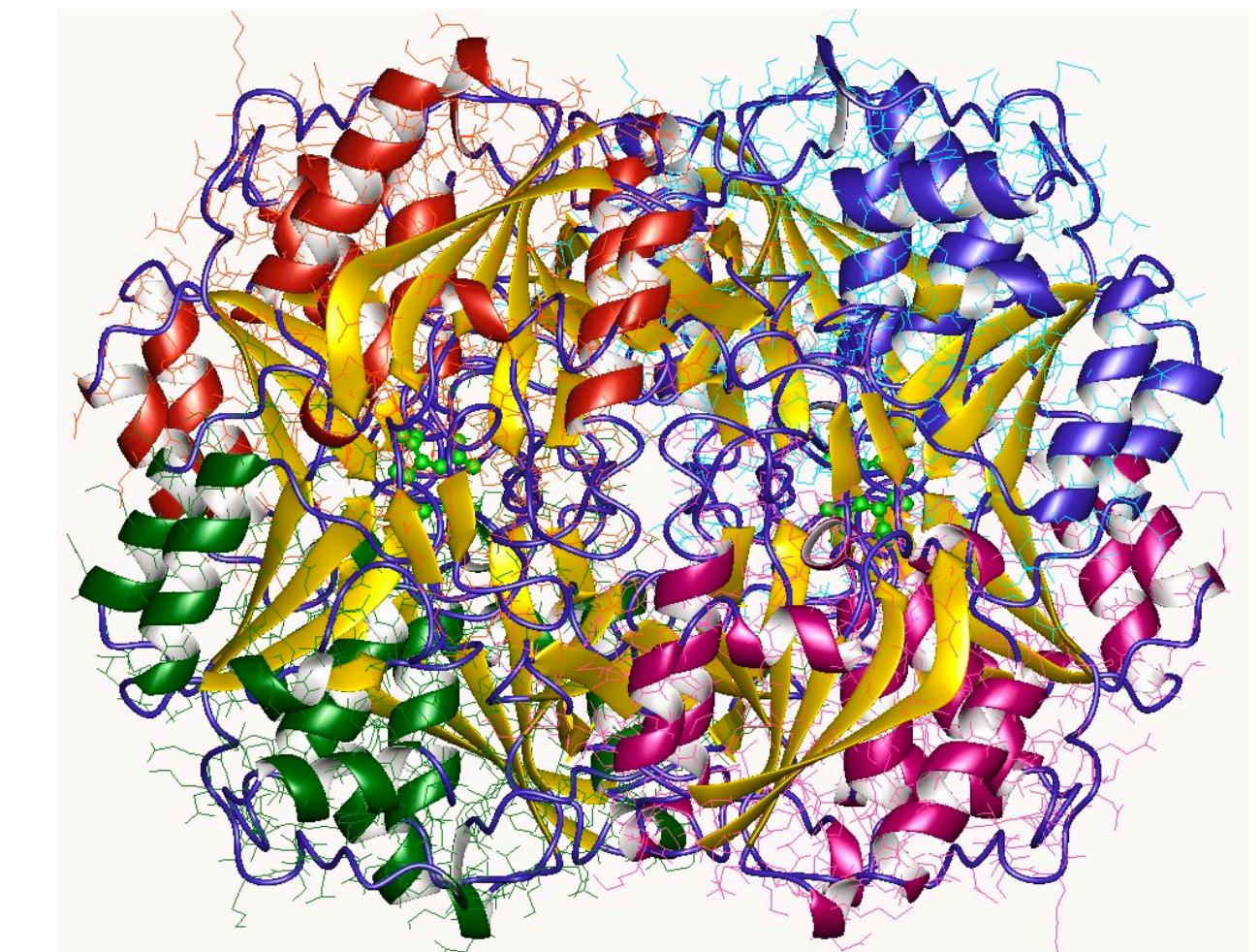
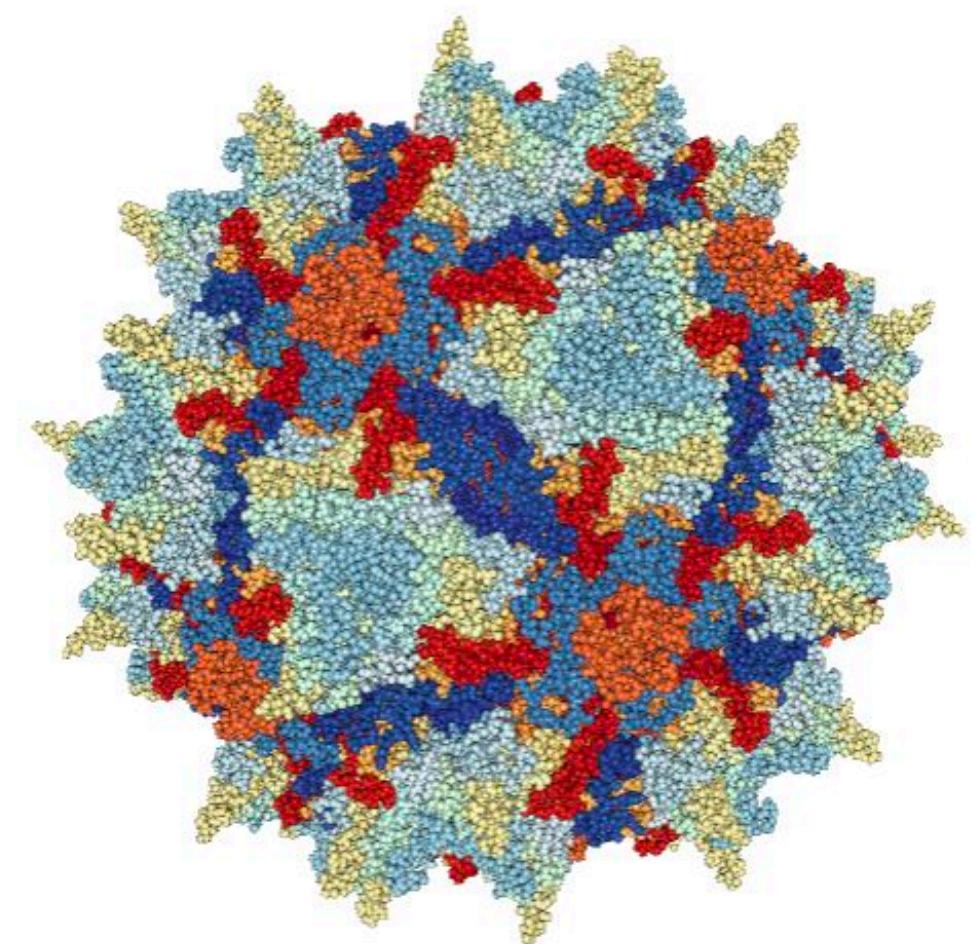
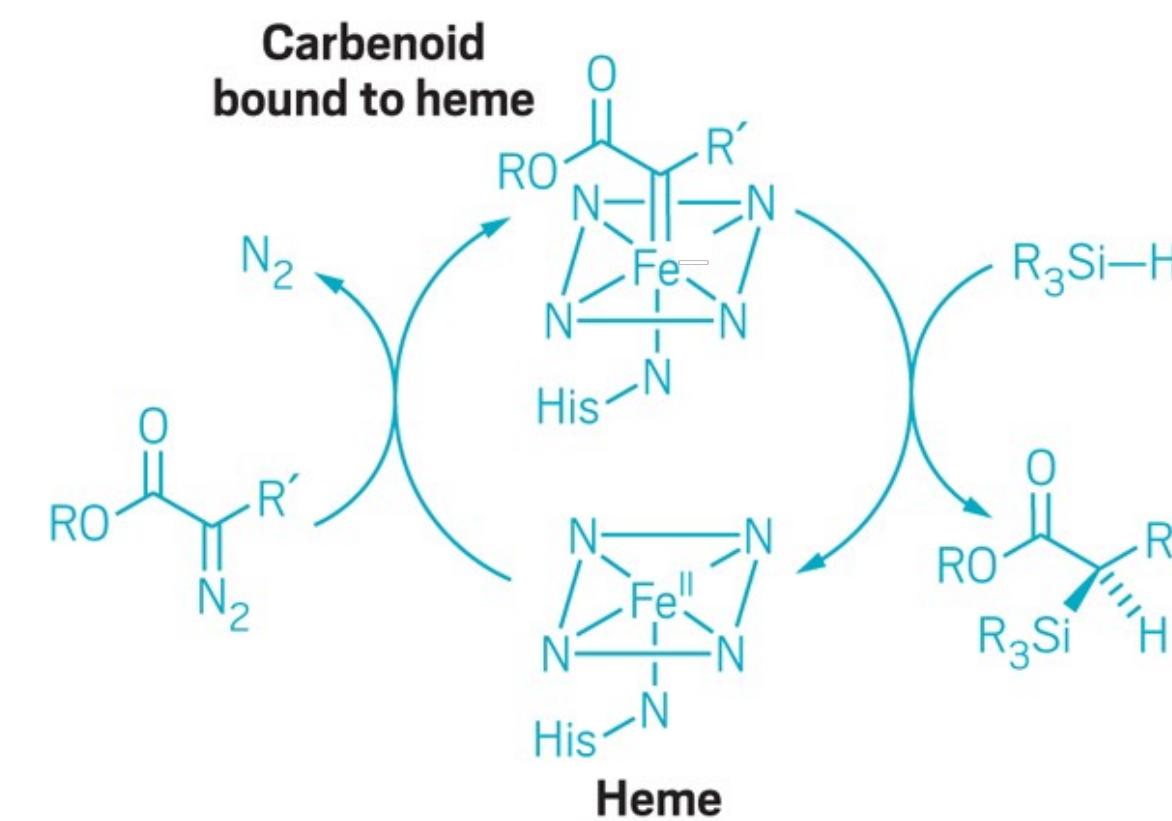
# Why design proteins?

- New chemistry
- Therapeutics
- To learn how they work



# Why design proteins?

- New chemistry
- Therapeutics
- To learn how they work
- Molecular tools



# The protein design problem

# The protein design problem

MGTGDHDD...

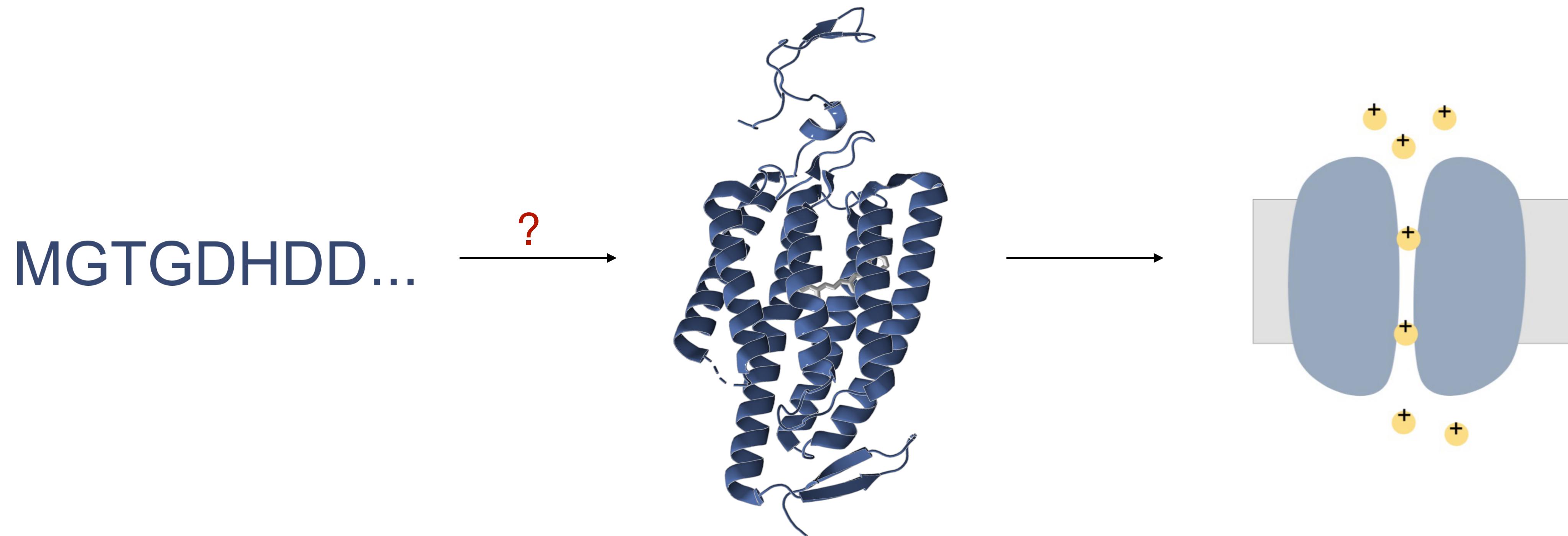
# The protein design problem



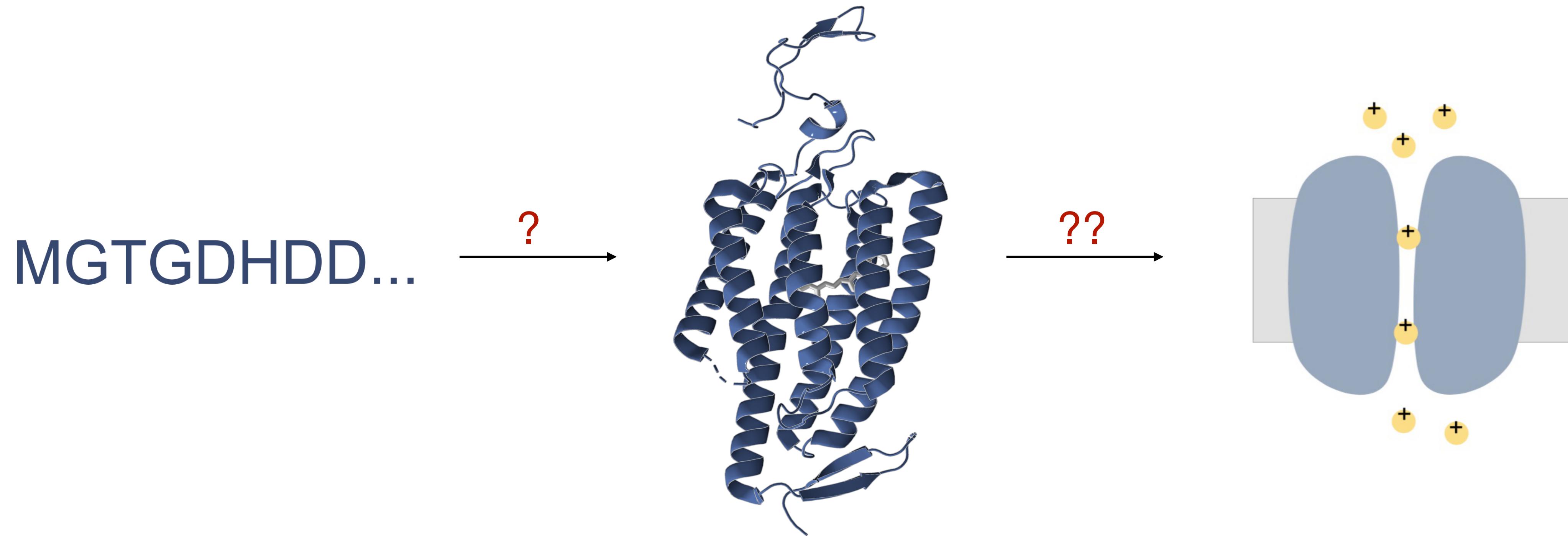
# The protein design problem



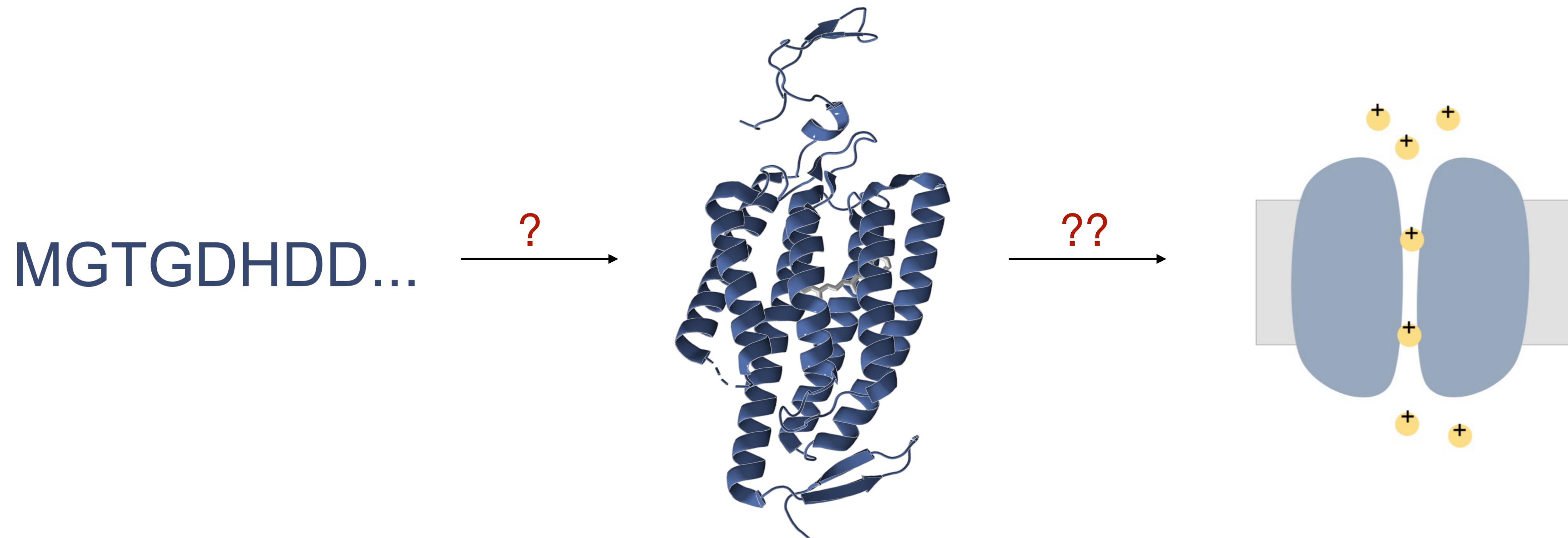
# The protein design problem



# The protein design problem



# The protein design problem



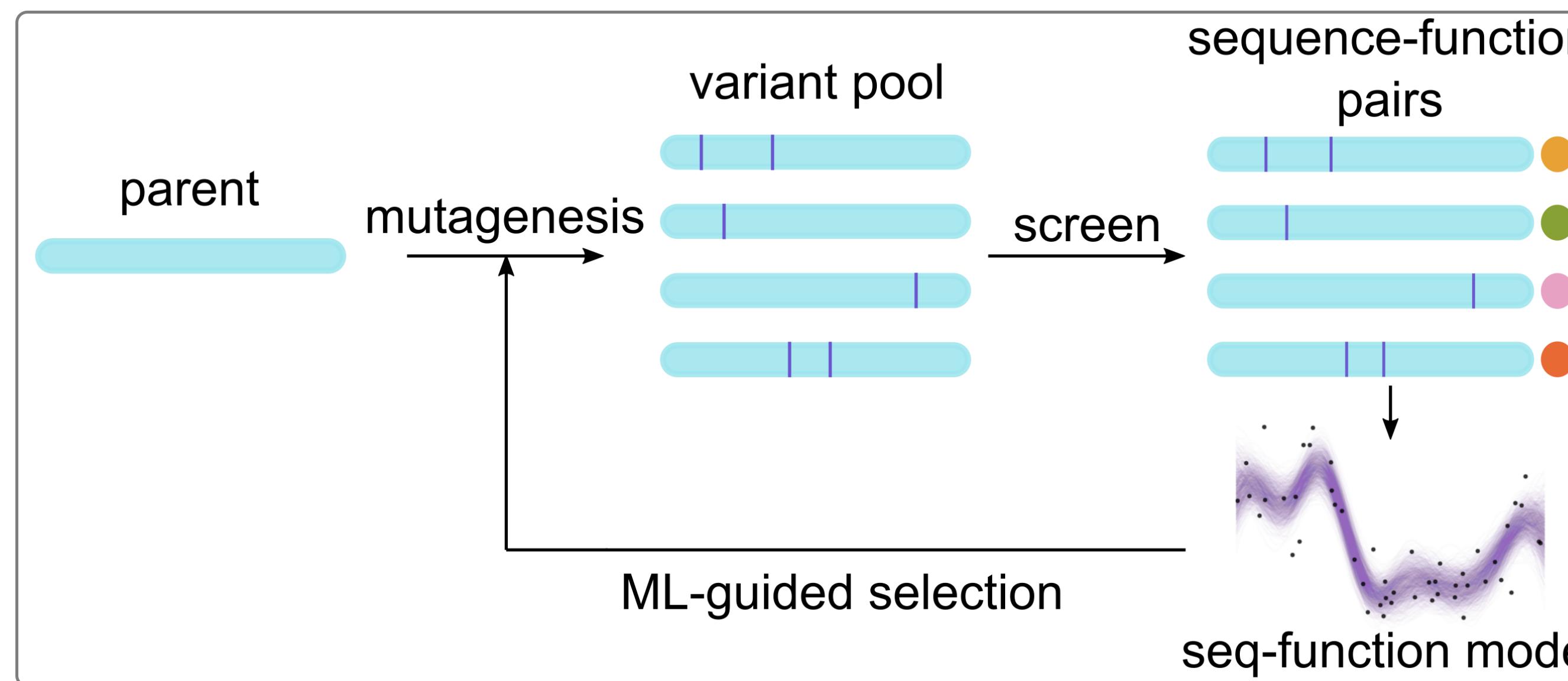
What sequence will give the desired function?

# The protein design problem

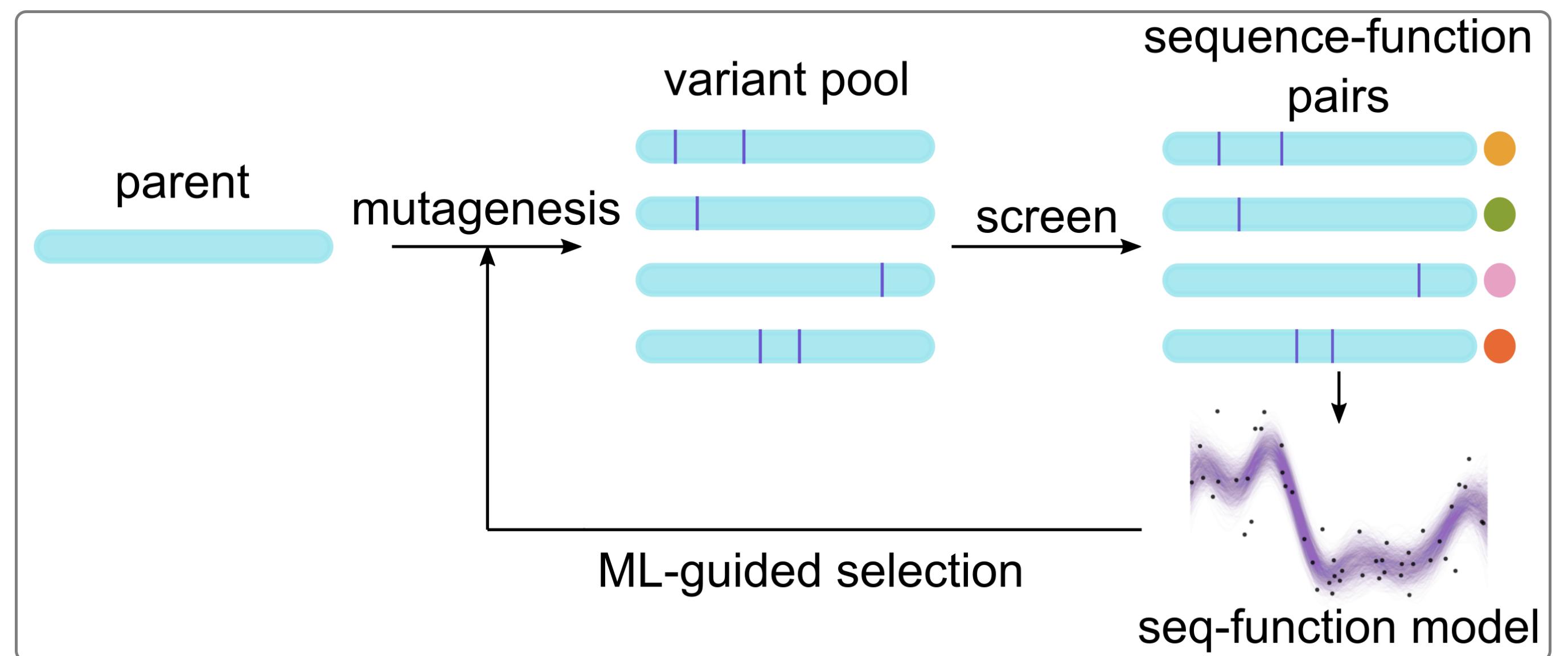


What sequence will give the desired function?

# Machine learning enables optimization with fewer measurements

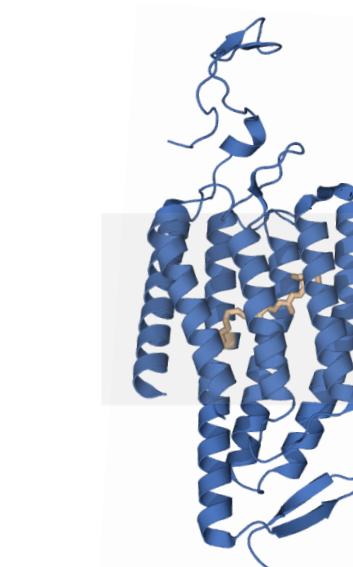
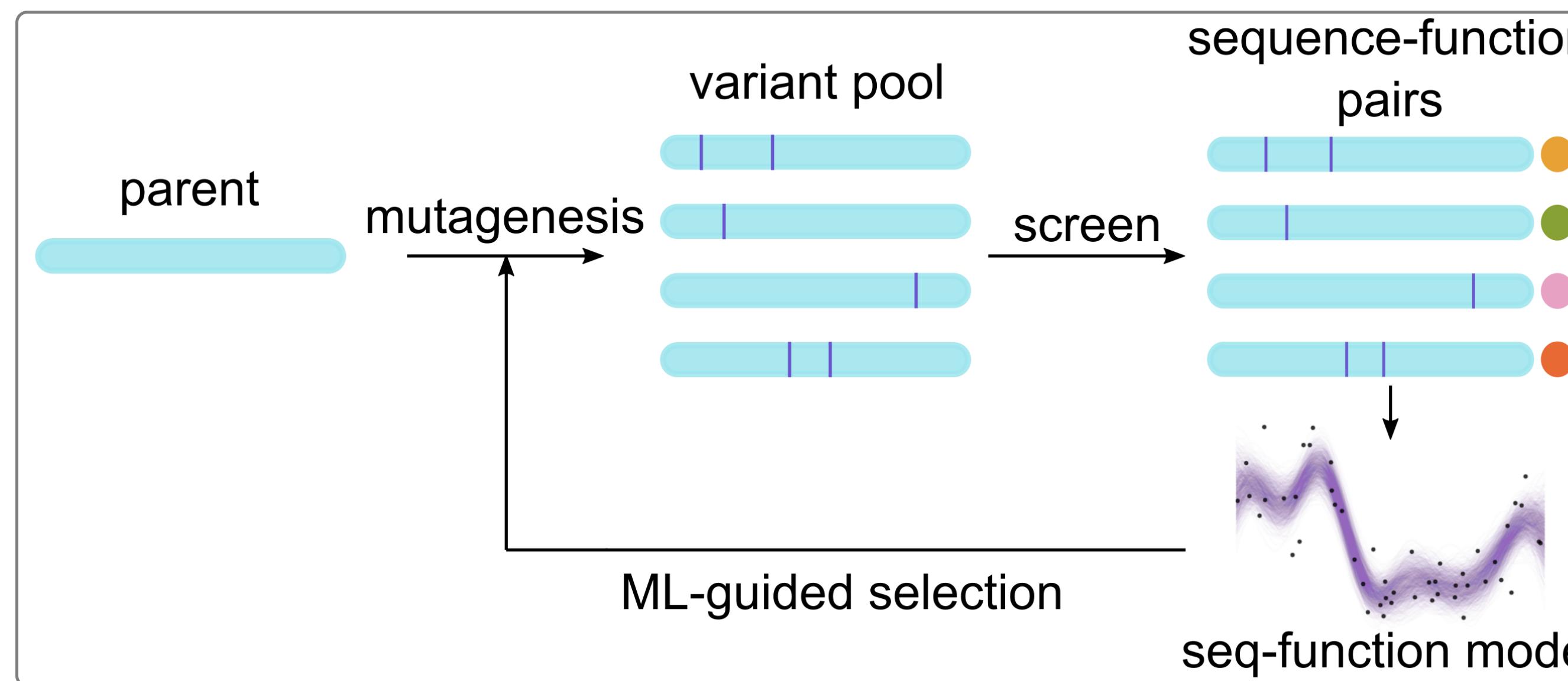


# Machine learning enables optimization with fewer measurements

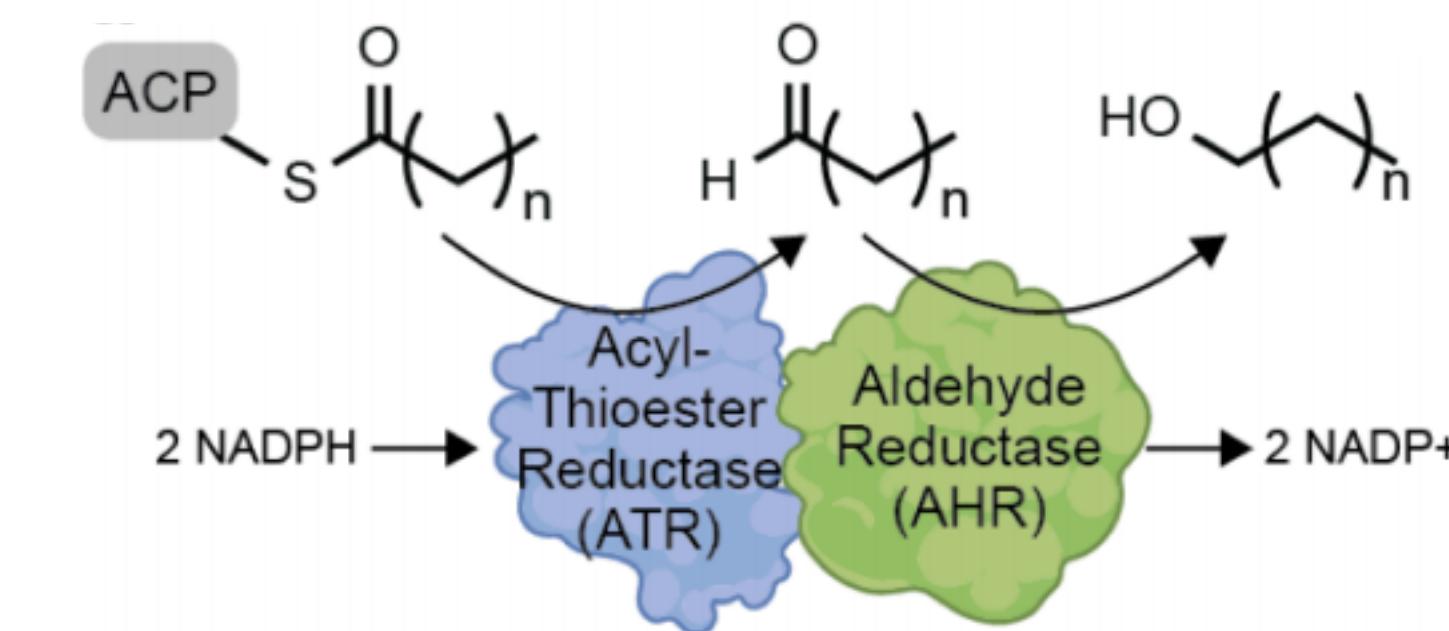


Channelrhodopsin current  
~200 measurements (Bedbrook 2019)

# Machine learning enables optimization with fewer measurements

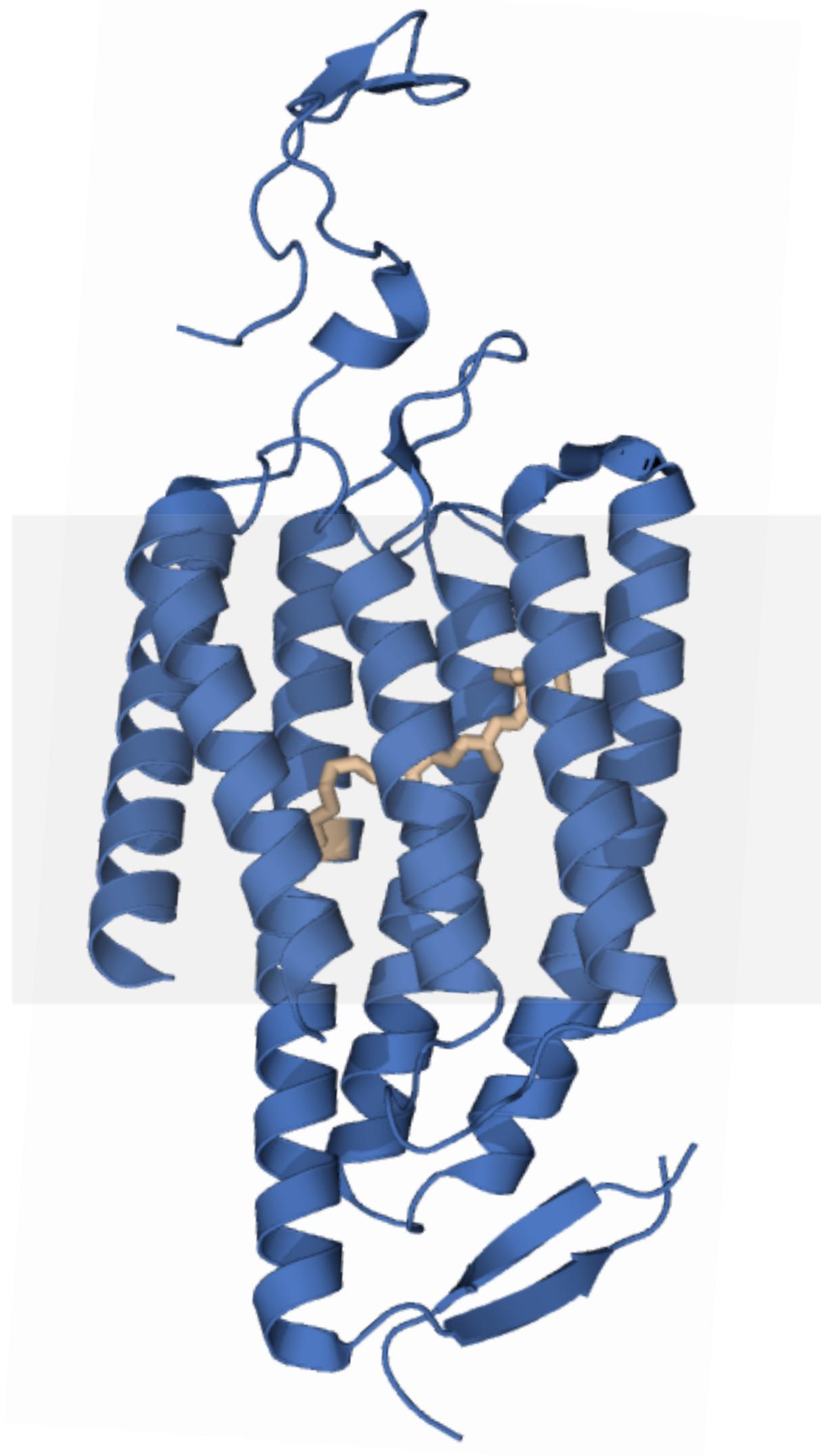


Channelrhodopsin current  
~200 measurements (Bedbrook 2019)

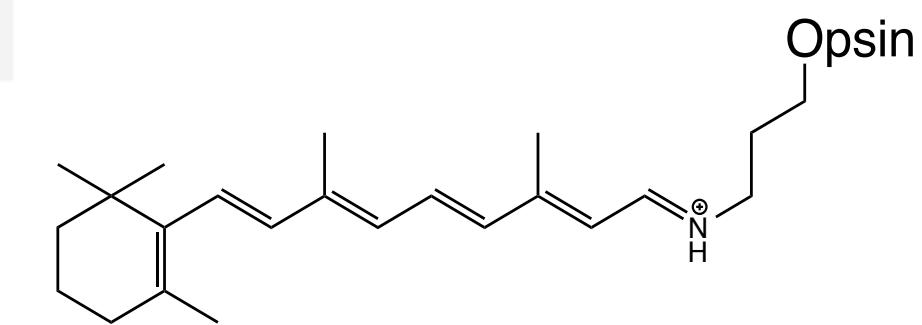
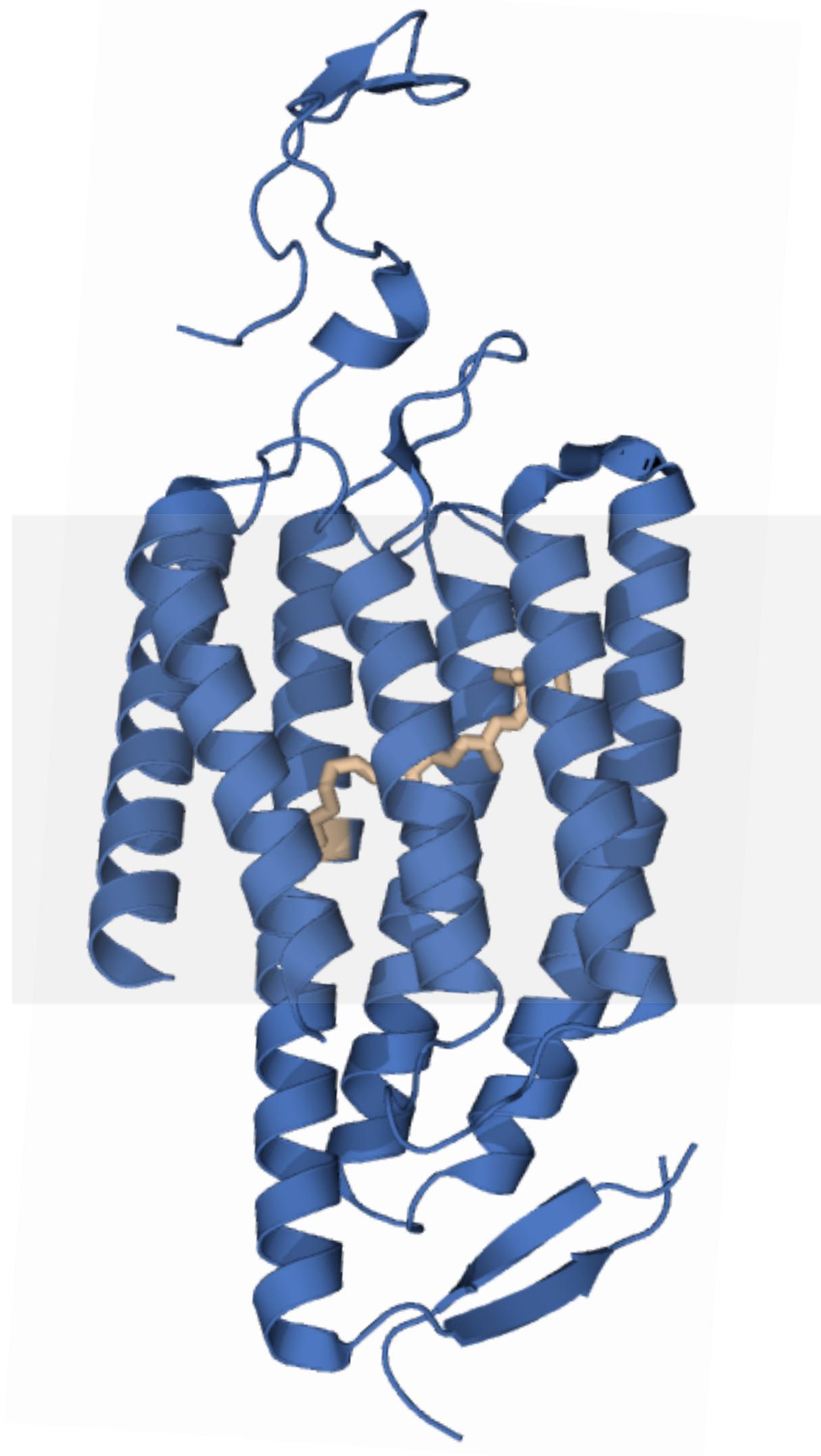


acyl-ACP reductase fatty alcohol synthesis  
~100 measurements (Greenhalgh 2021)

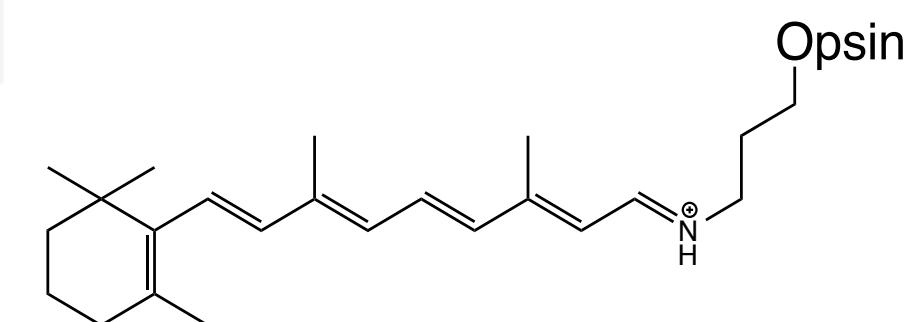
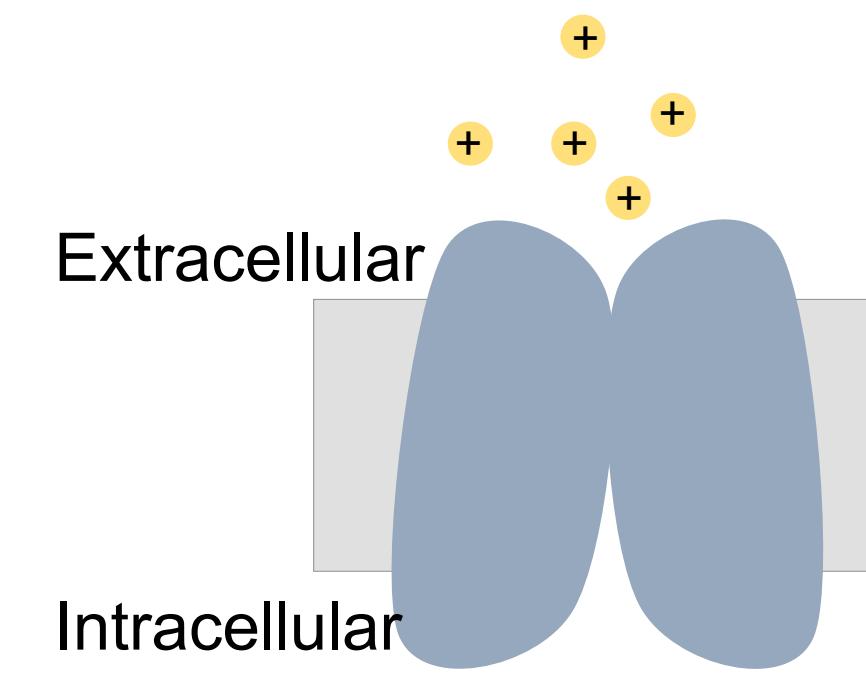
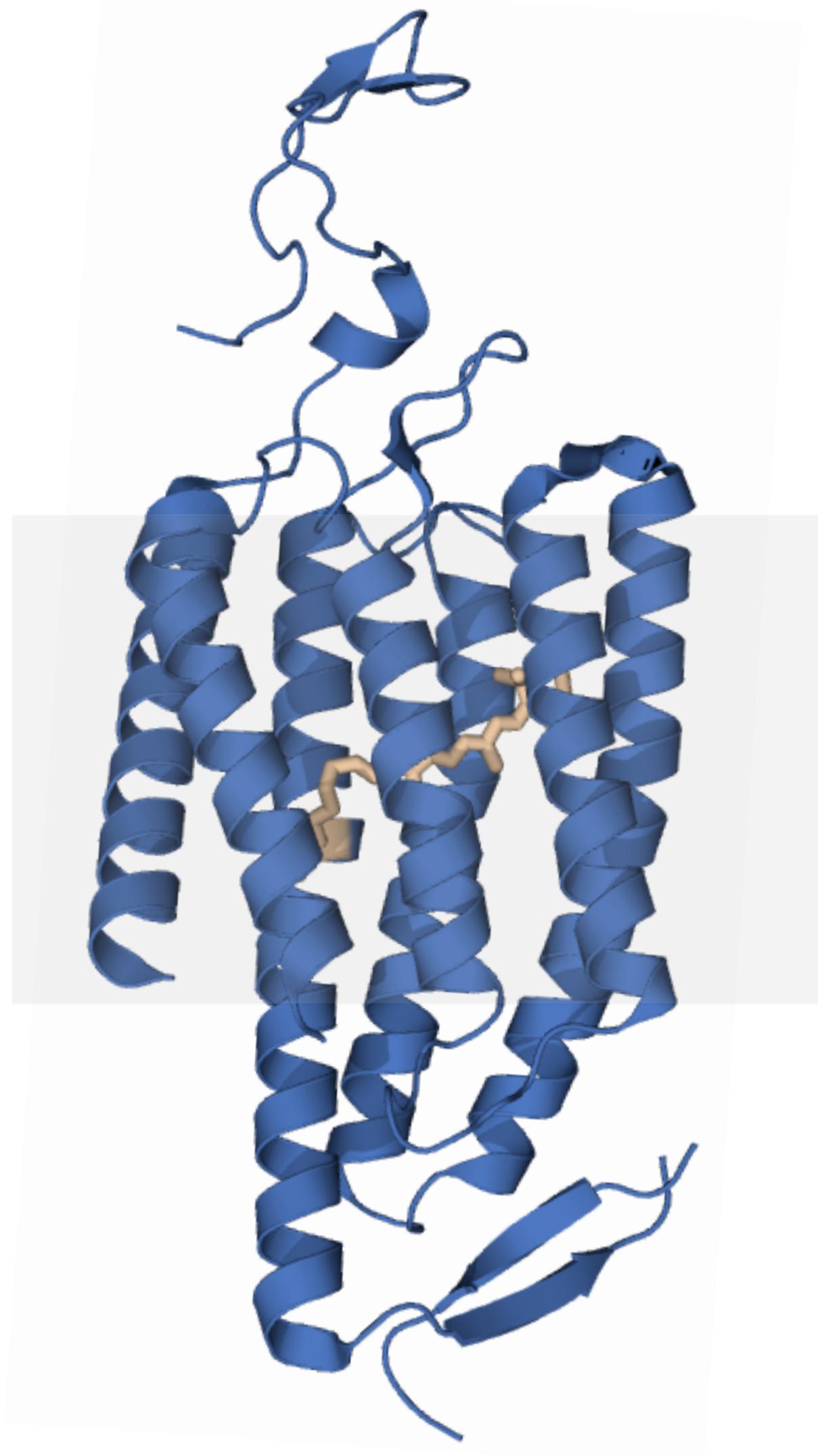
# Channelrhodopsins: light-gated ion channels



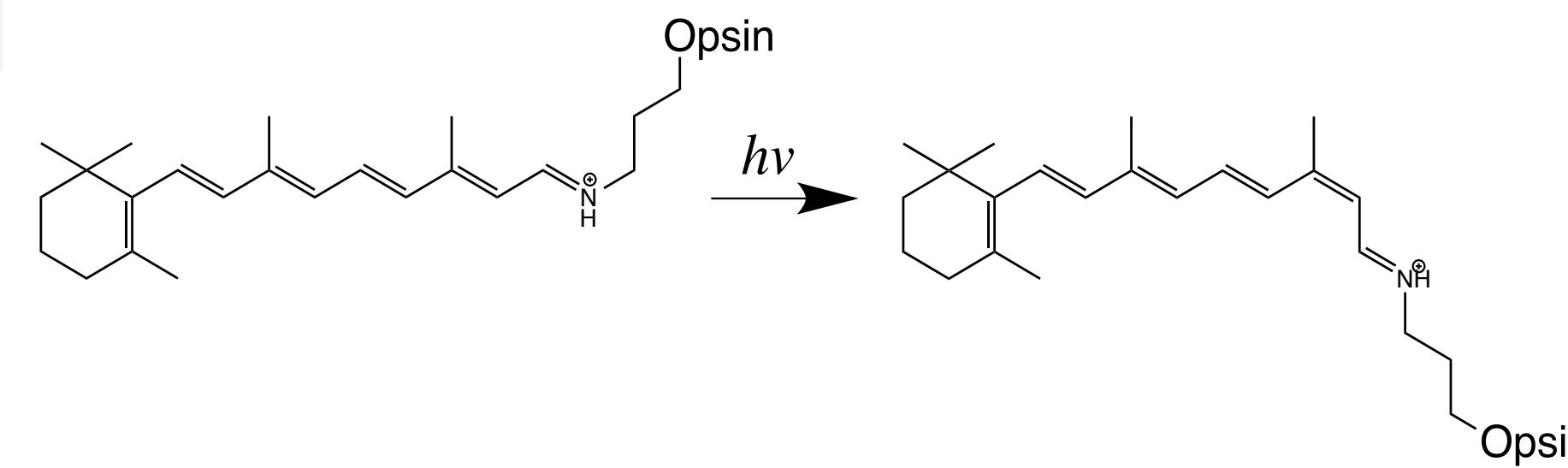
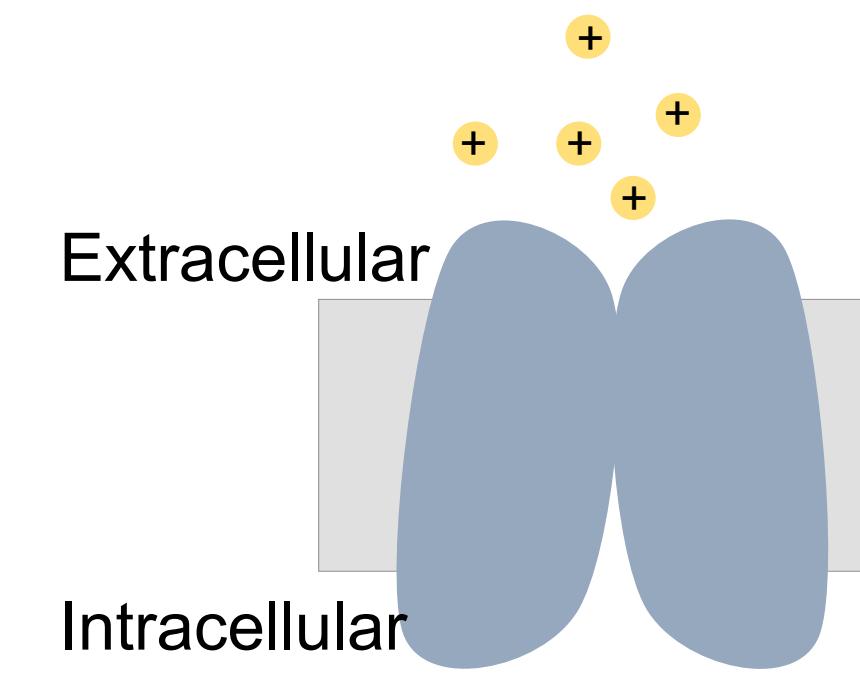
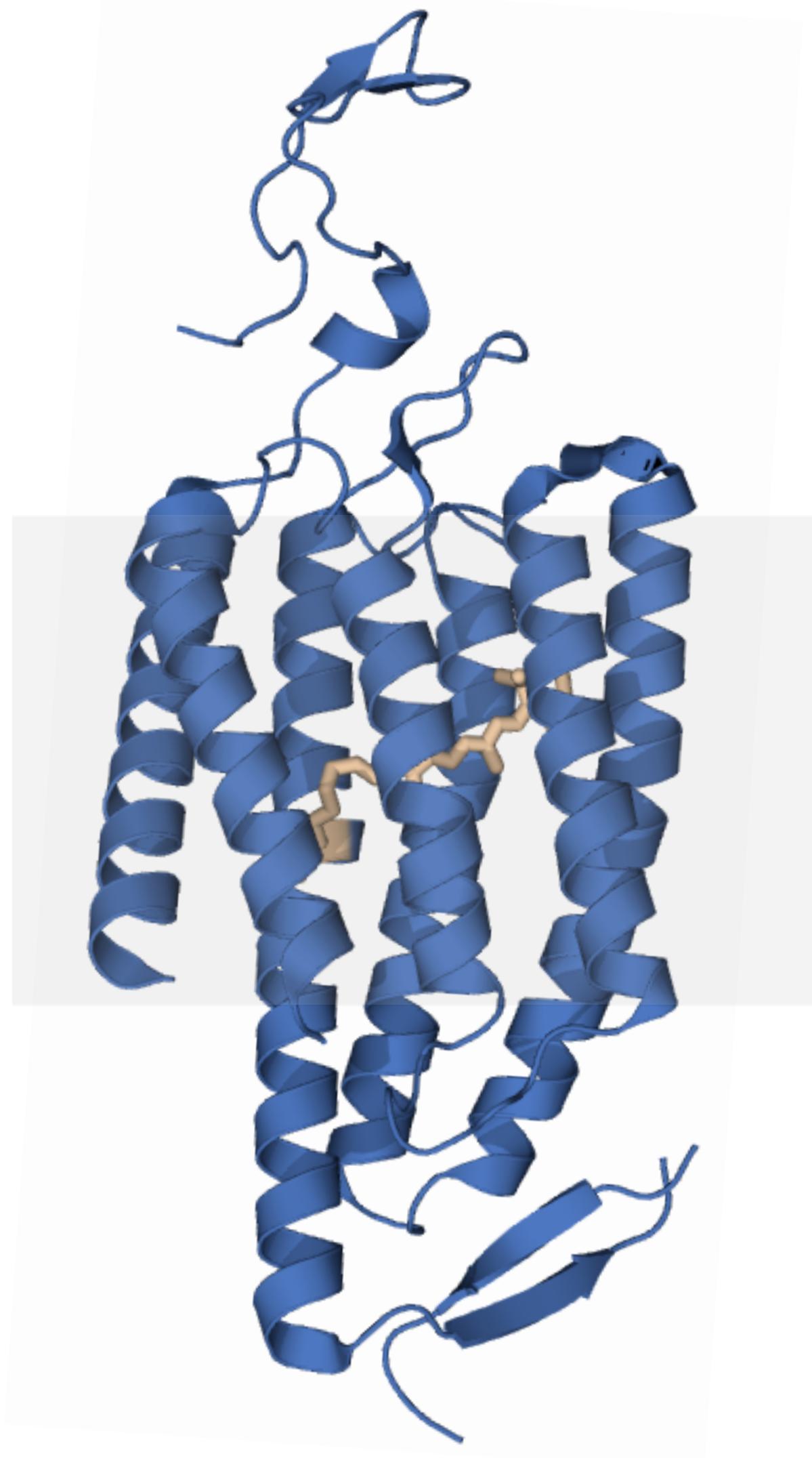
# Channelrhodopsins: light-gated ion channels



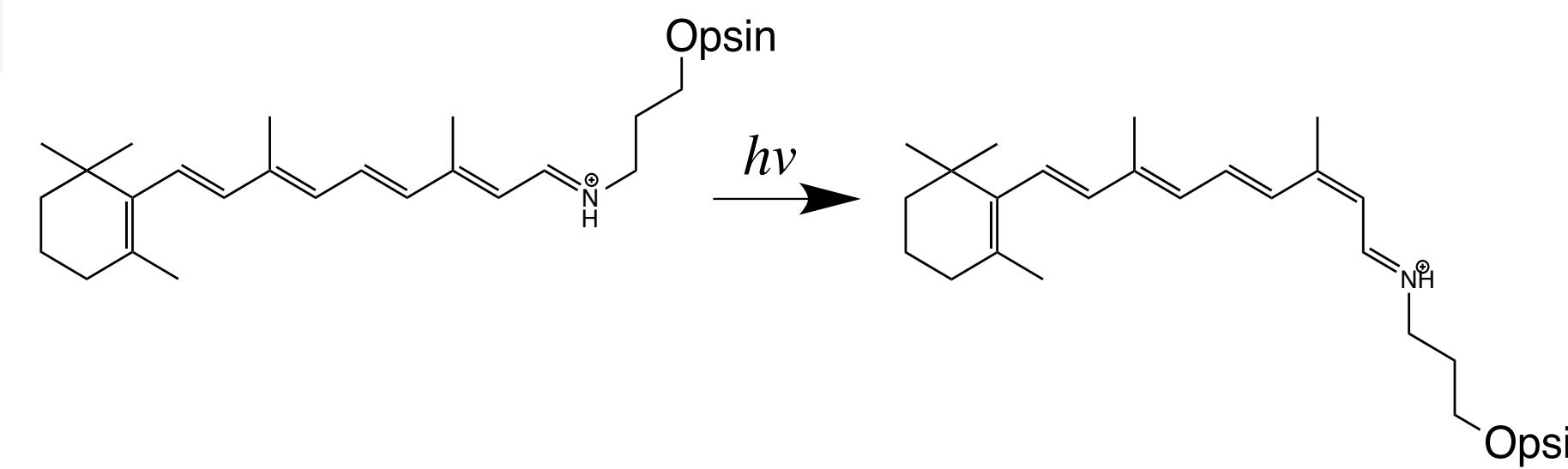
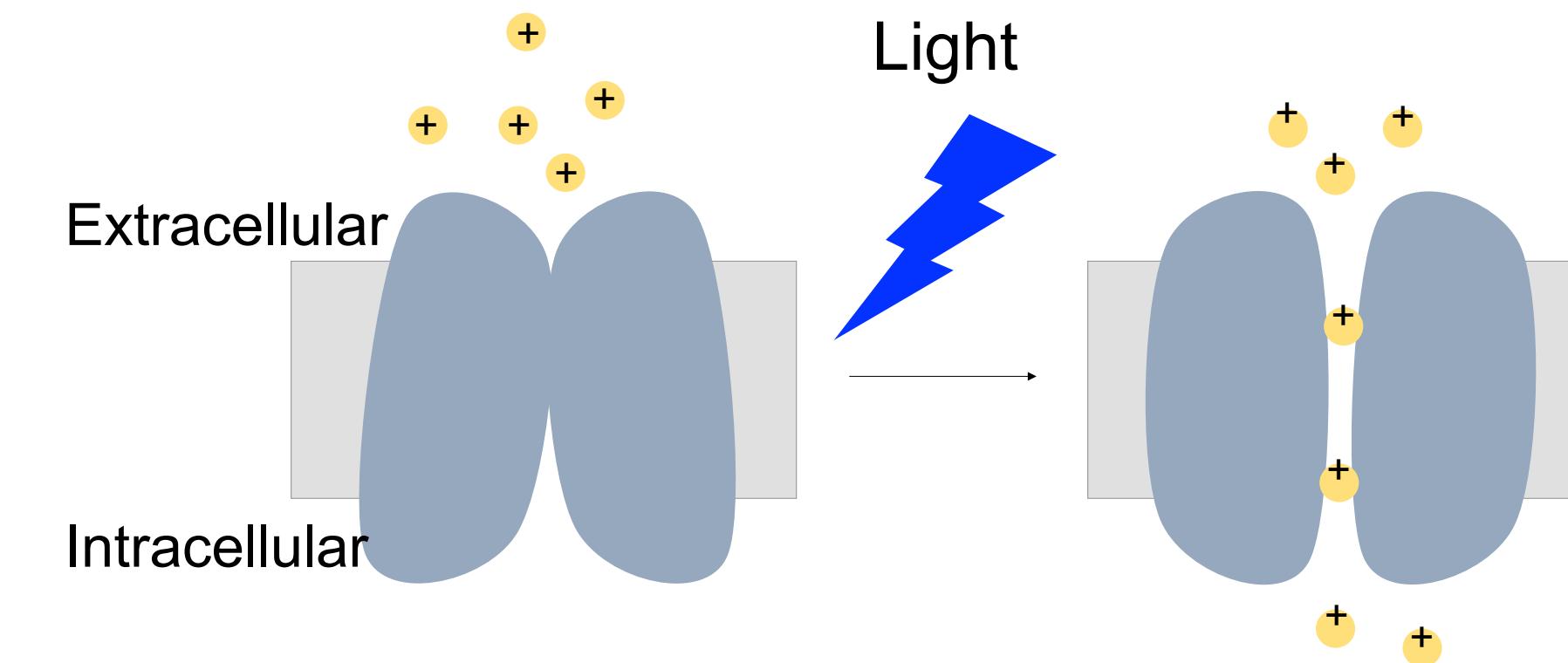
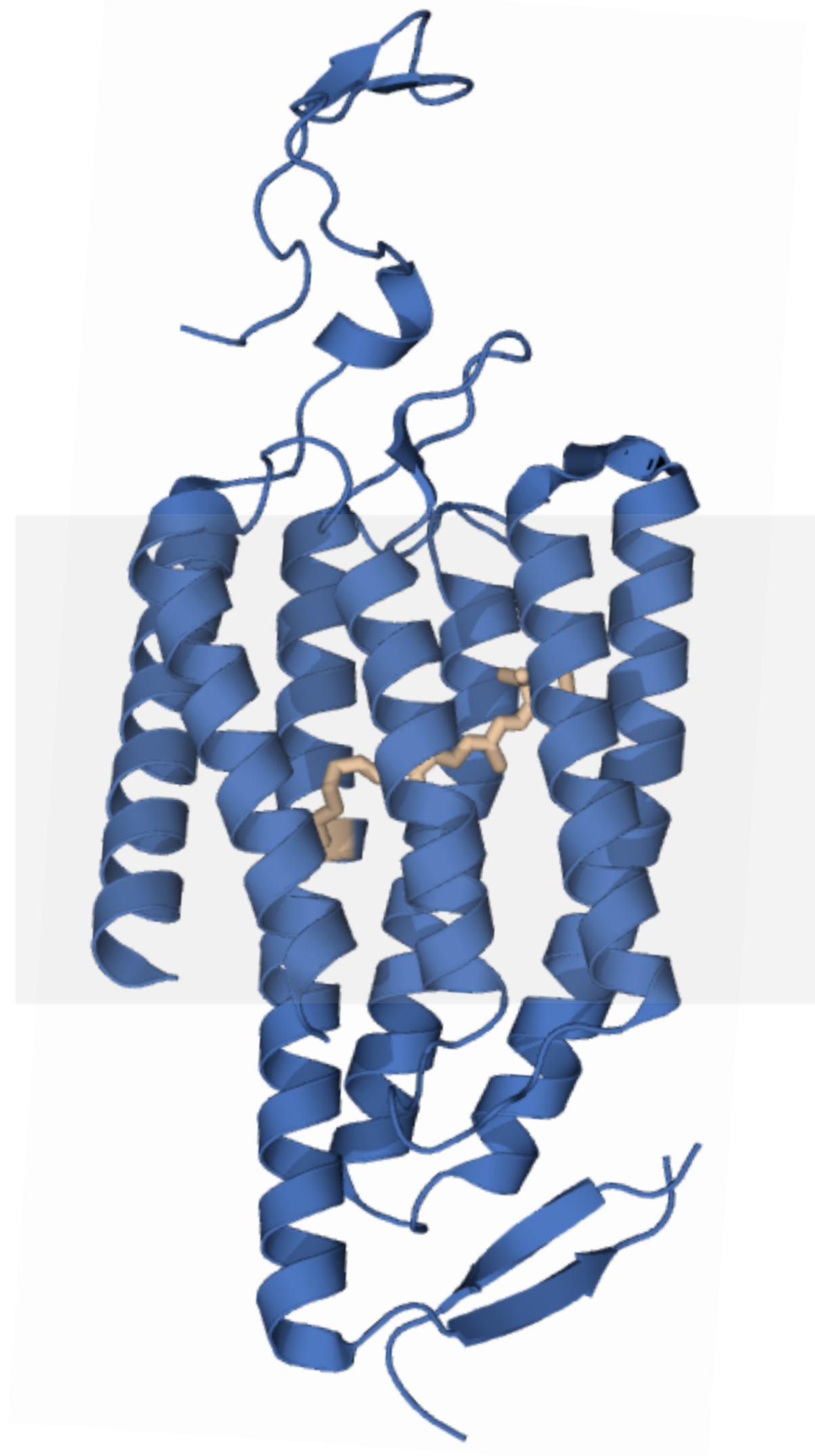
# Channelrhodopsins: light-gated ion channels



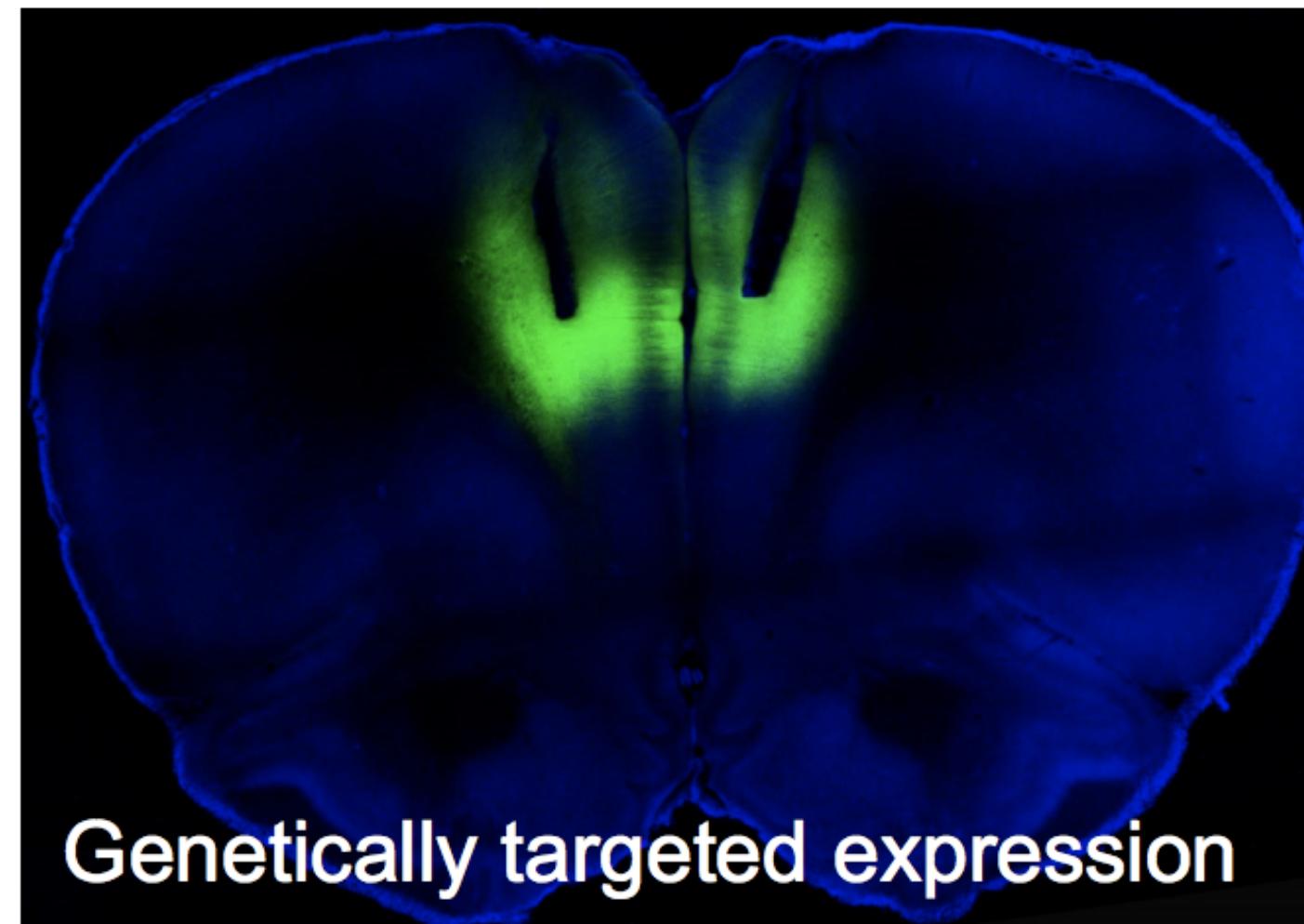
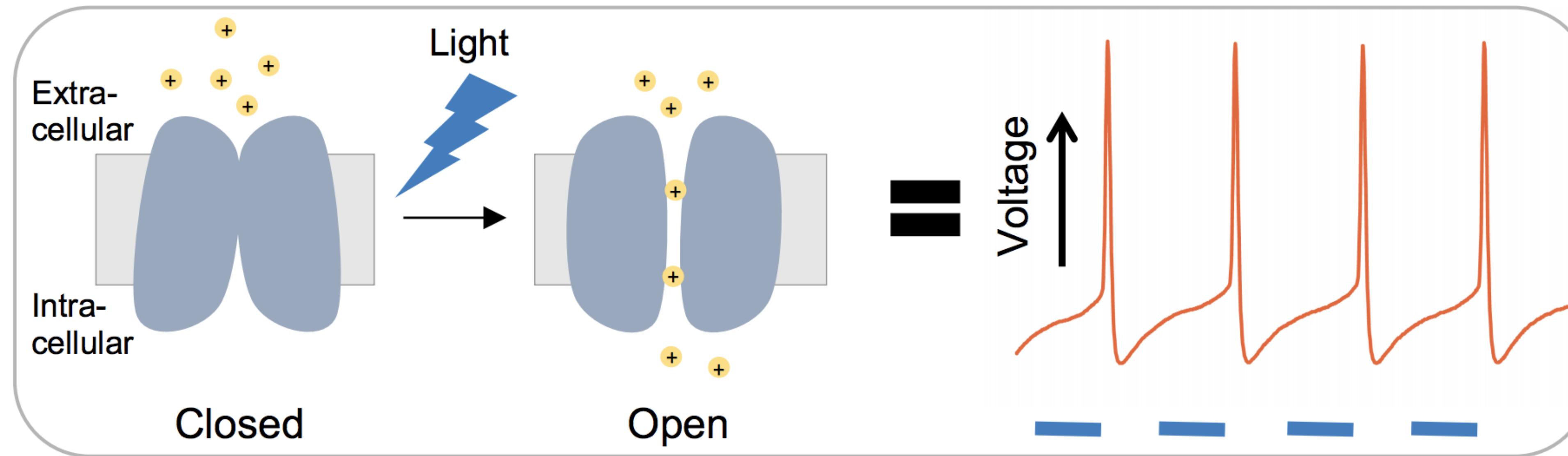
# Channelrhodopsins: light-gated ion channels



# Channelrhodopsins: light-gated ion channels



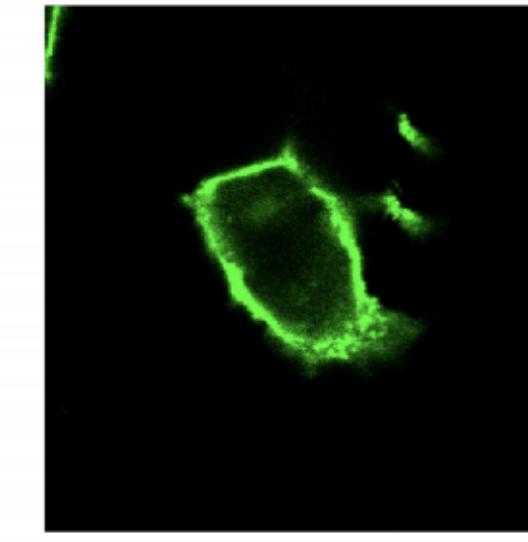
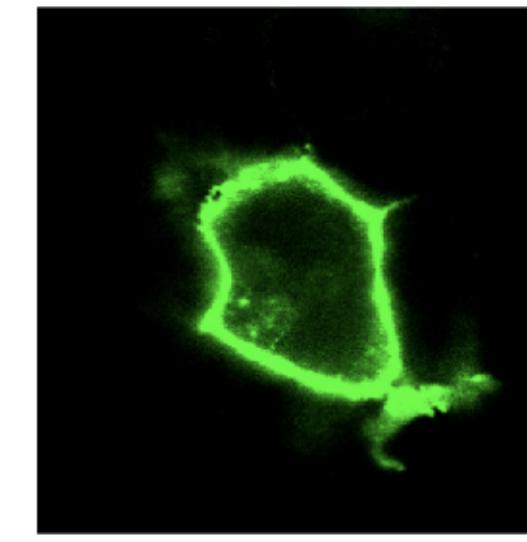
# ChRs are optogenetic tools



# Multiple engineering goals

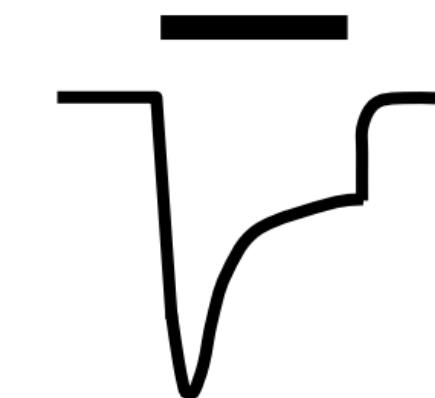
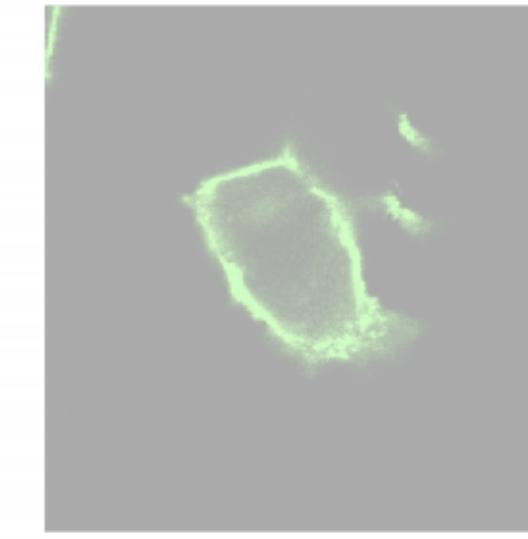
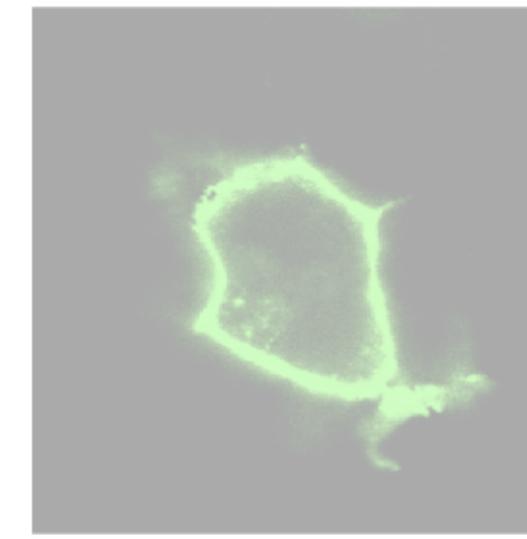
# Multiple engineering goals

- Heterologous membrane localization



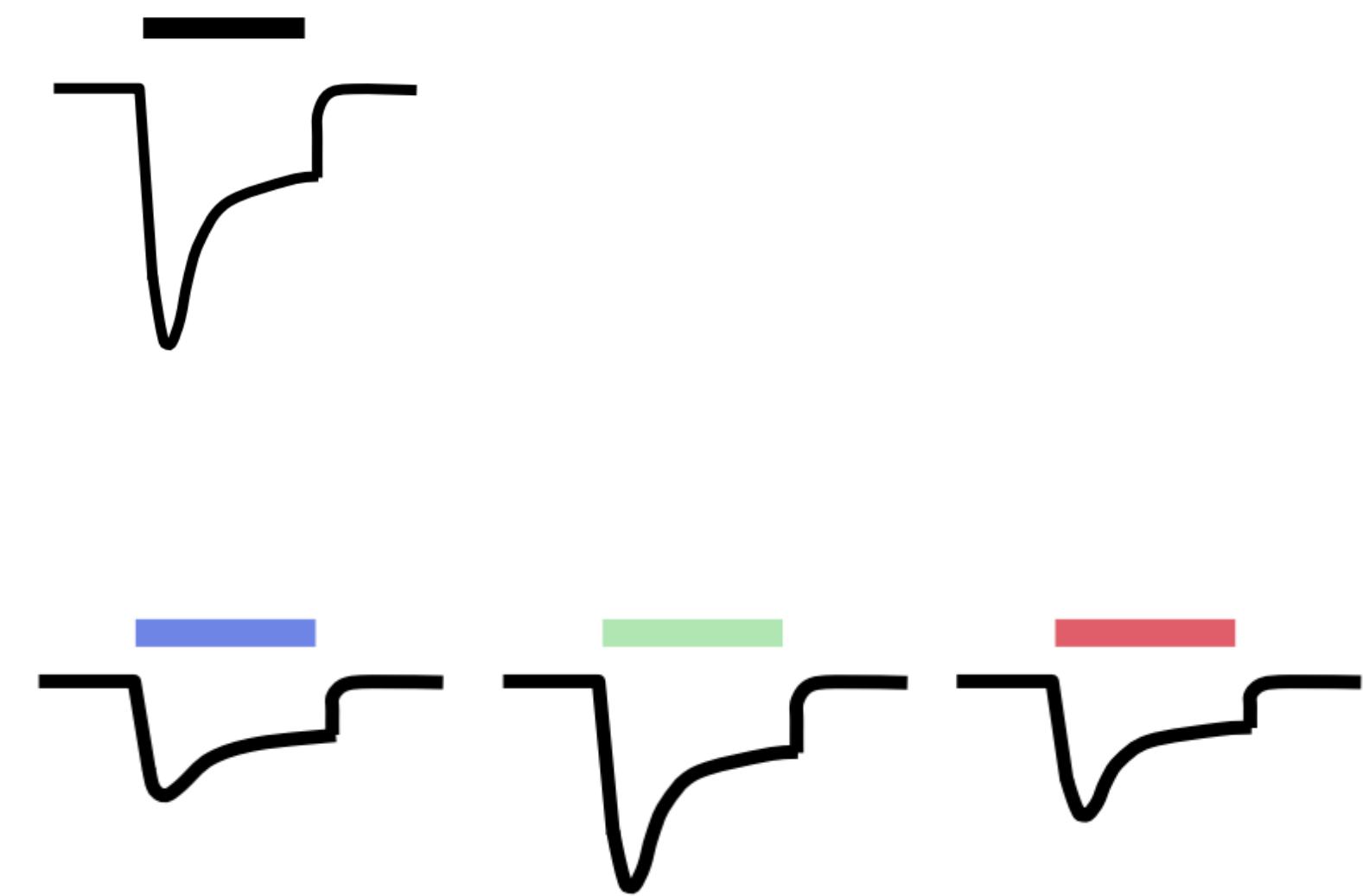
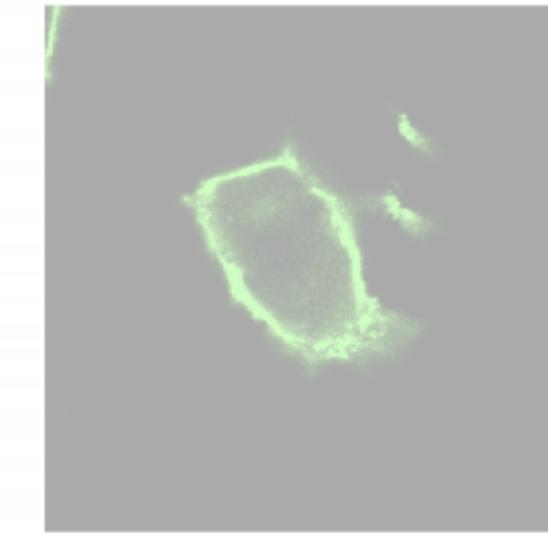
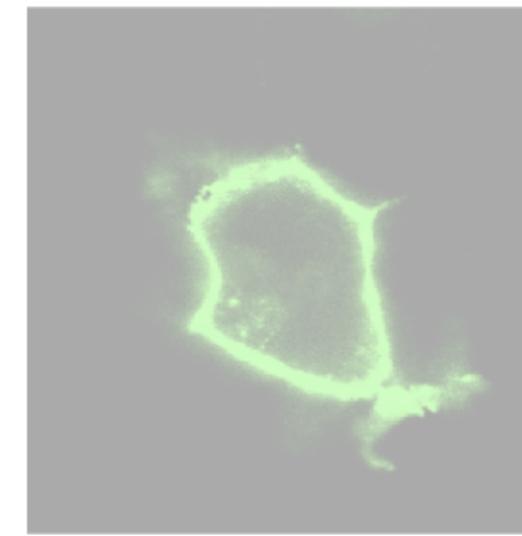
# Multiple engineering goals

- Heterologous membrane localization
- Increased sensitivity



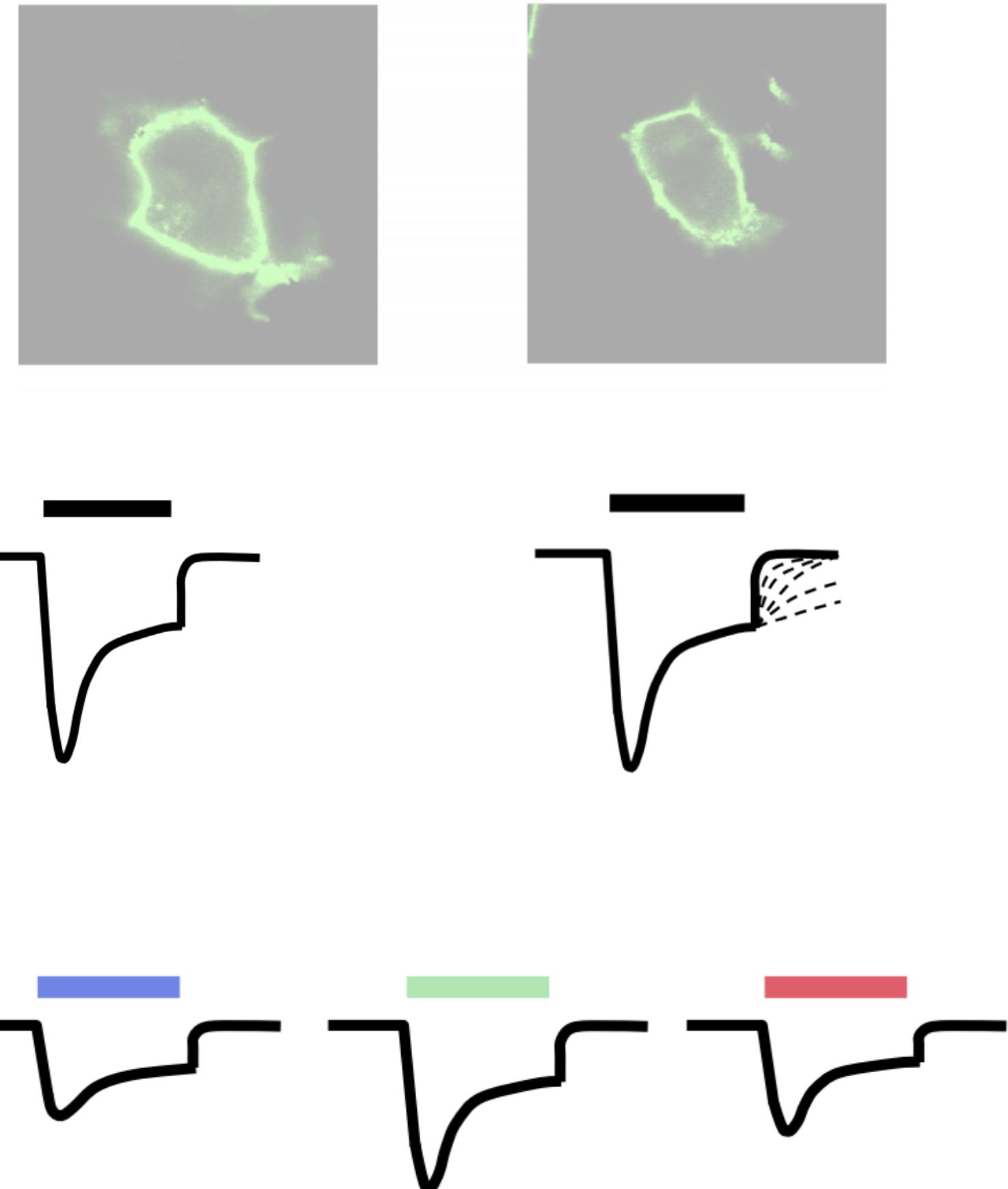
# Multiple engineering goals

- Heterologous membrane localization
- Increased sensitivity
- Different activation wavelengths

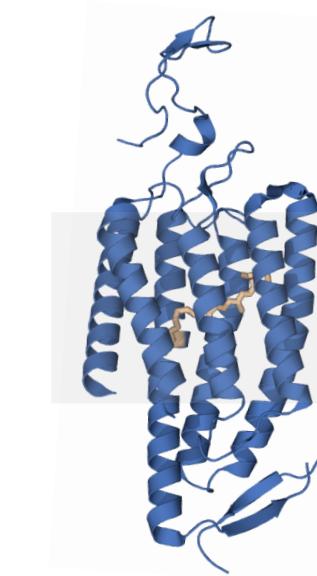
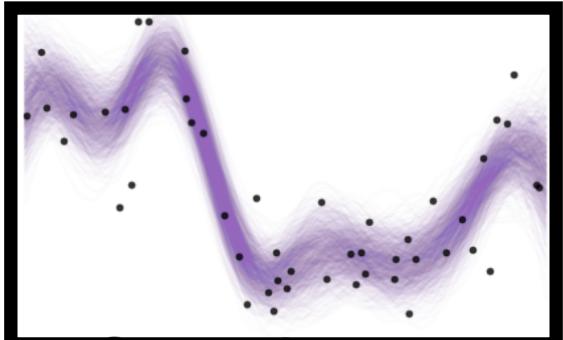


# Multiple engineering goals

- Heterologous membrane localization
- Increased sensitivity
- Different activation wavelengths
- Different on/off kinetics

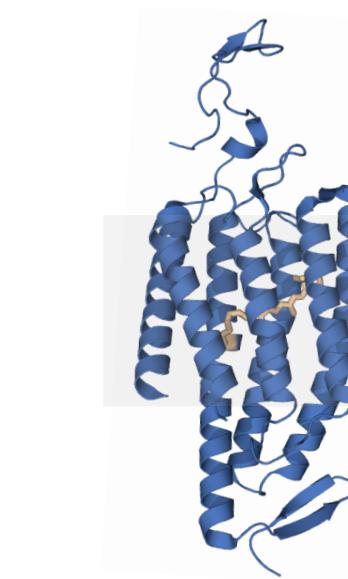
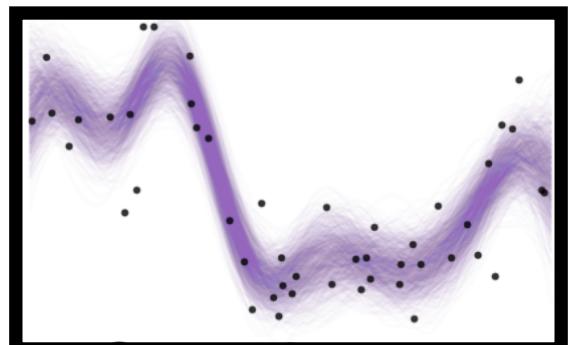
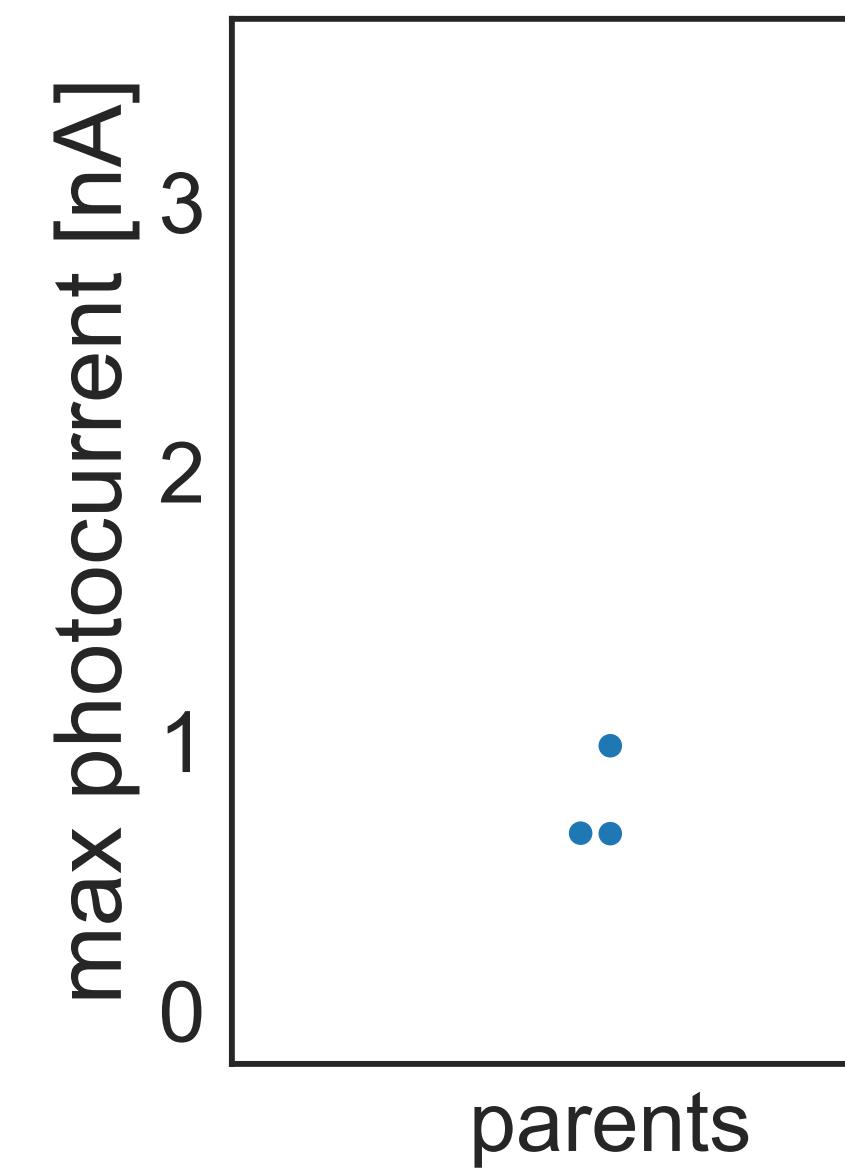


# Machine learning enables optimization with fewer measurements



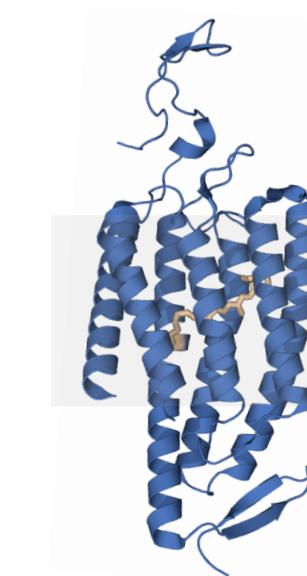
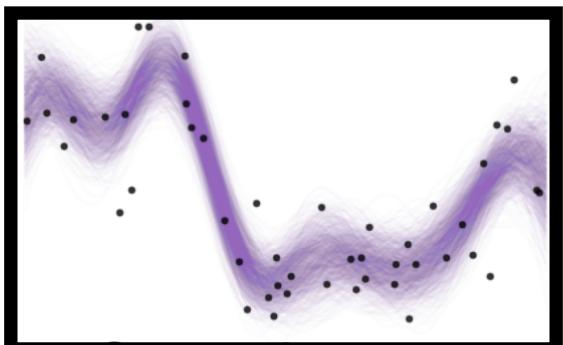
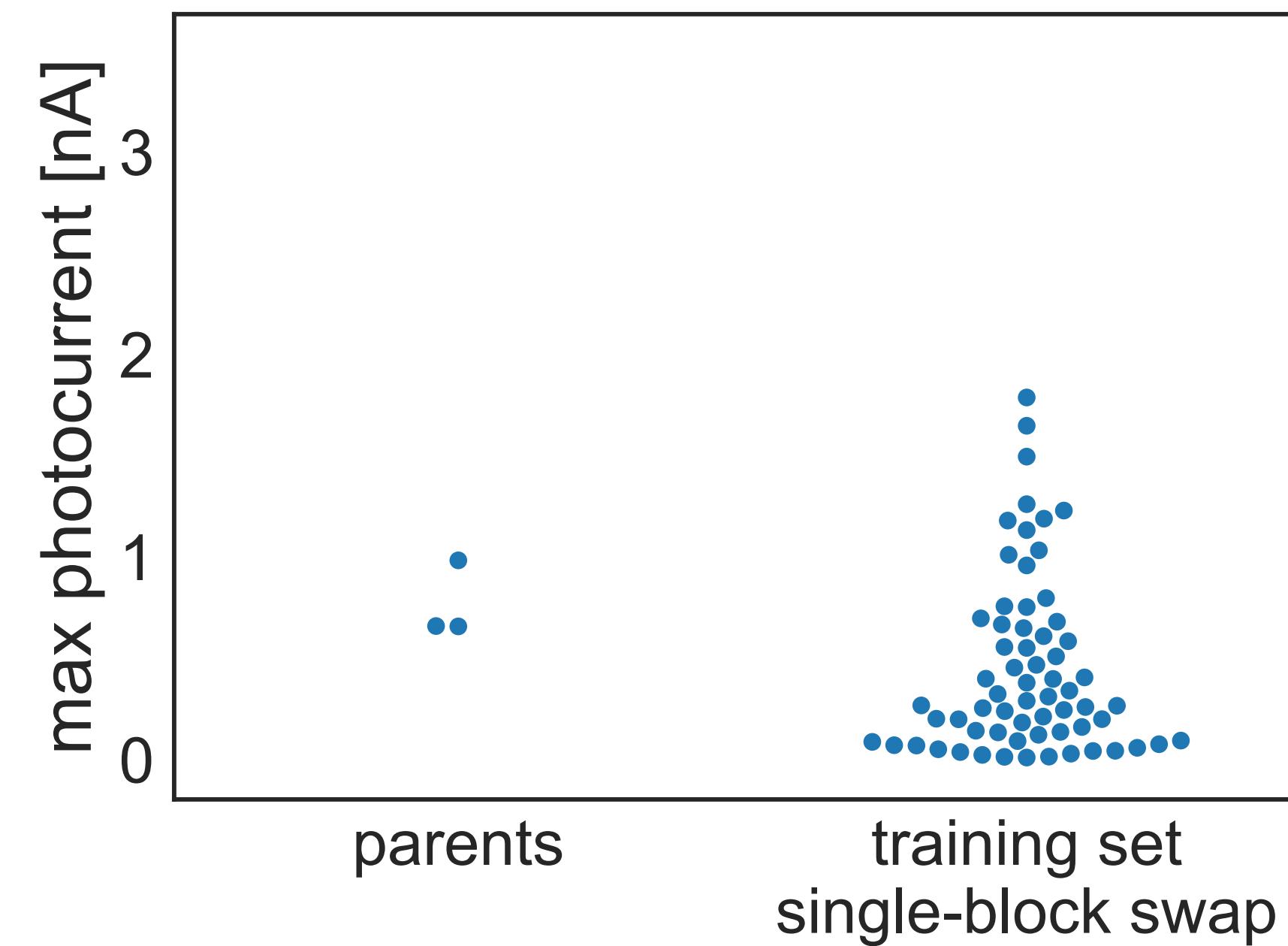
Channelrhodopsin current  
~200 measurements (Bedbrook 2019)

# Machine learning enables optimization with fewer measurements



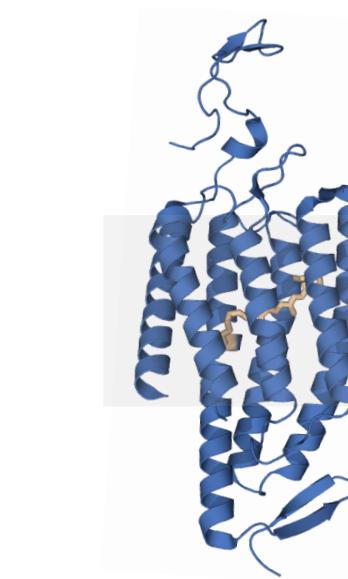
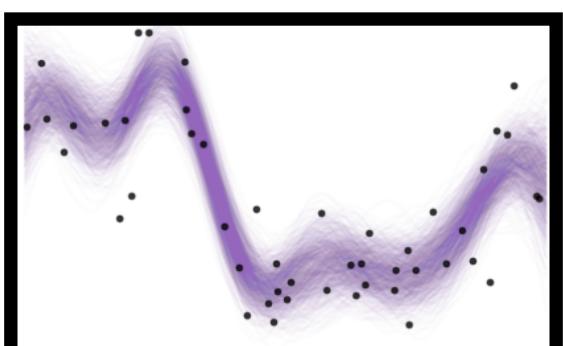
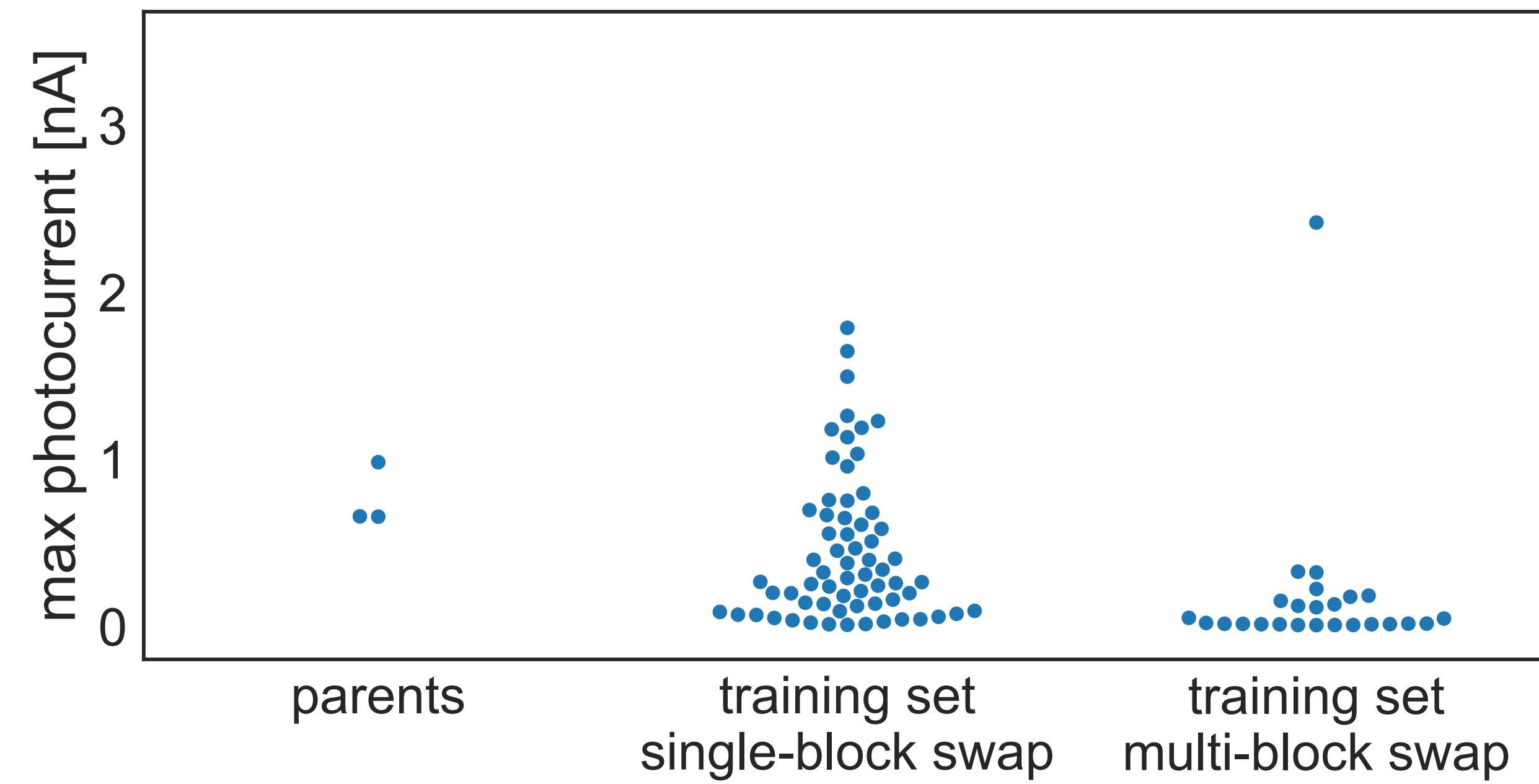
Channelrhodopsin current  
~200 measurements (Bedbrook 2019)

# Machine learning enables optimization with fewer measurements



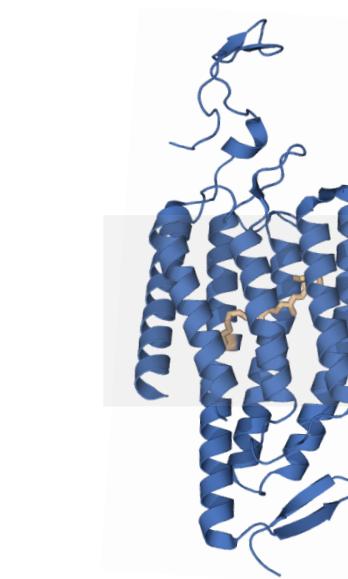
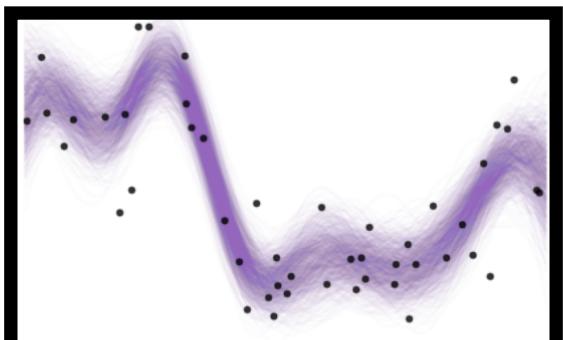
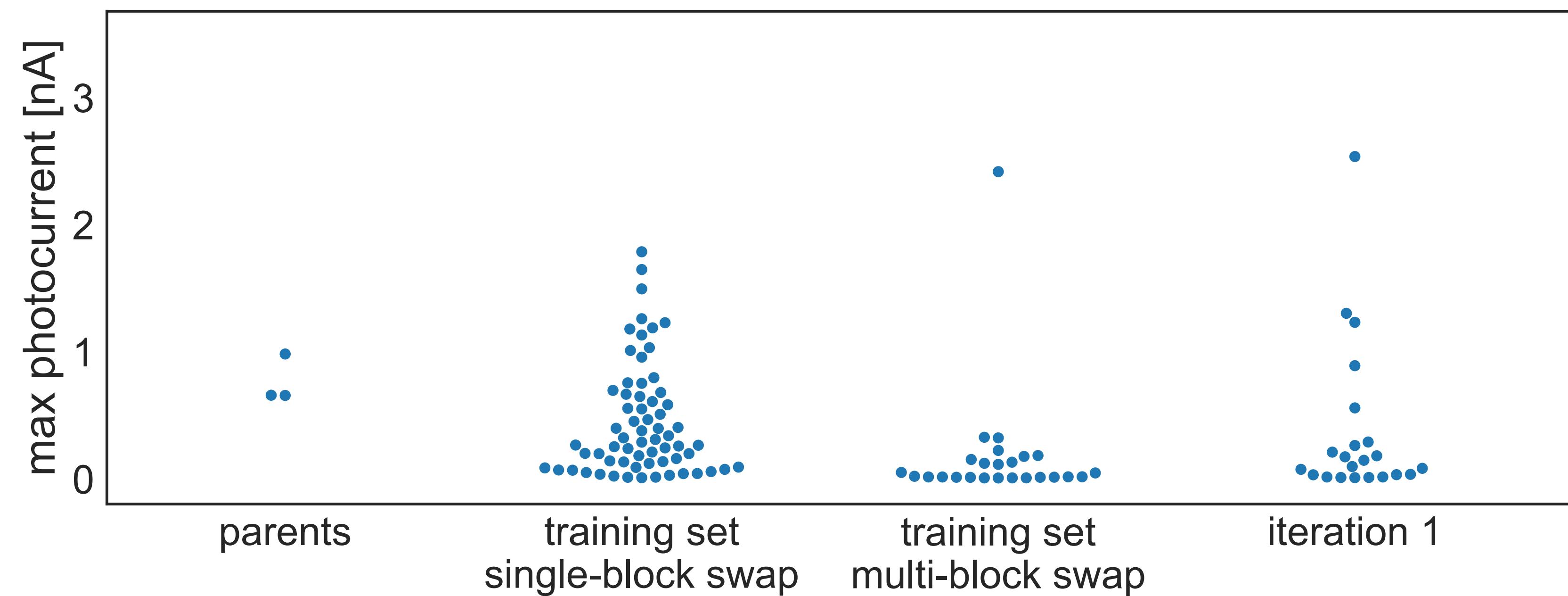
Channelrhodopsin current  
~200 measurements (Bedbrook 2019)

# Machine learning enables optimization with fewer measurements



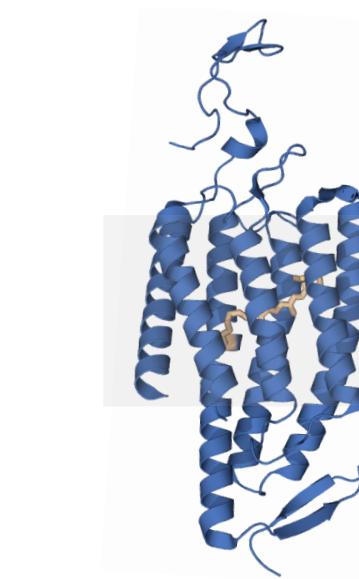
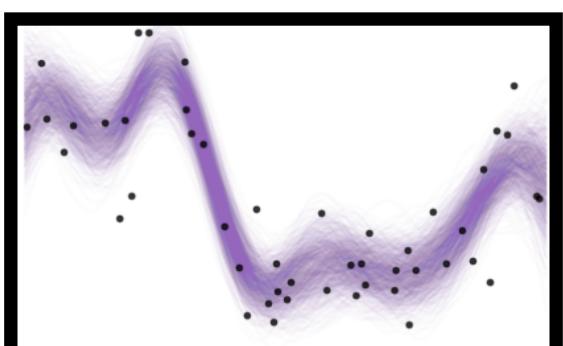
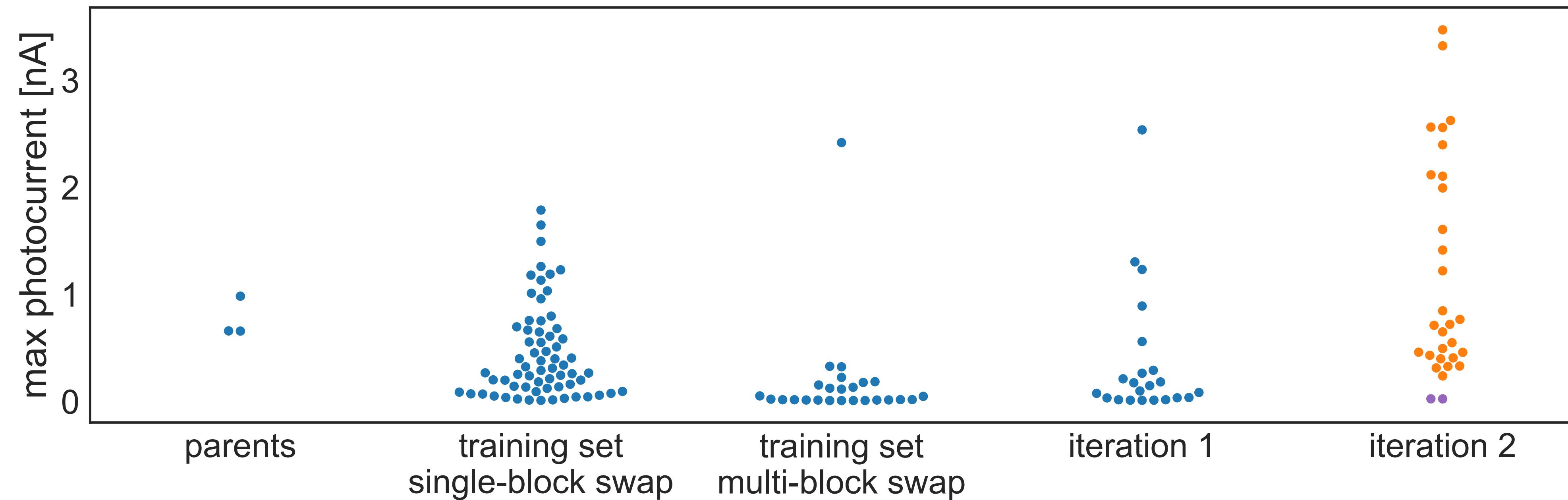
Channelrhodopsin current  
~200 measurements (Bedbrook 2019)

# Machine learning enables optimization with fewer measurements



Channelrhodopsin current  
~200 measurements (Bedbrook 2019)

# Machine learning enables optimization with fewer measurements



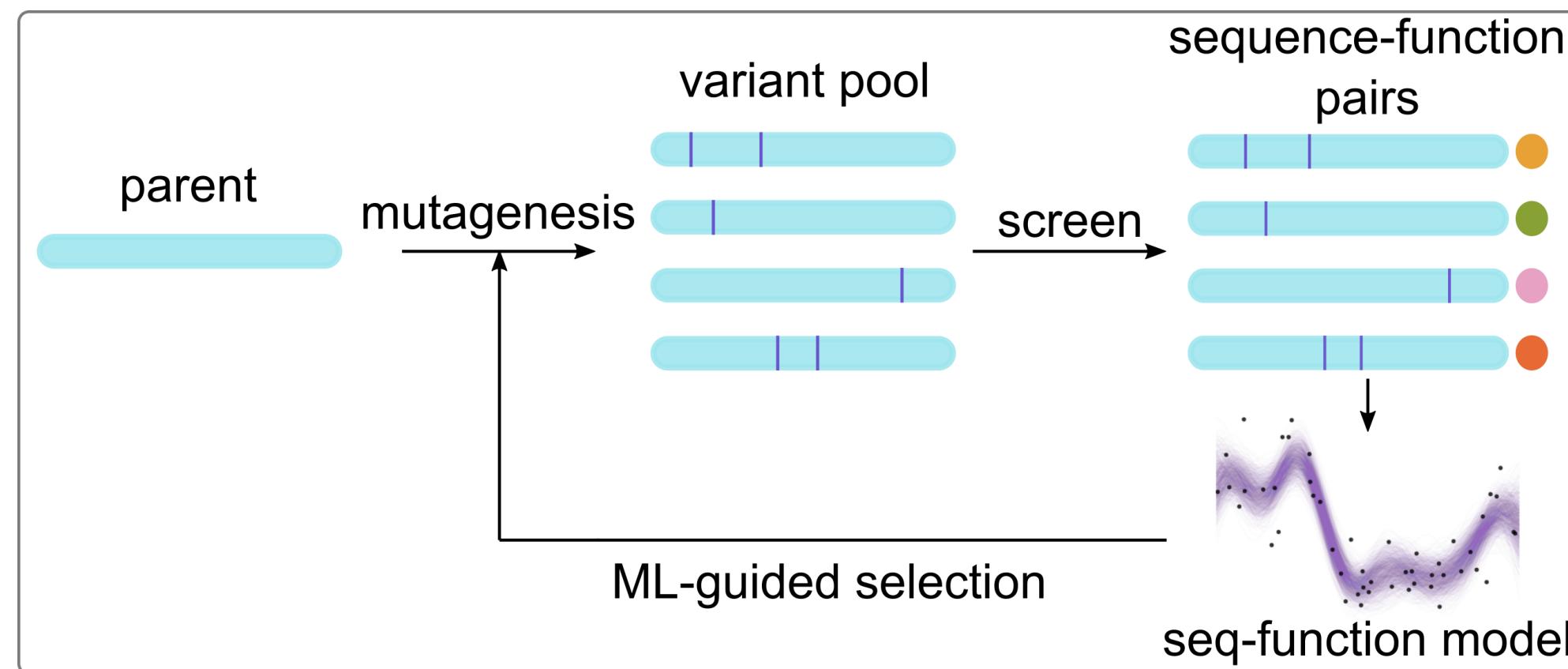
Channelrhodopsin current  
~200 measurements (Bedbrook 2019)

# Engineered ChRs control mouse neurons without skull removal

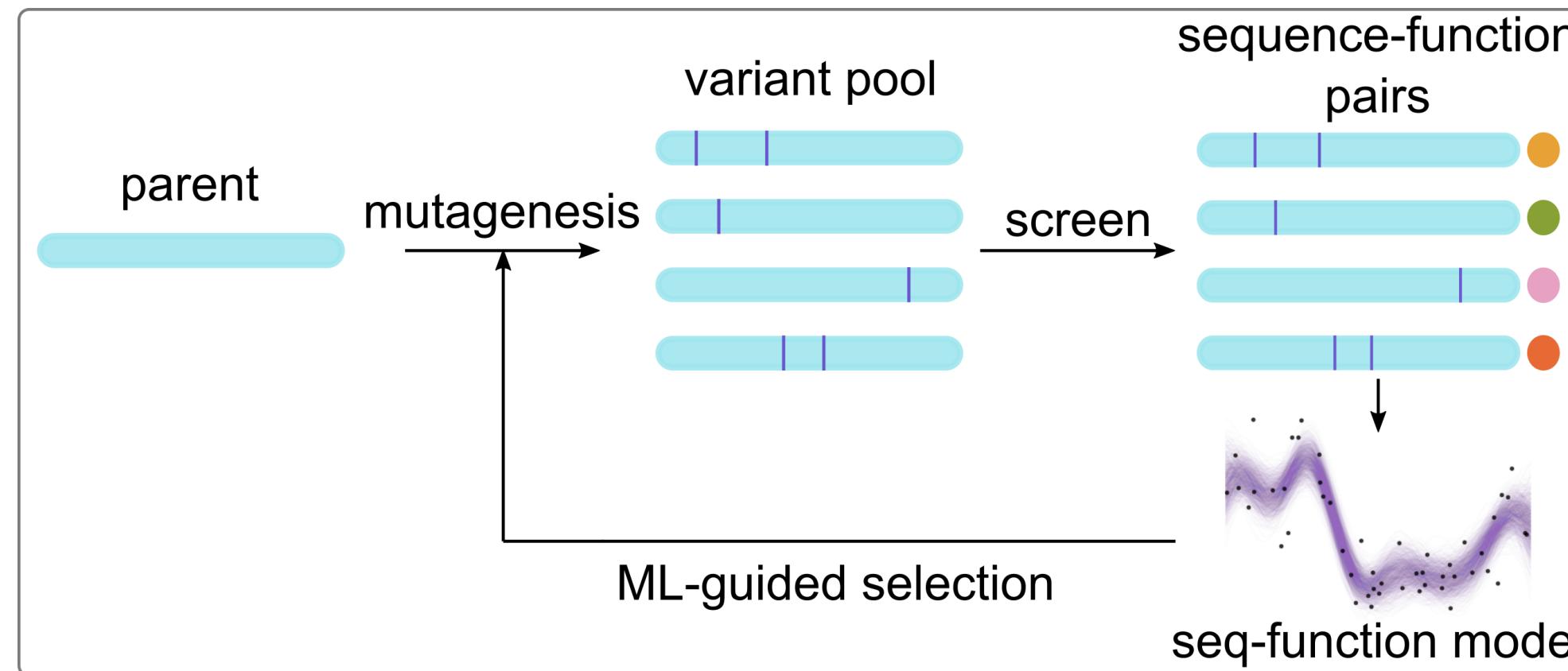
# Engineered ChRs control mouse neurons without skull removal



# MLDE: successes and limitations

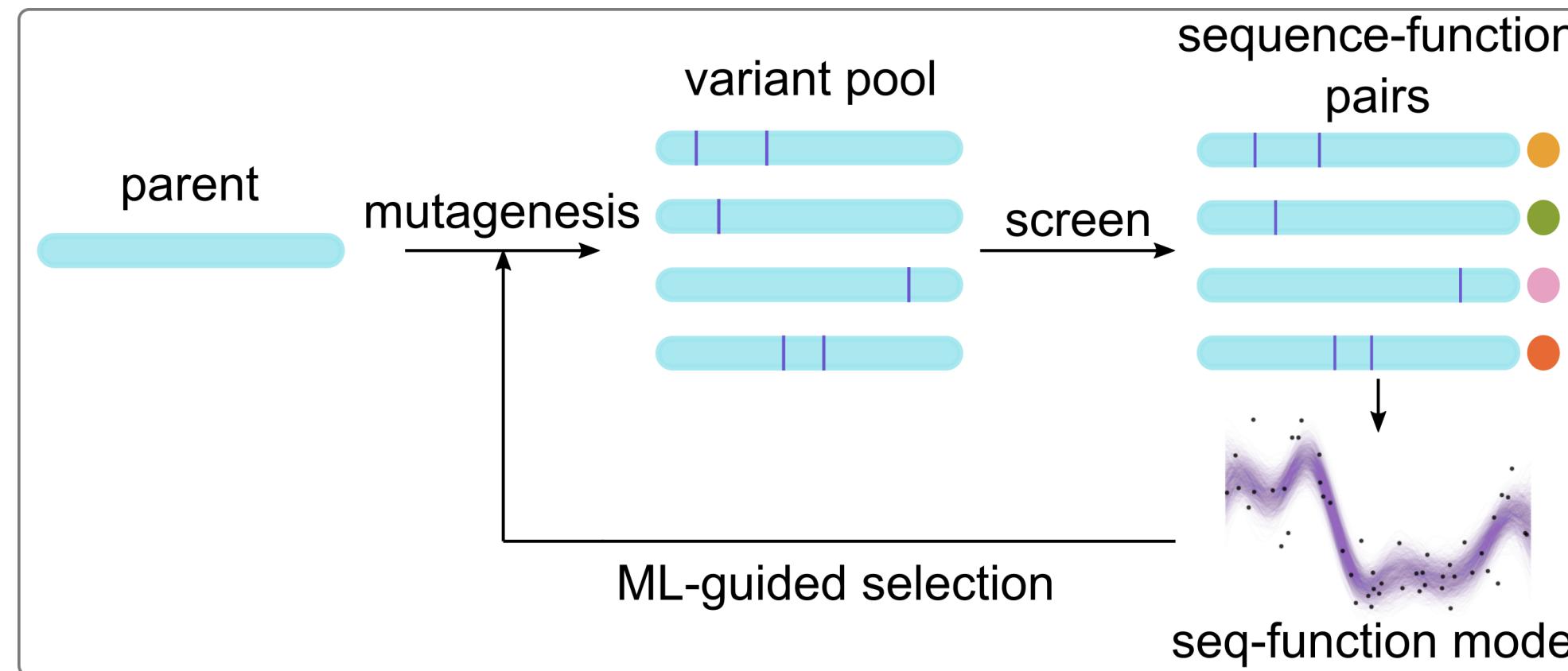


# MLDE: successes and limitations



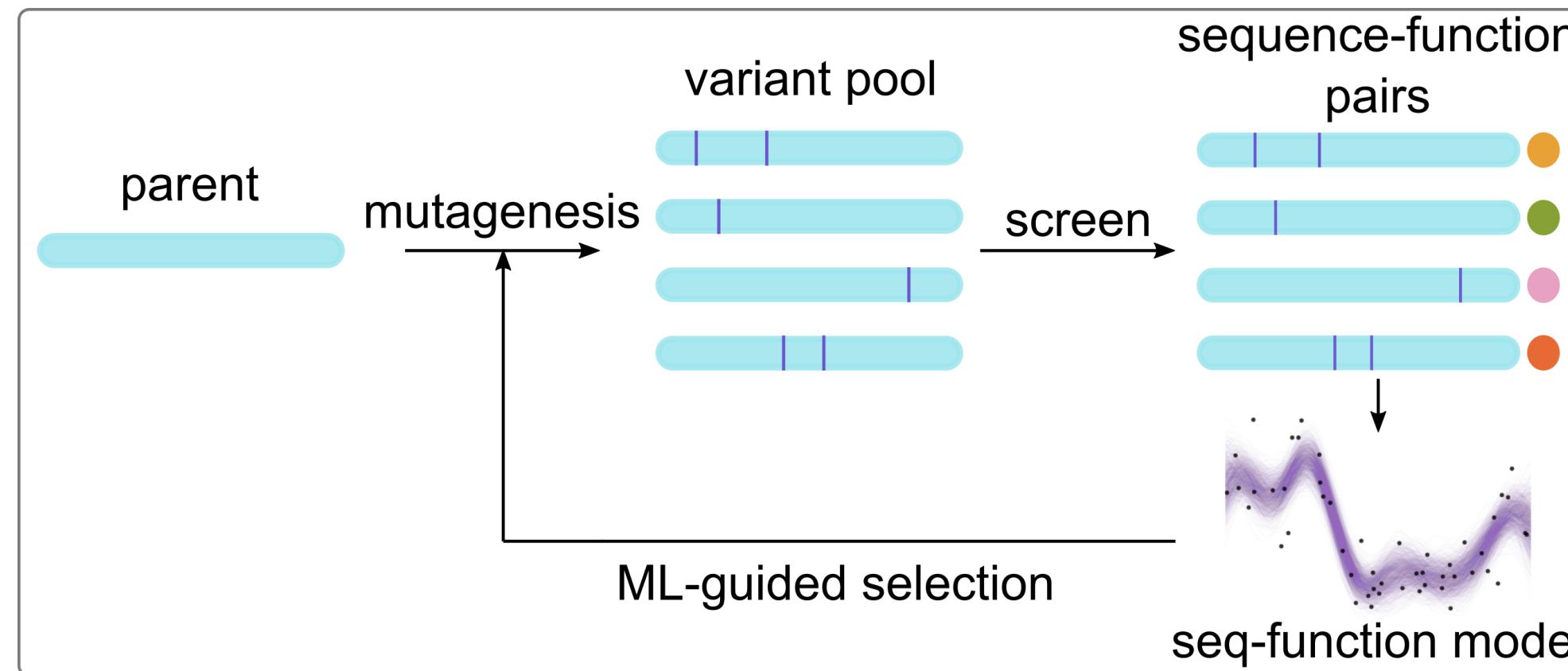
- Accurate models with < 200 measurements

# MLDE: successes and limitations



- Accurate models with < 200 measurements
- Find rare, improved sequences with < 200 measurements

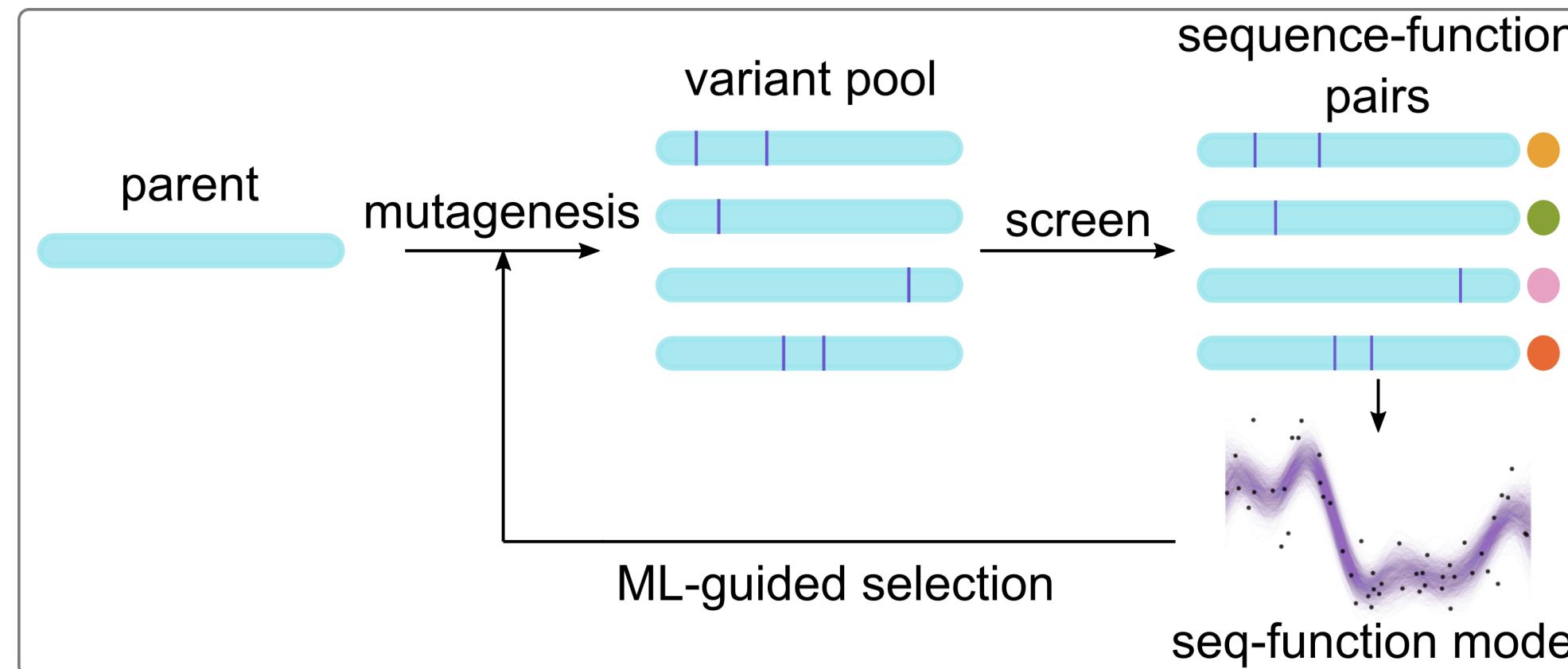
# MLDE: successes and limitations



- Accurate models with < 200 measurements
- Find rare, improved sequences with < 200 measurements

Still requires starting point!

# MLDE: successes and limitations



- Accurate models with < 200 measurements
- Find rare, improved sequences with < 200 measurements

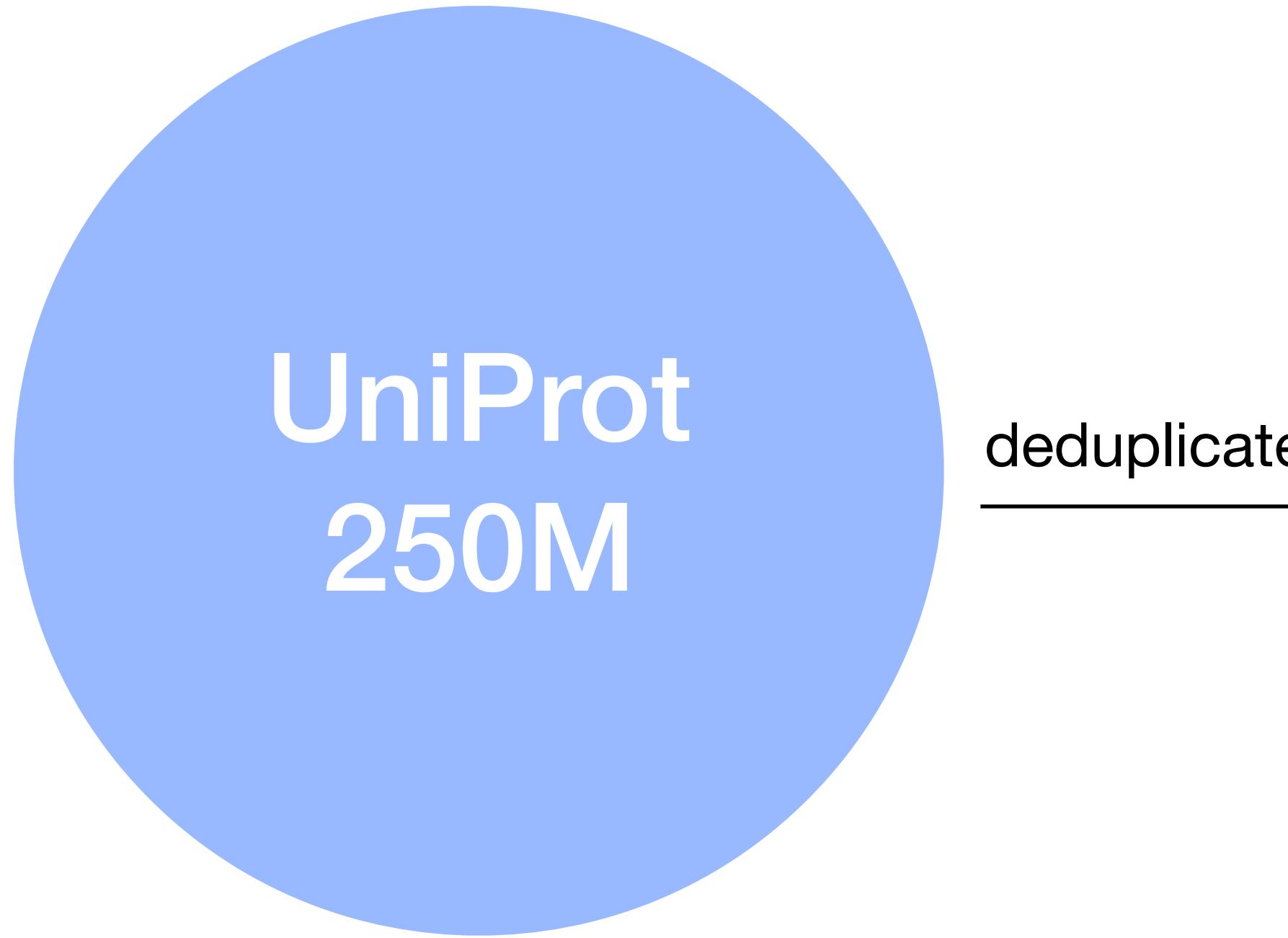
Still requires starting point!  
Ignores mountains of protein data

**Pretraining can leverage large sequence and  
structure databases**

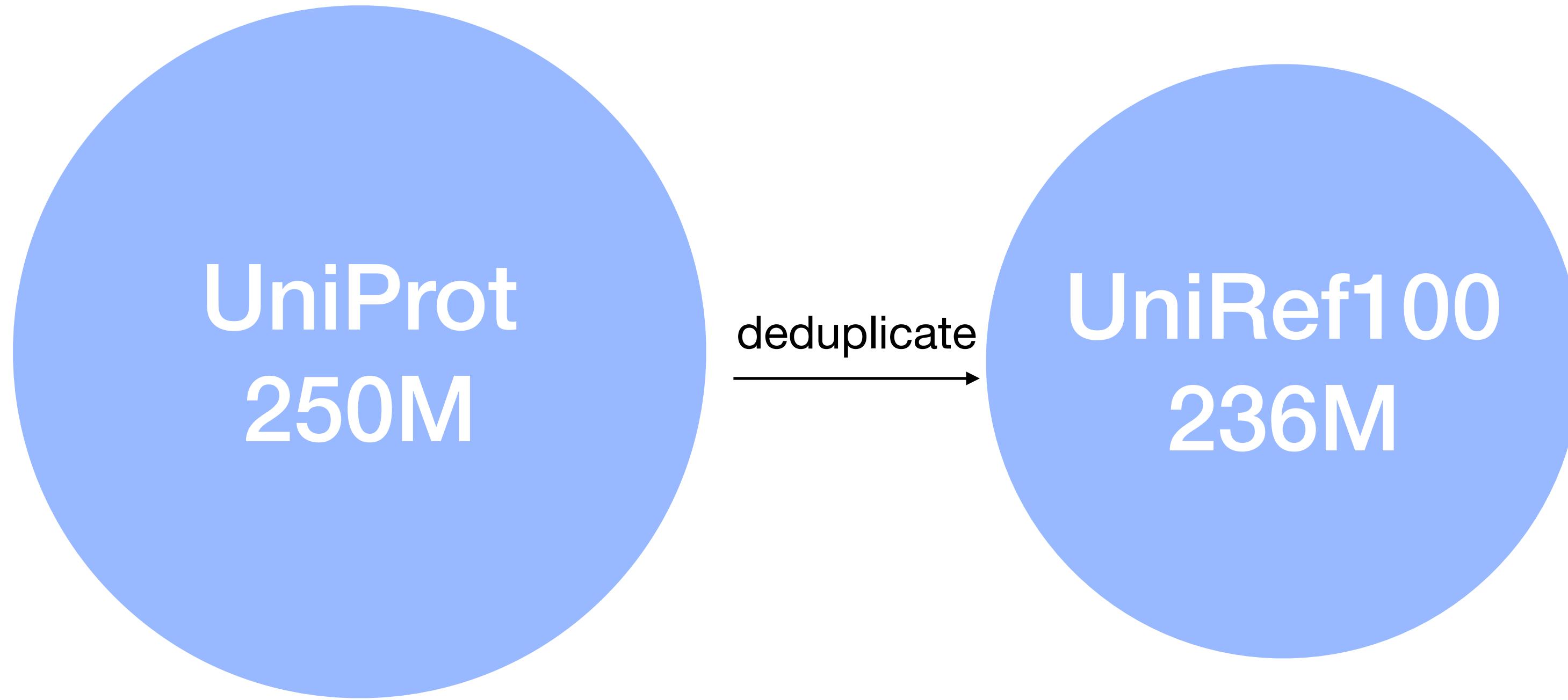
# Pretraining can leverage large sequence and structure databases



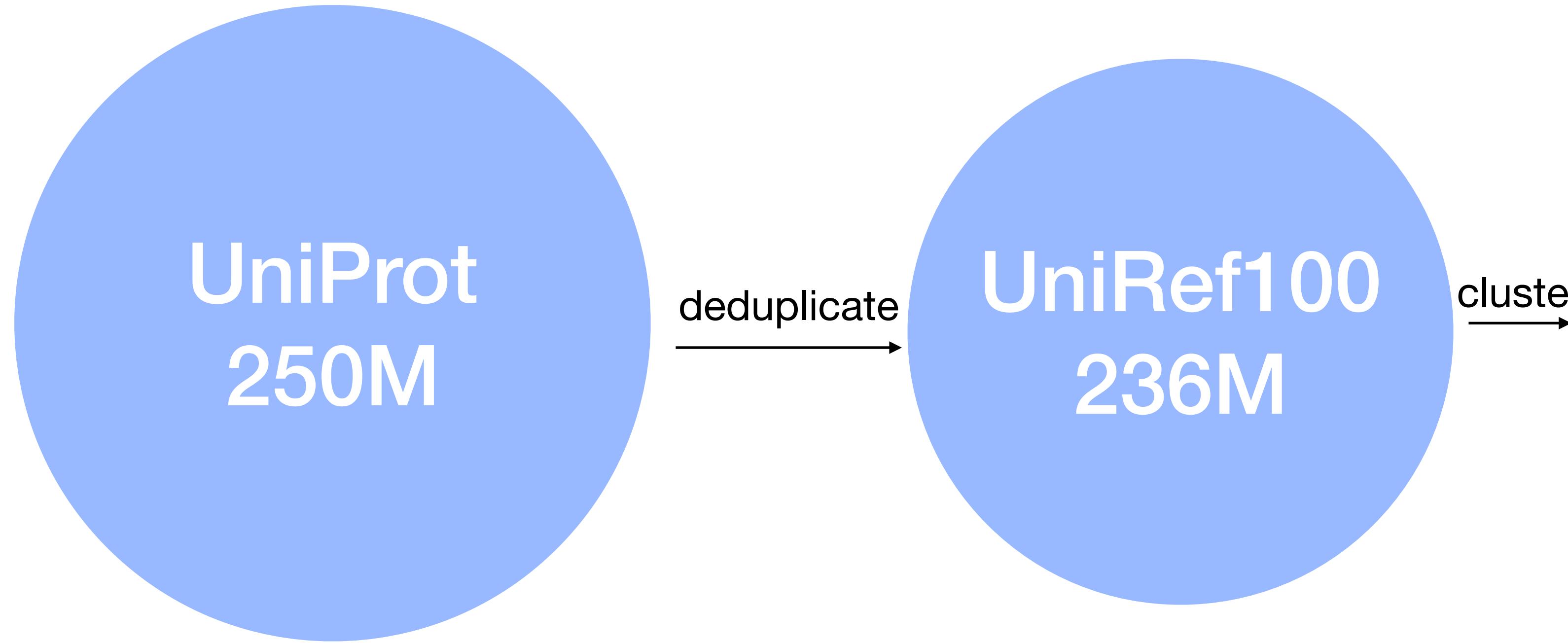
# Pretraining can leverage large sequence and structure databases



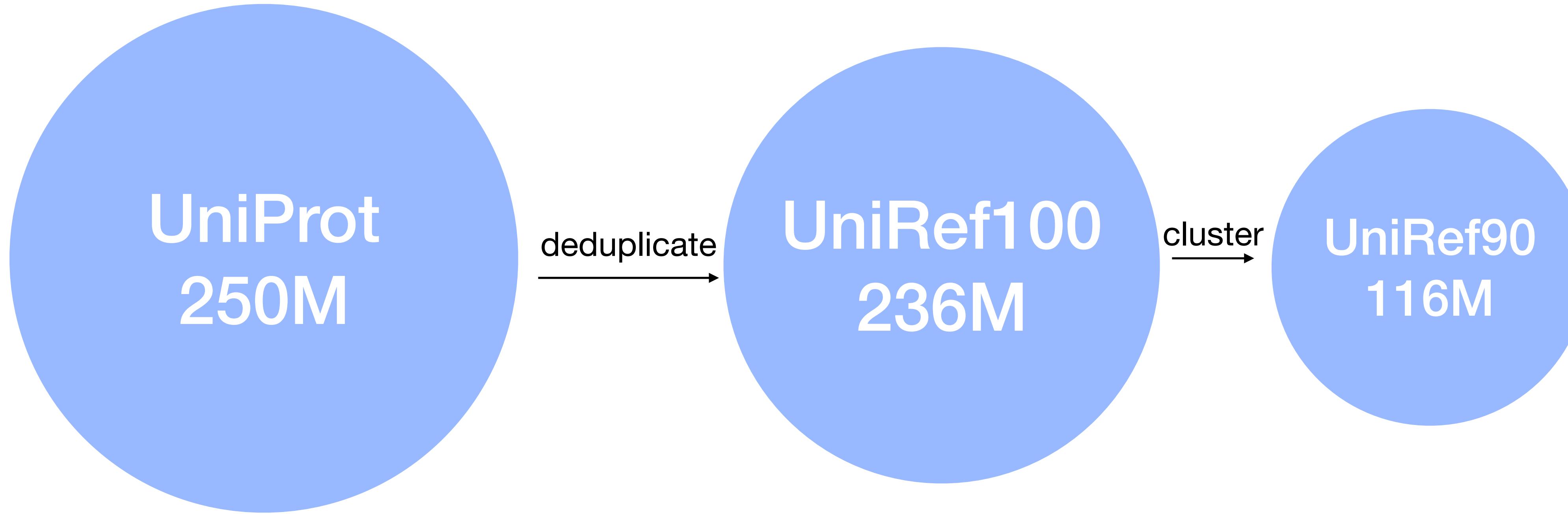
# Pretraining can leverage large sequence and structure databases



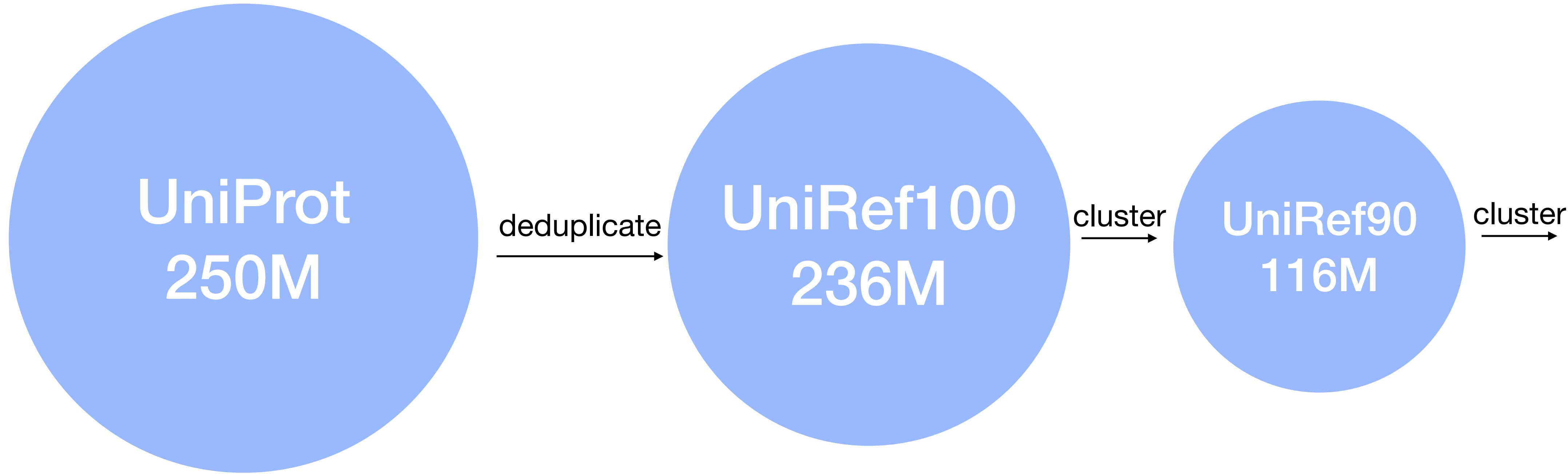
# Pretraining can leverage large sequence and structure databases



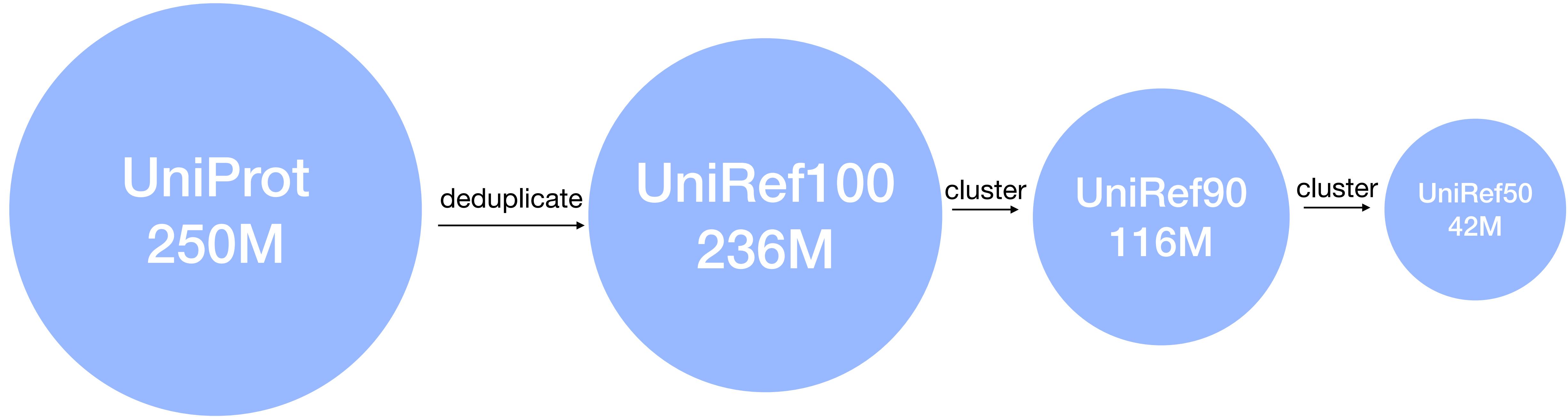
# Pretraining can leverage large sequence and structure databases



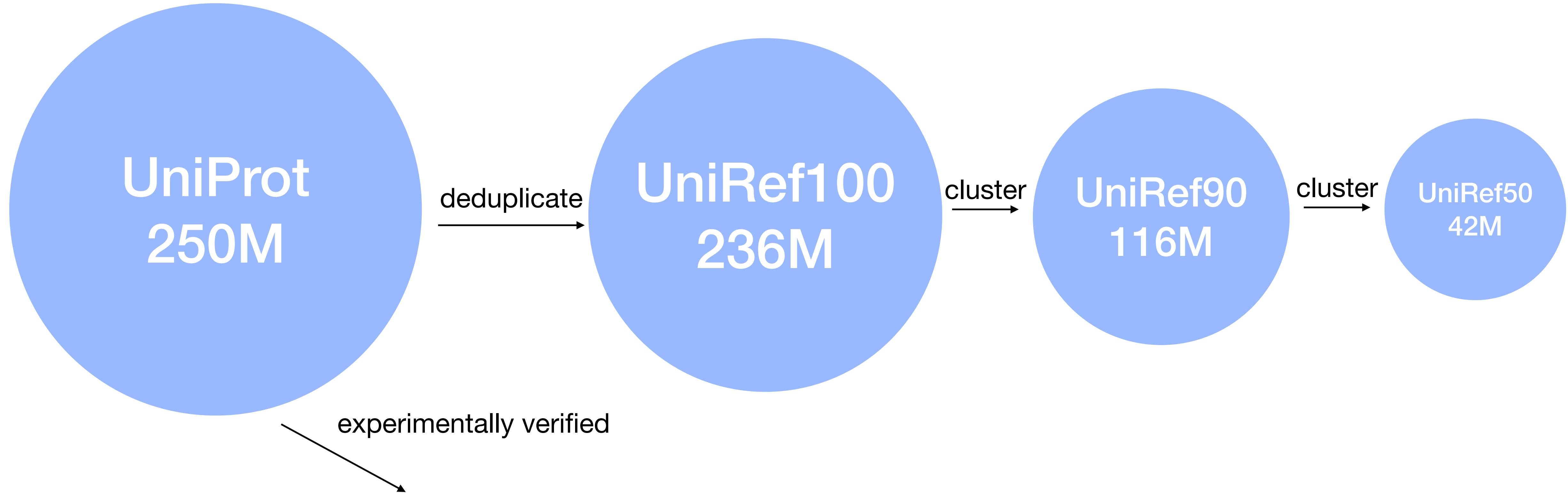
# Pretraining can leverage large sequence and structure databases



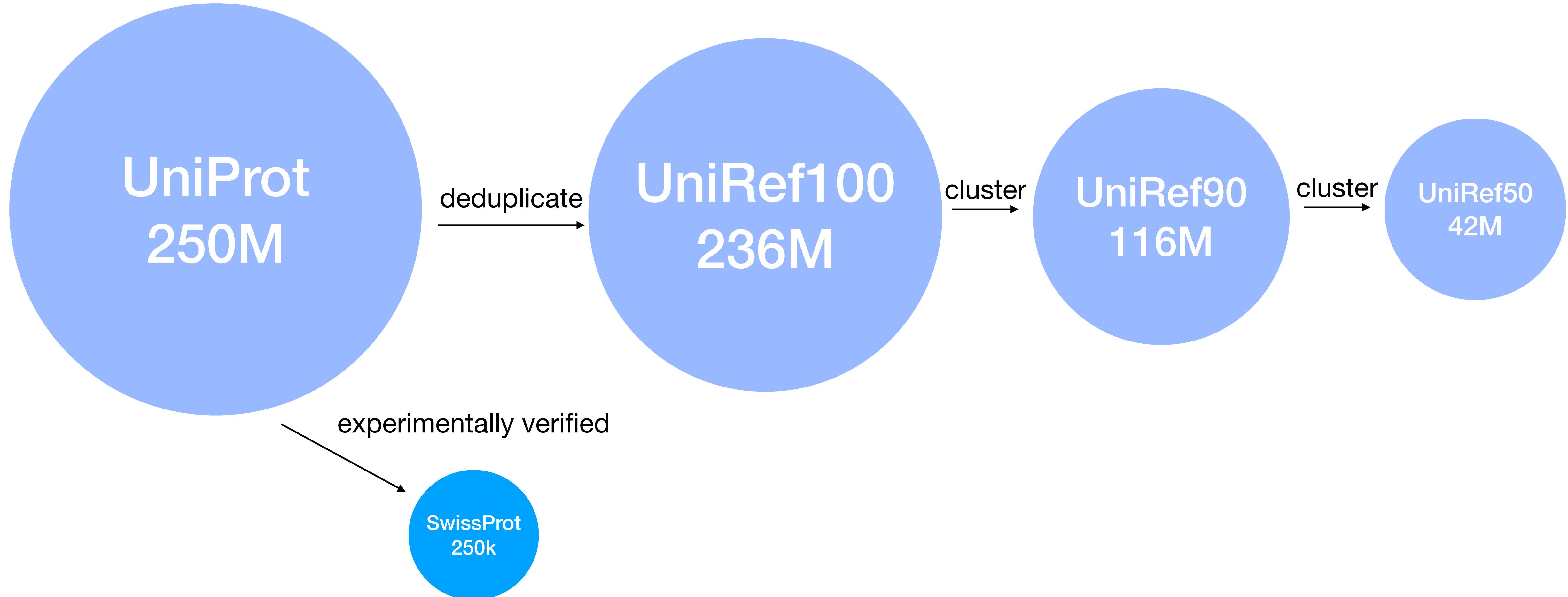
# Pretraining can leverage large sequence and structure databases



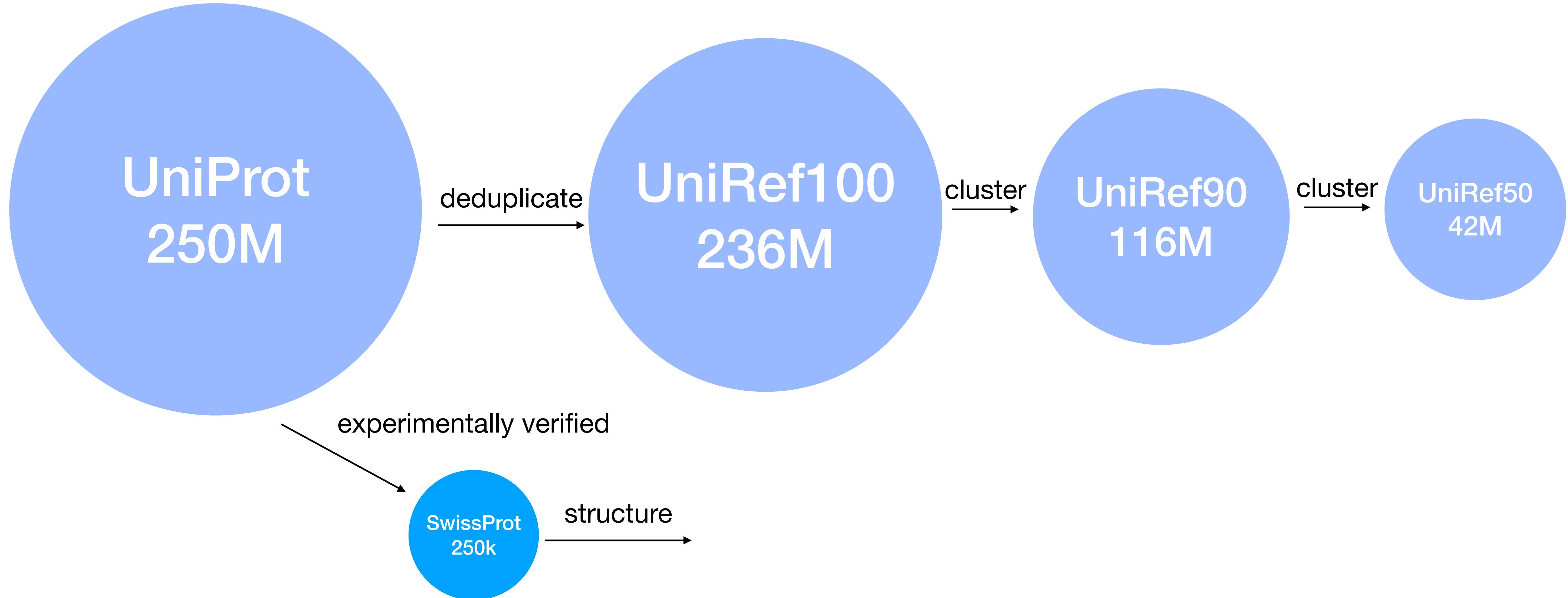
# Pretraining can leverage large sequence and structure databases



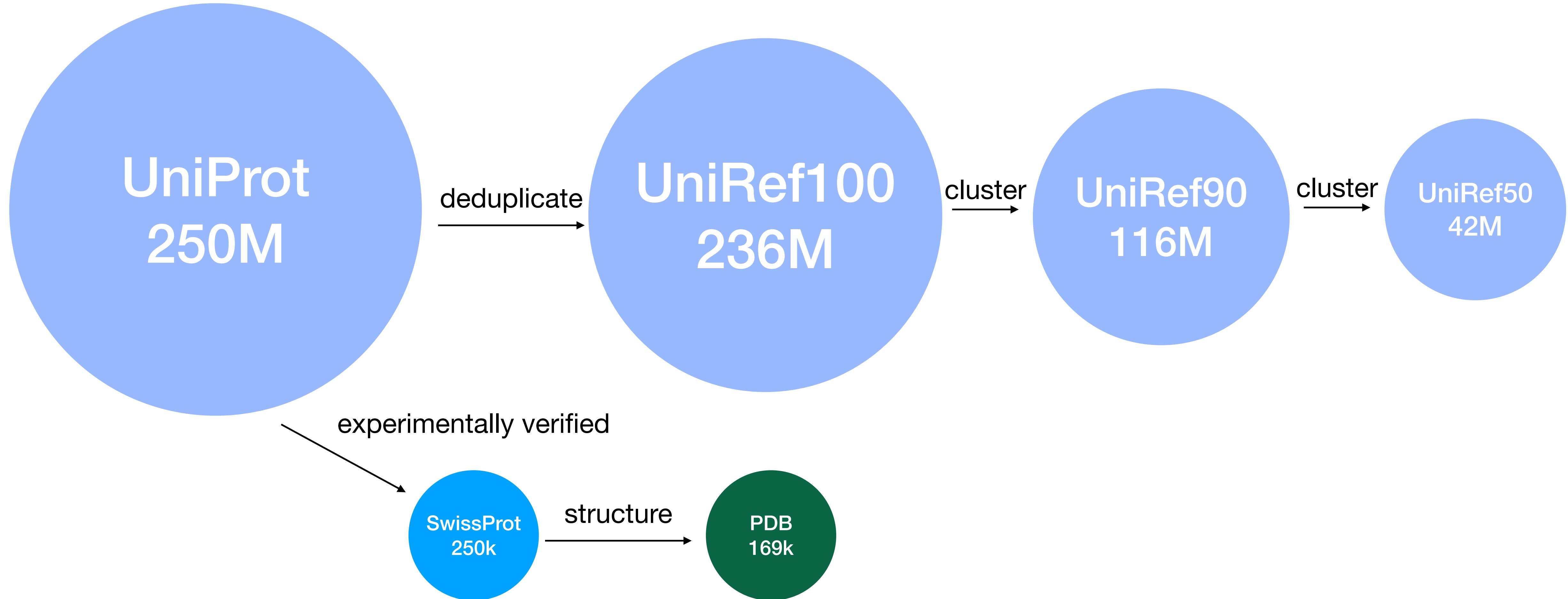
# Pretraining can leverage large sequence and structure databases



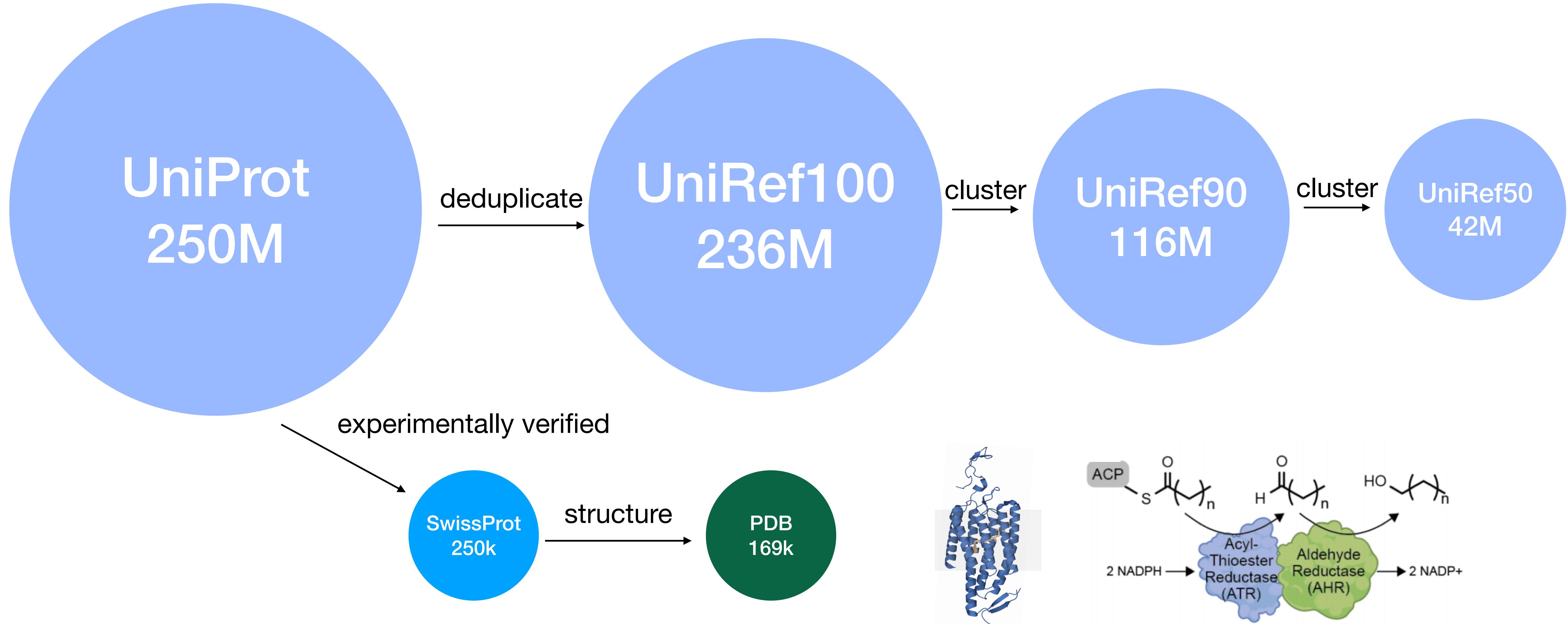
# Pretraining can leverage large sequence and structure databases



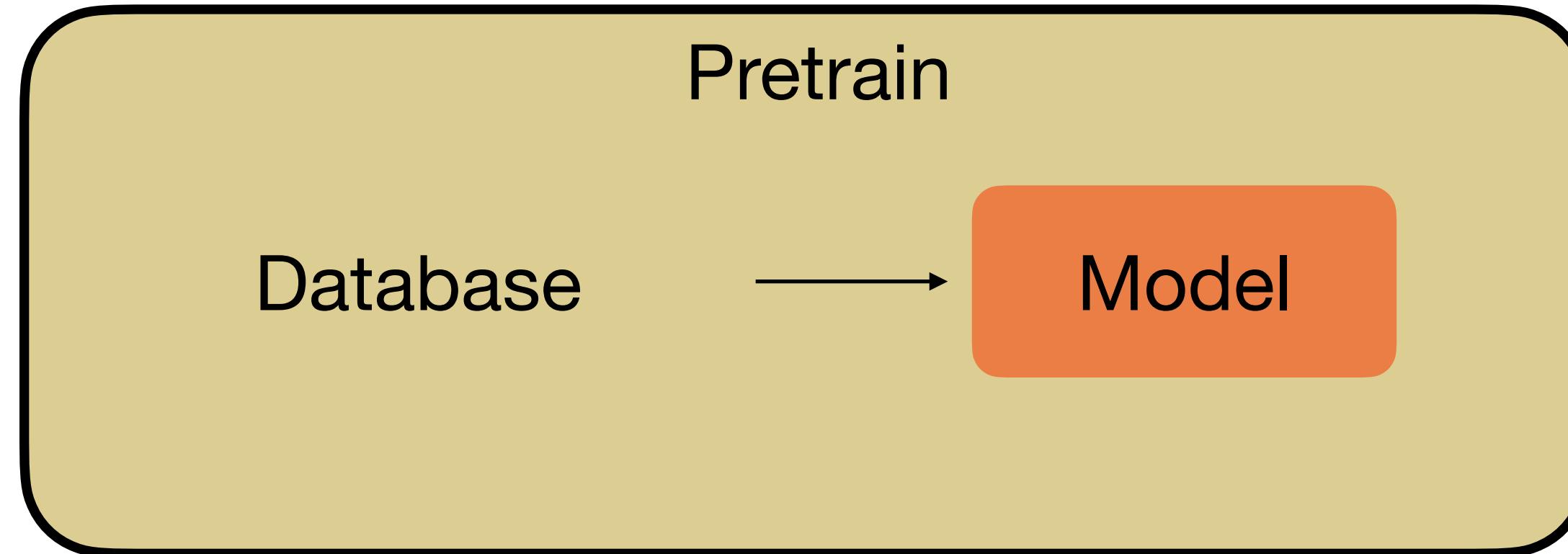
# Pretraining can leverage large sequence and structure databases



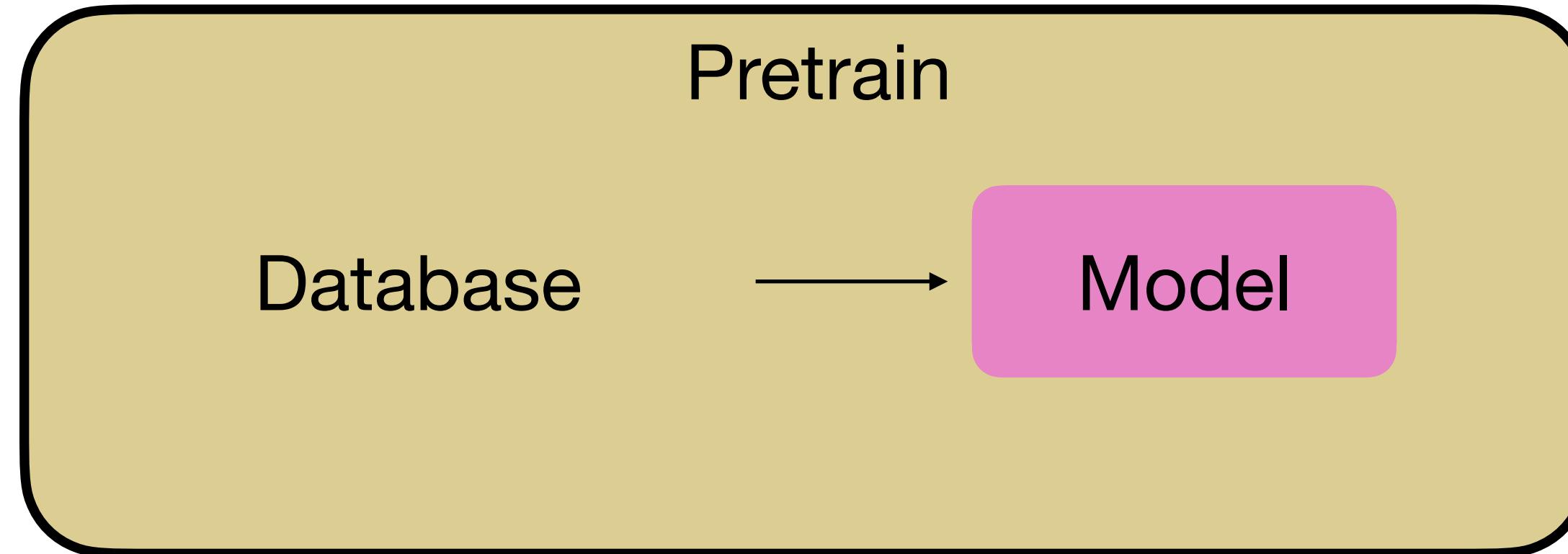
# Pretraining can leverage large sequence and structure databases



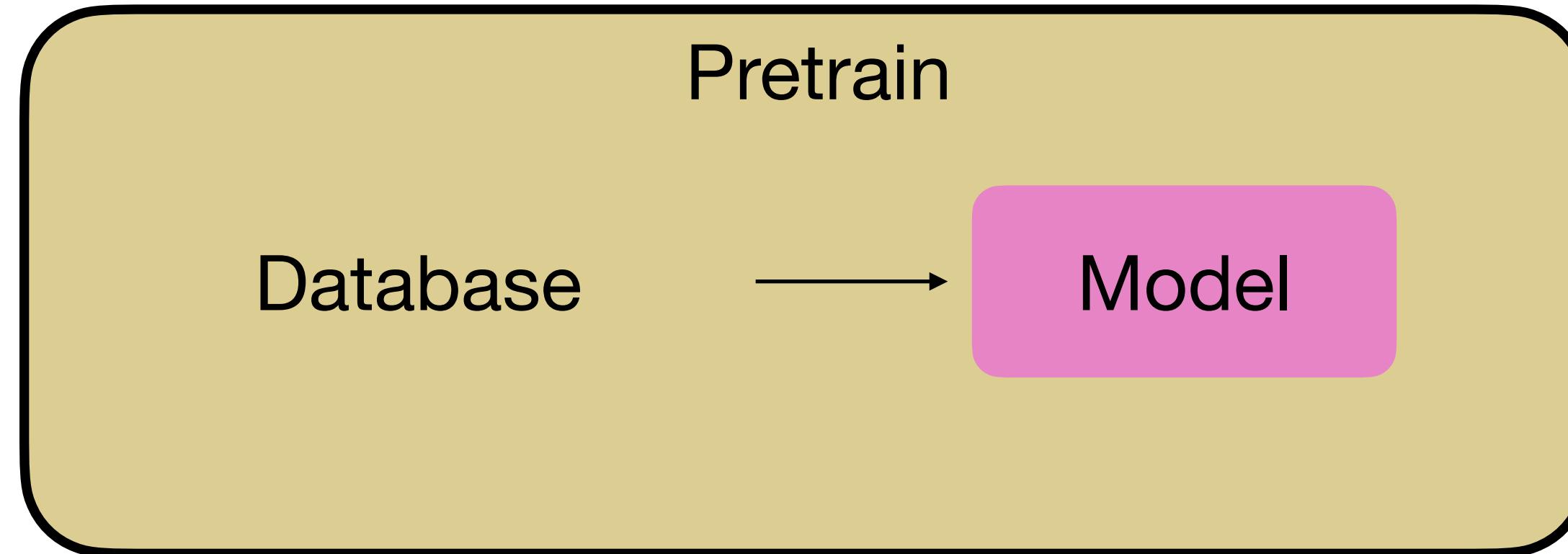
# Pretraining allows us to transfer knowledge



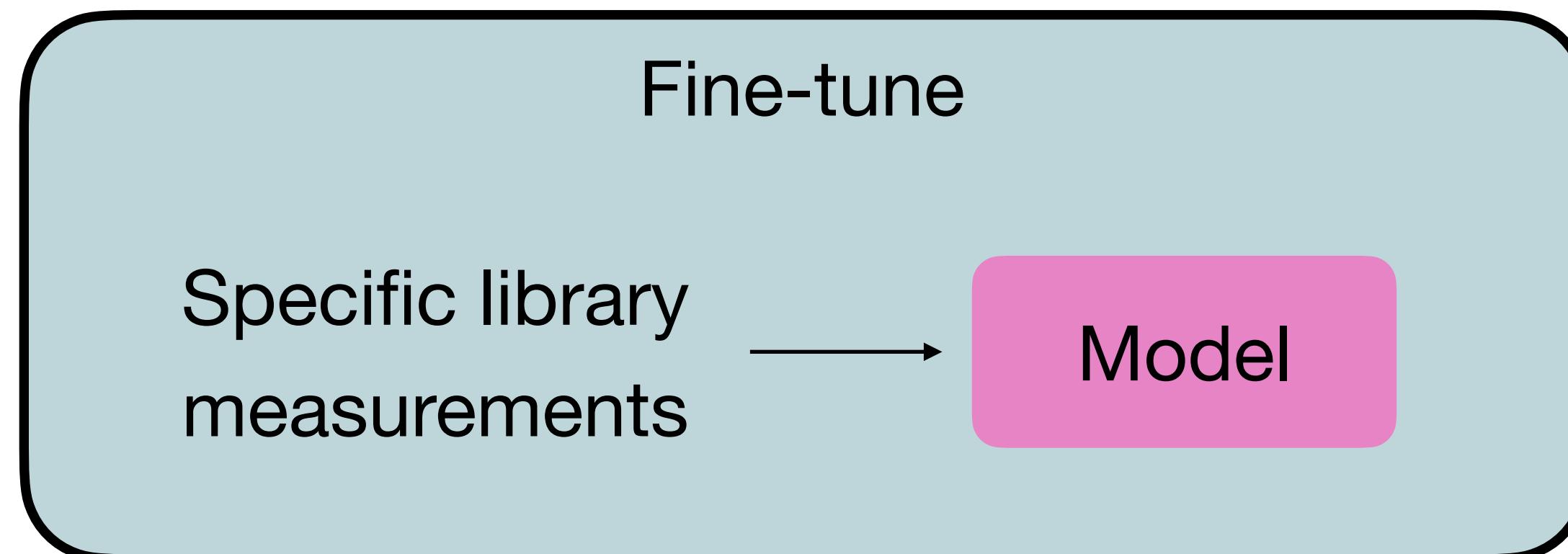
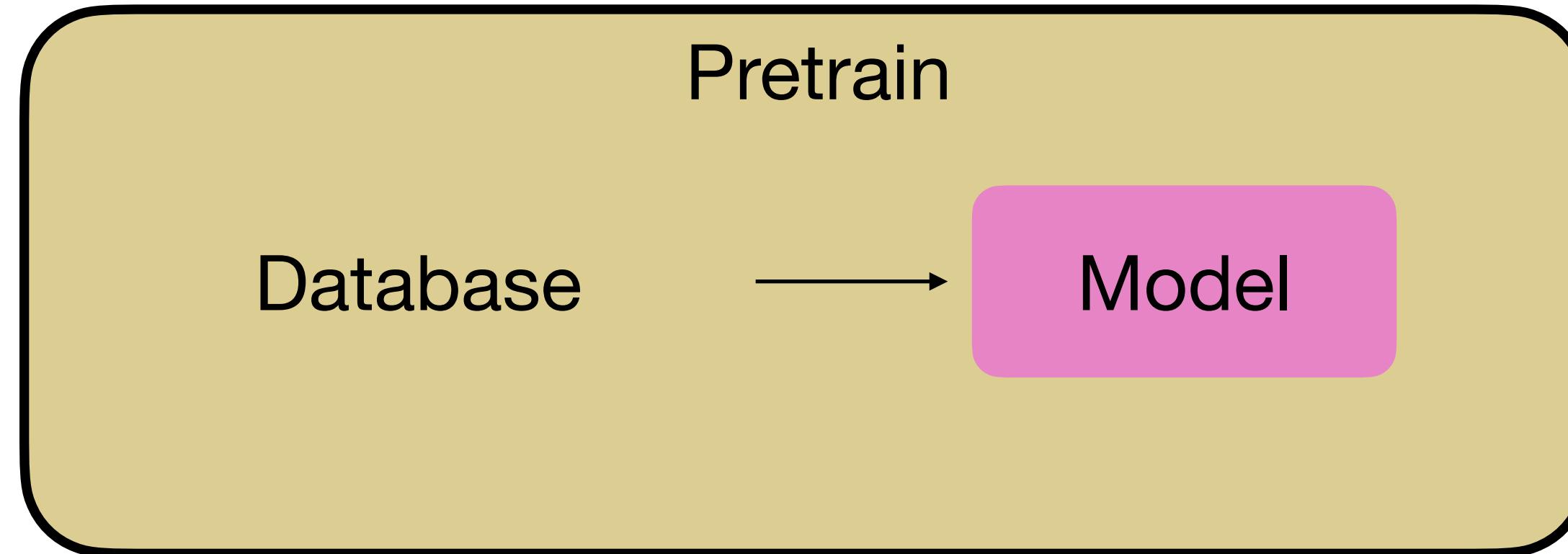
# Pretraining allows us to transfer knowledge



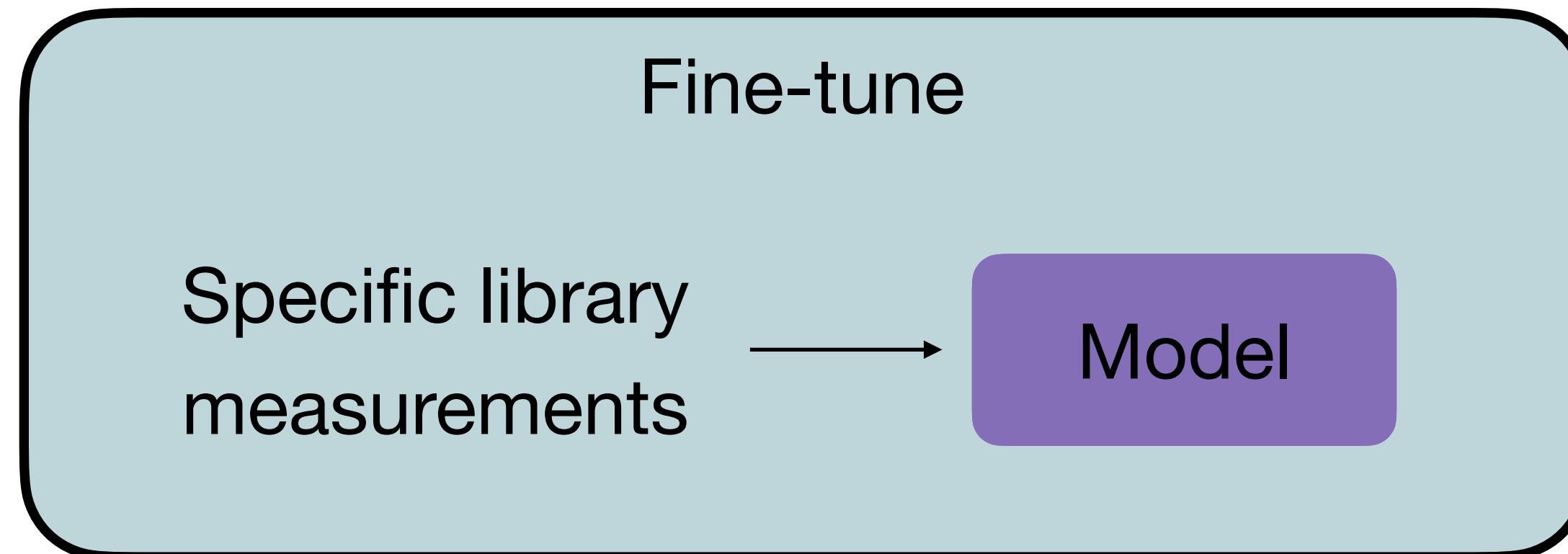
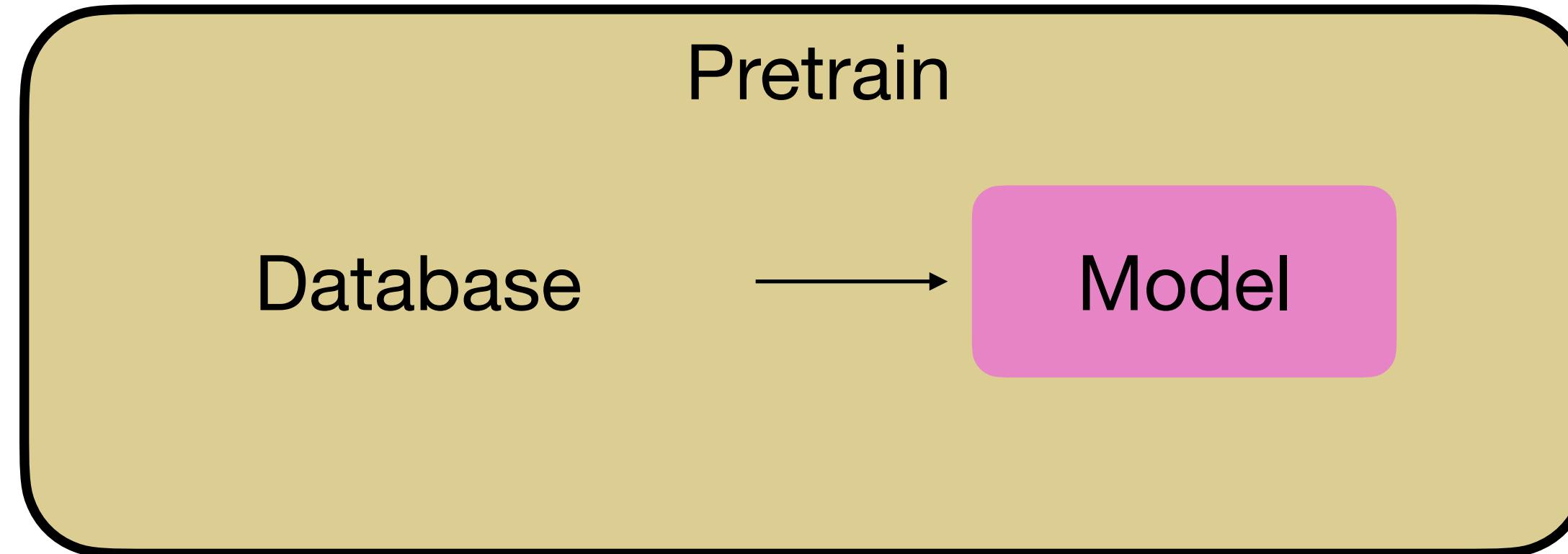
# Pretraining allows us to transfer knowledge



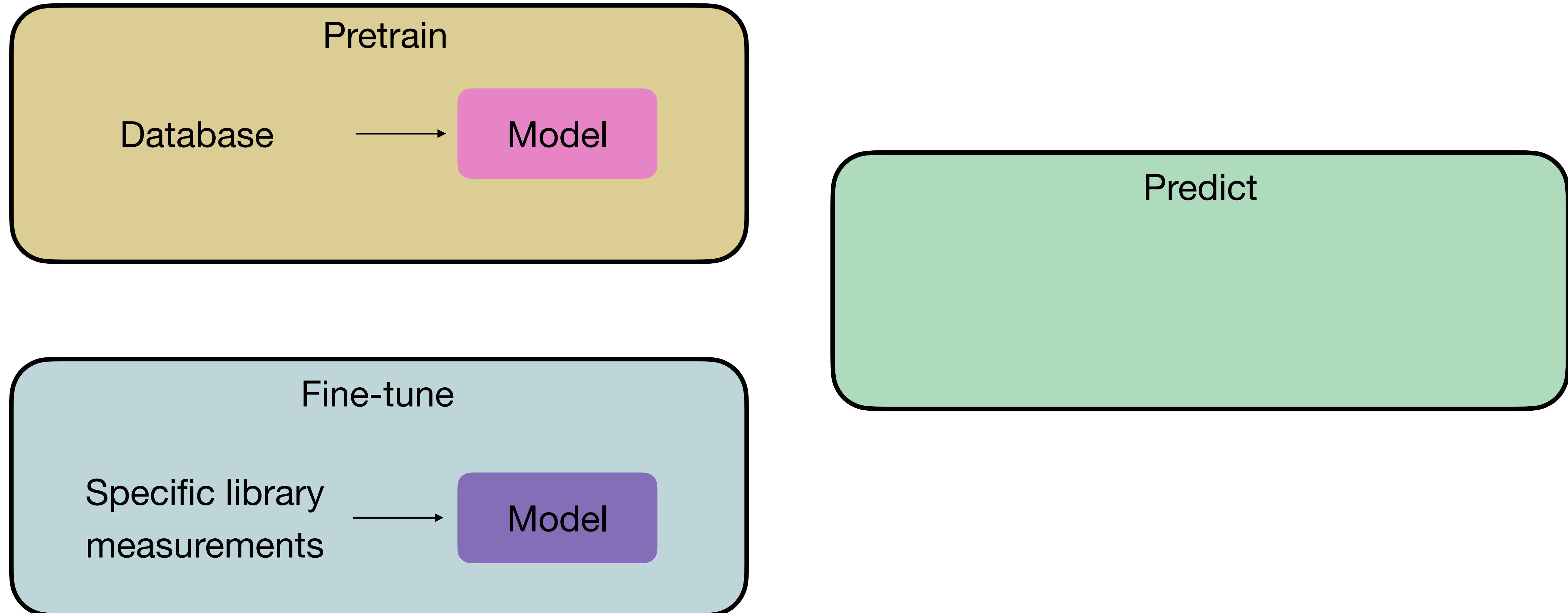
# Pretraining allows us to transfer knowledge



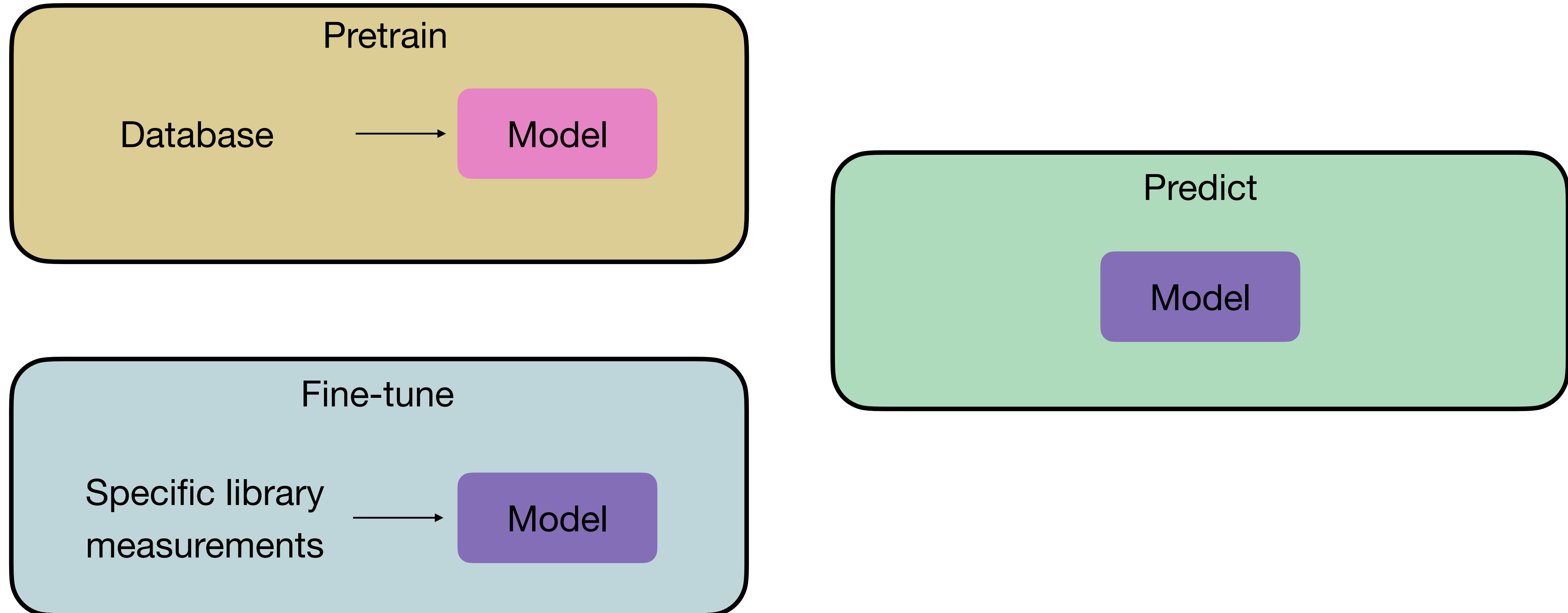
# Pretraining allows us to transfer knowledge



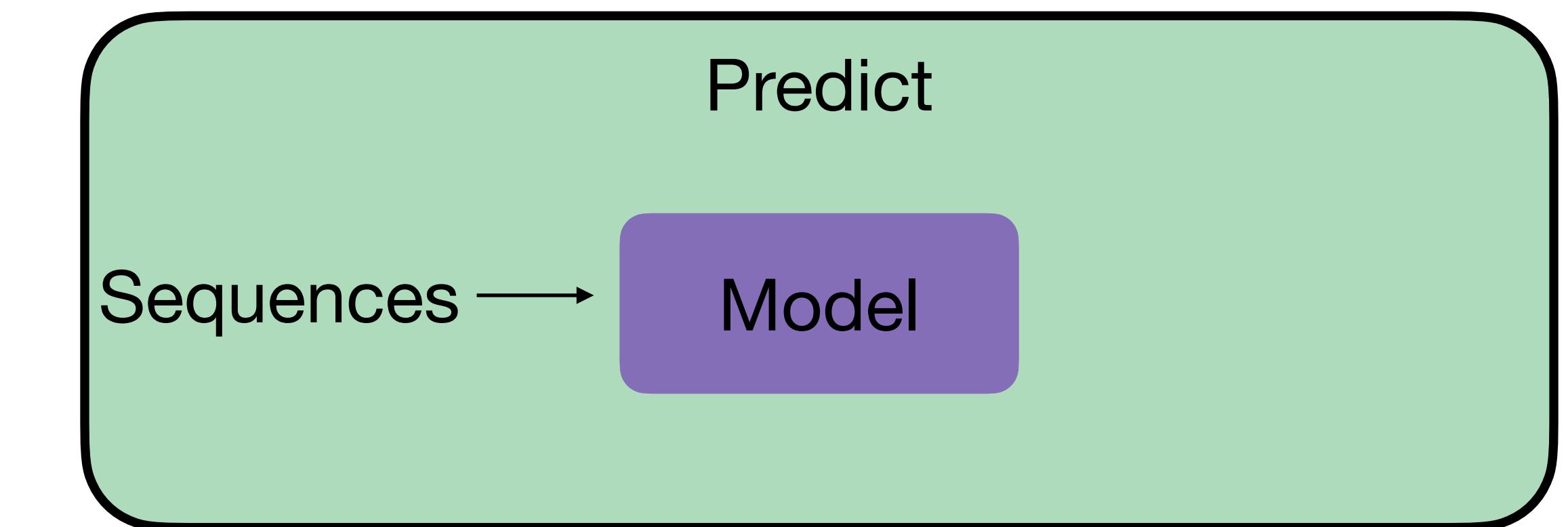
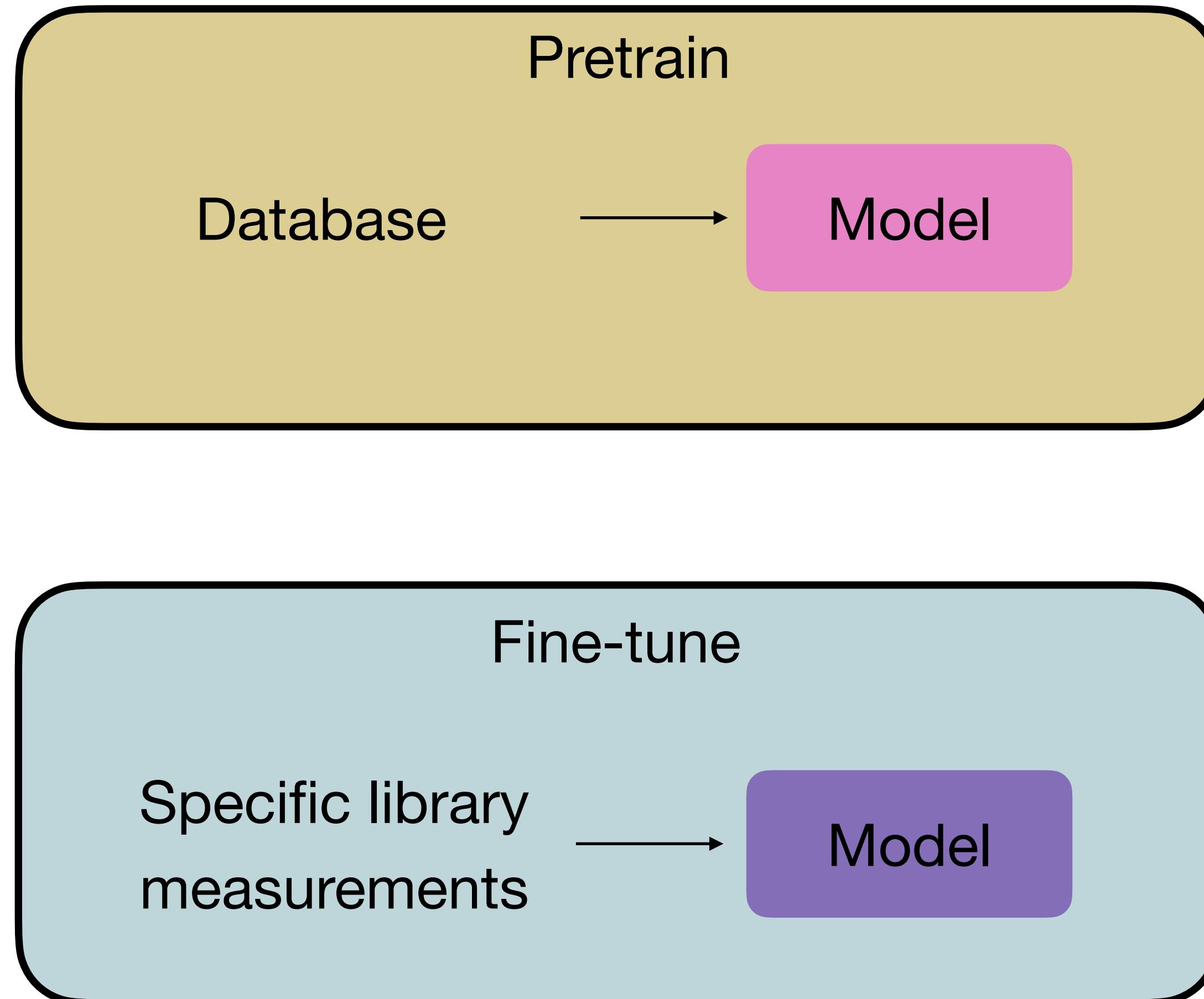
# Pretraining allows us to transfer knowledge



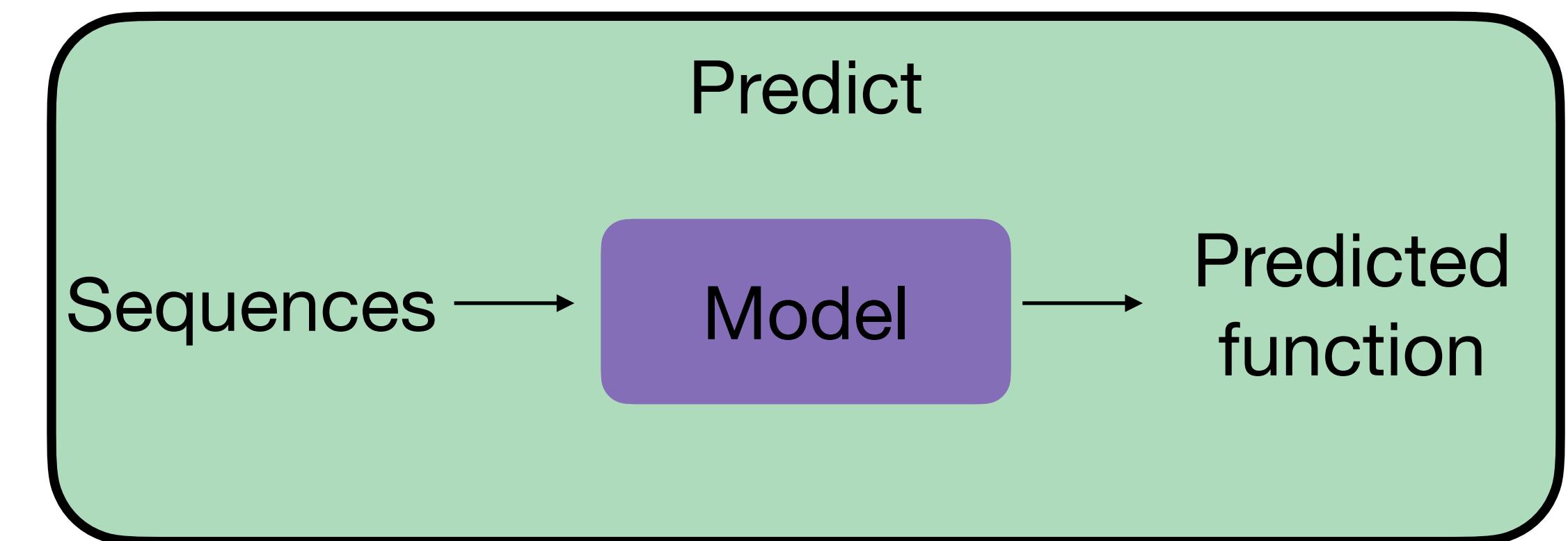
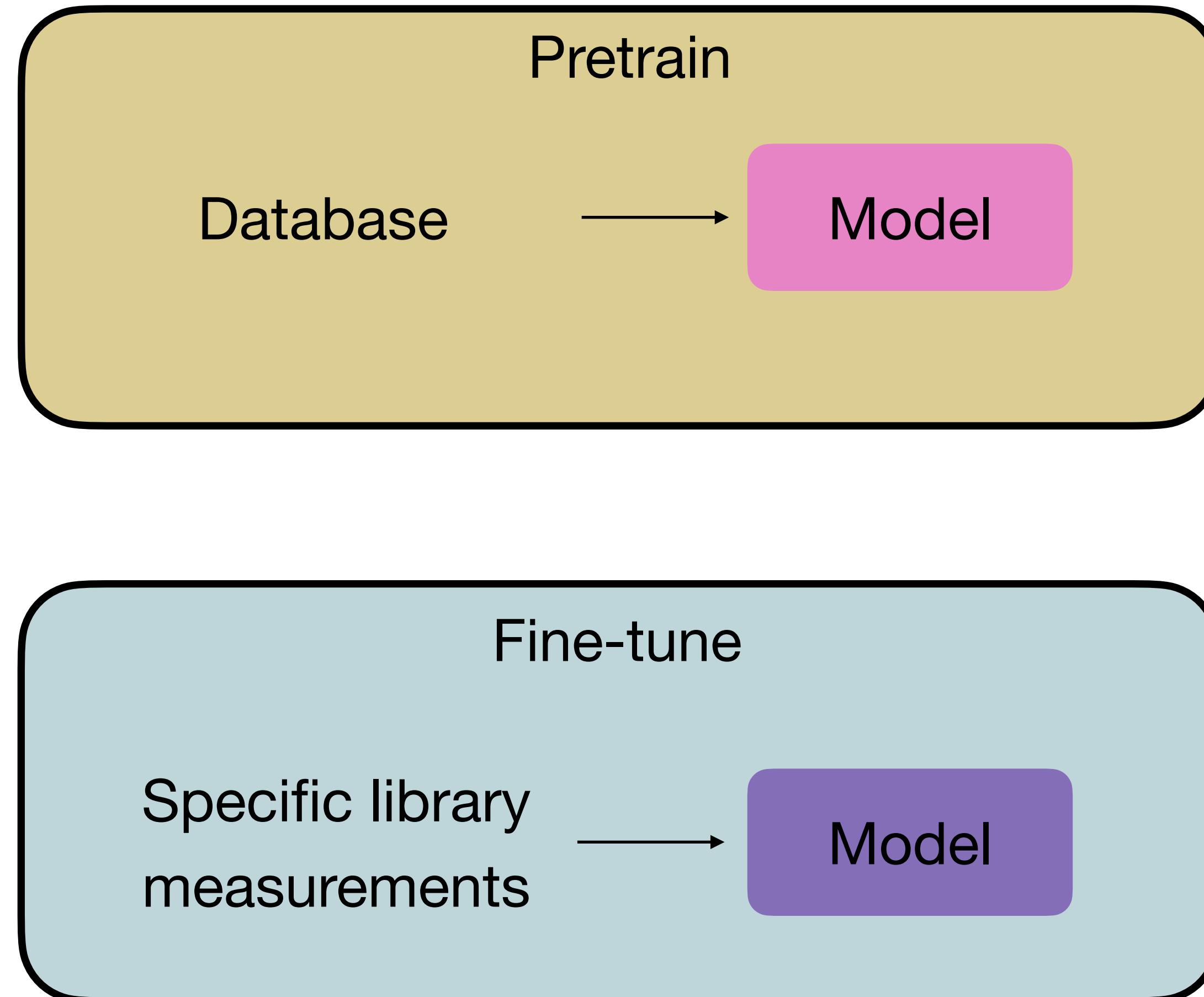
# Pretraining allows us to transfer knowledge



# Pretraining allows us to transfer knowledge



# Pretraining allows us to transfer knowledge



# Many methods pretend proteins are language

MFTGNDAGH

# Many methods pretend proteins are language

MFTGNDAGH

Treat amino acids like tokens

# Many methods pretend proteins are language

MFTGNDAGH

Treat amino acids like tokens

Use tools originally developed for language

# Many methods pretend proteins are language

MFTGNDAGH

# Many methods pretend proteins are language

MFTGNDAGH

# Many methods pretend proteins are language

Word2Vec

MFTGNDAGH

# Many methods pretend proteins are language

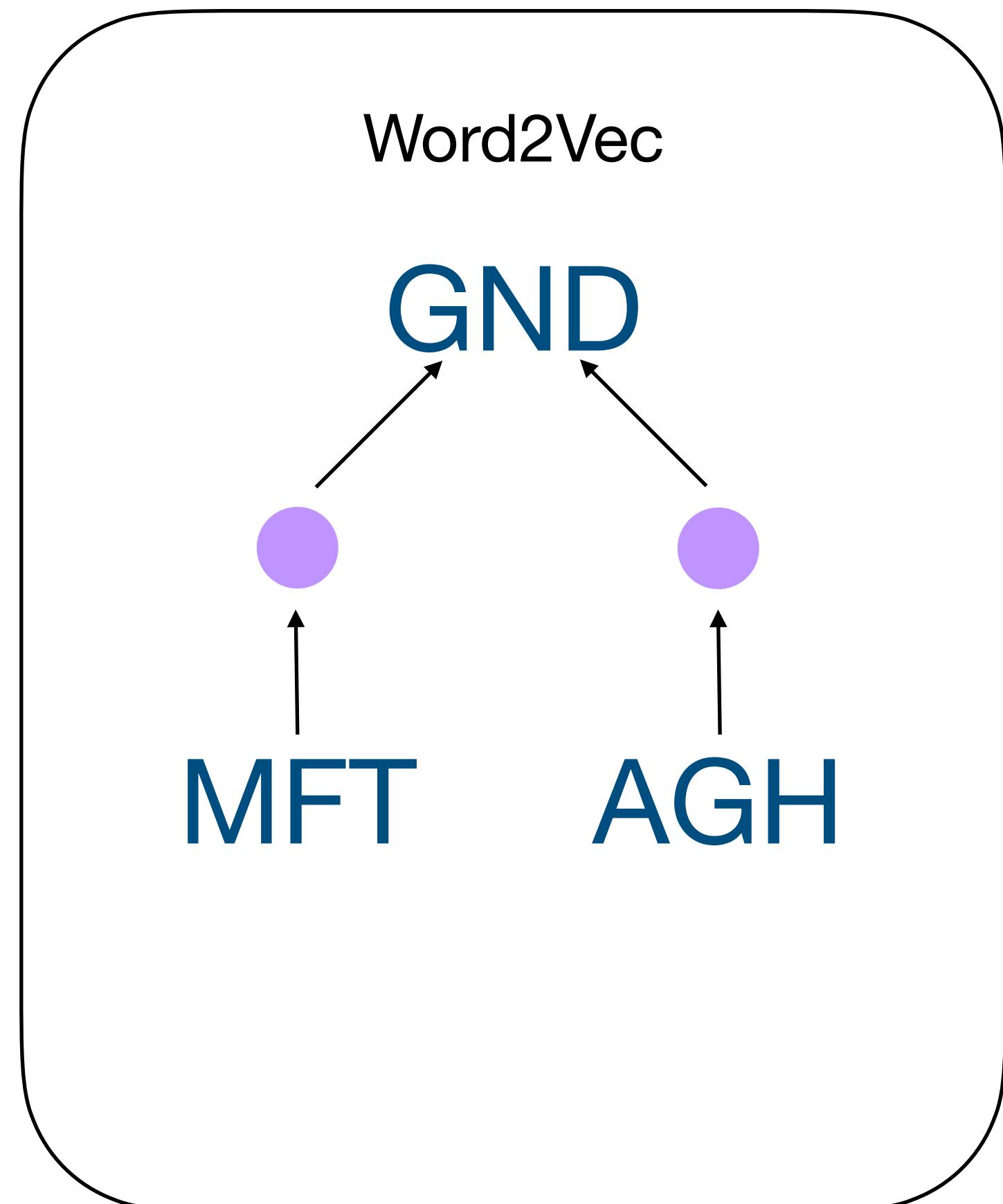
Word2Vec

MFT

AGH

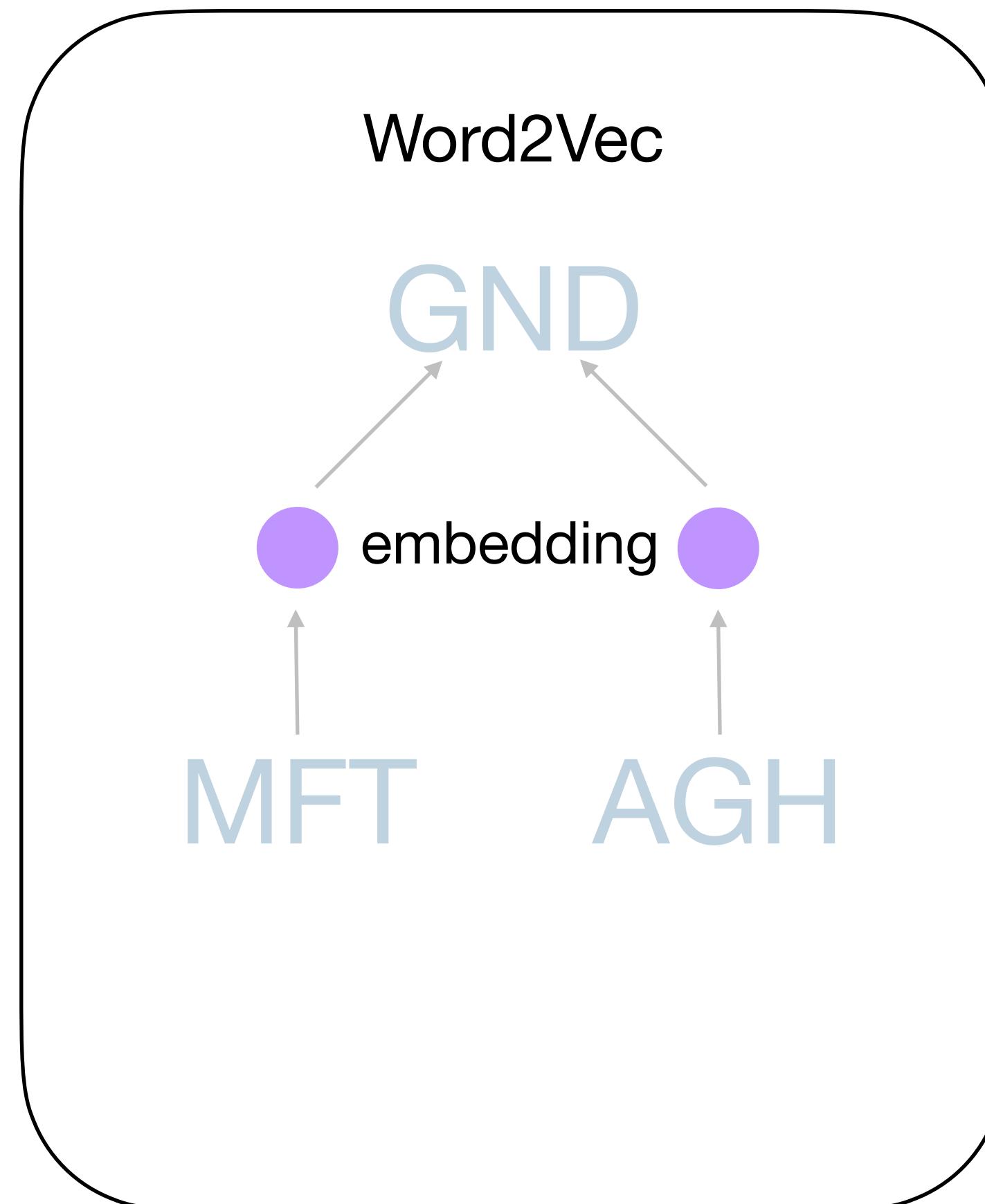
MFTGNDAGH

# Many methods pretend proteins are language



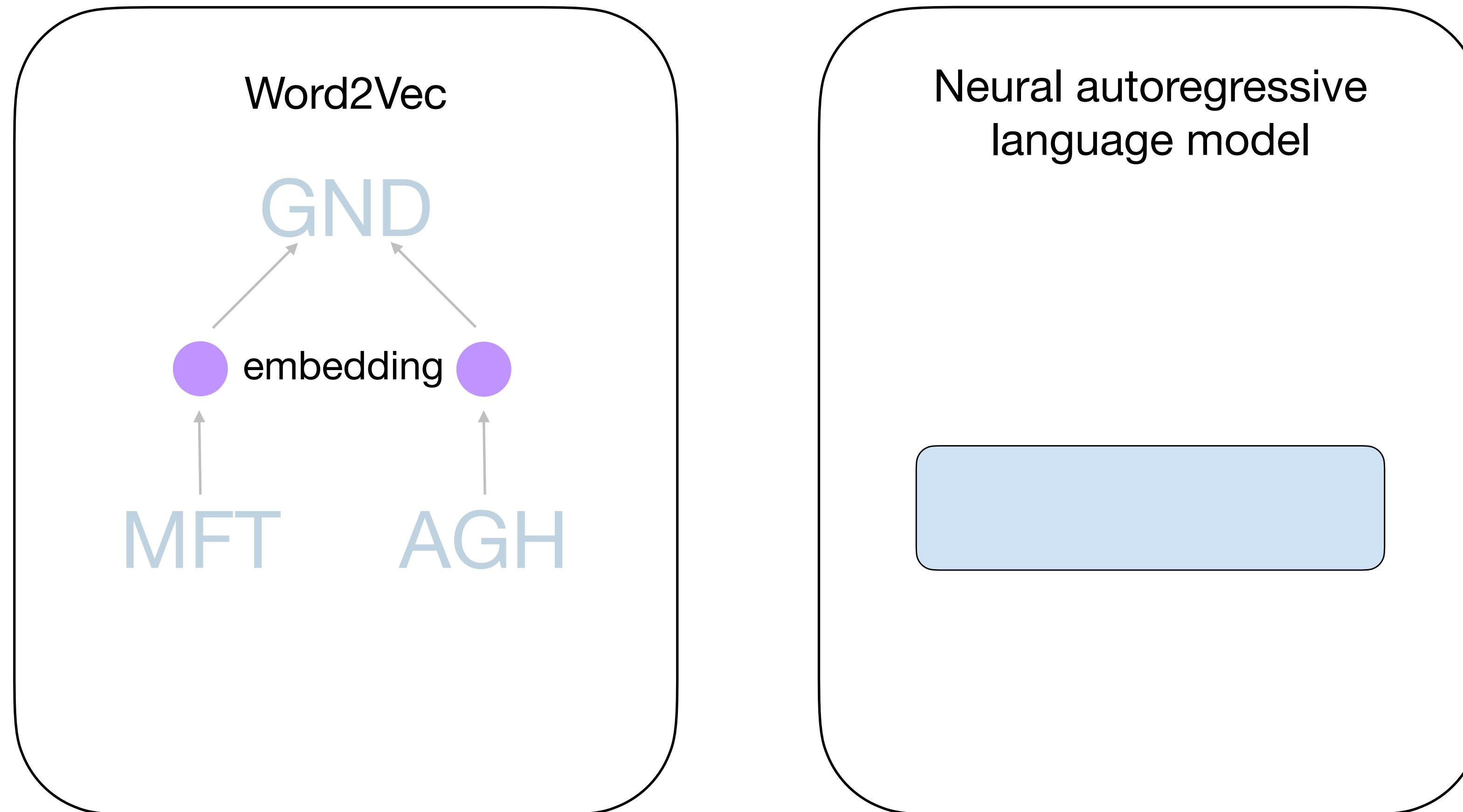
MFTGNDAGH

# Many methods pretend proteins are language



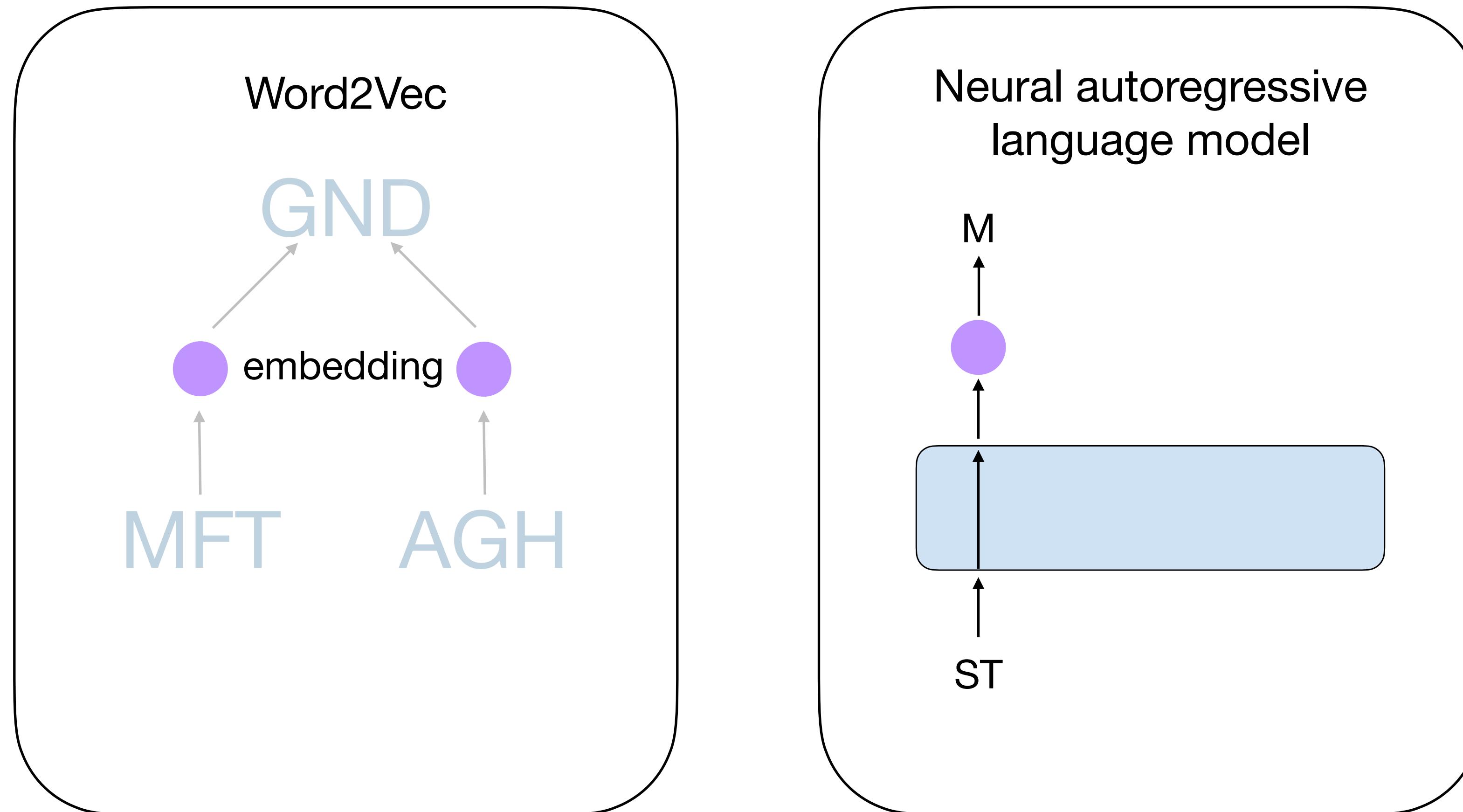
MFTGNDAGH

# Many methods pretend proteins are language



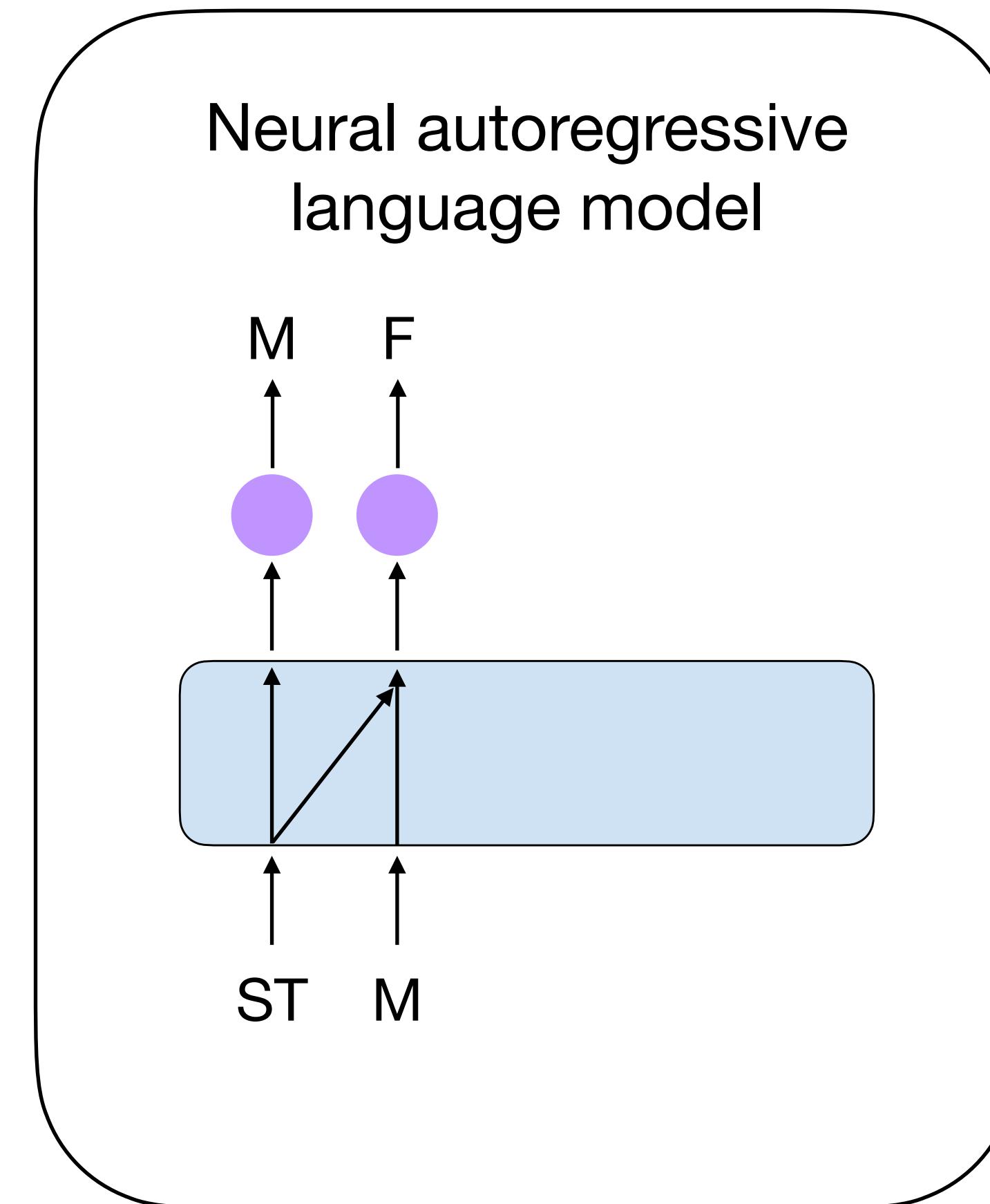
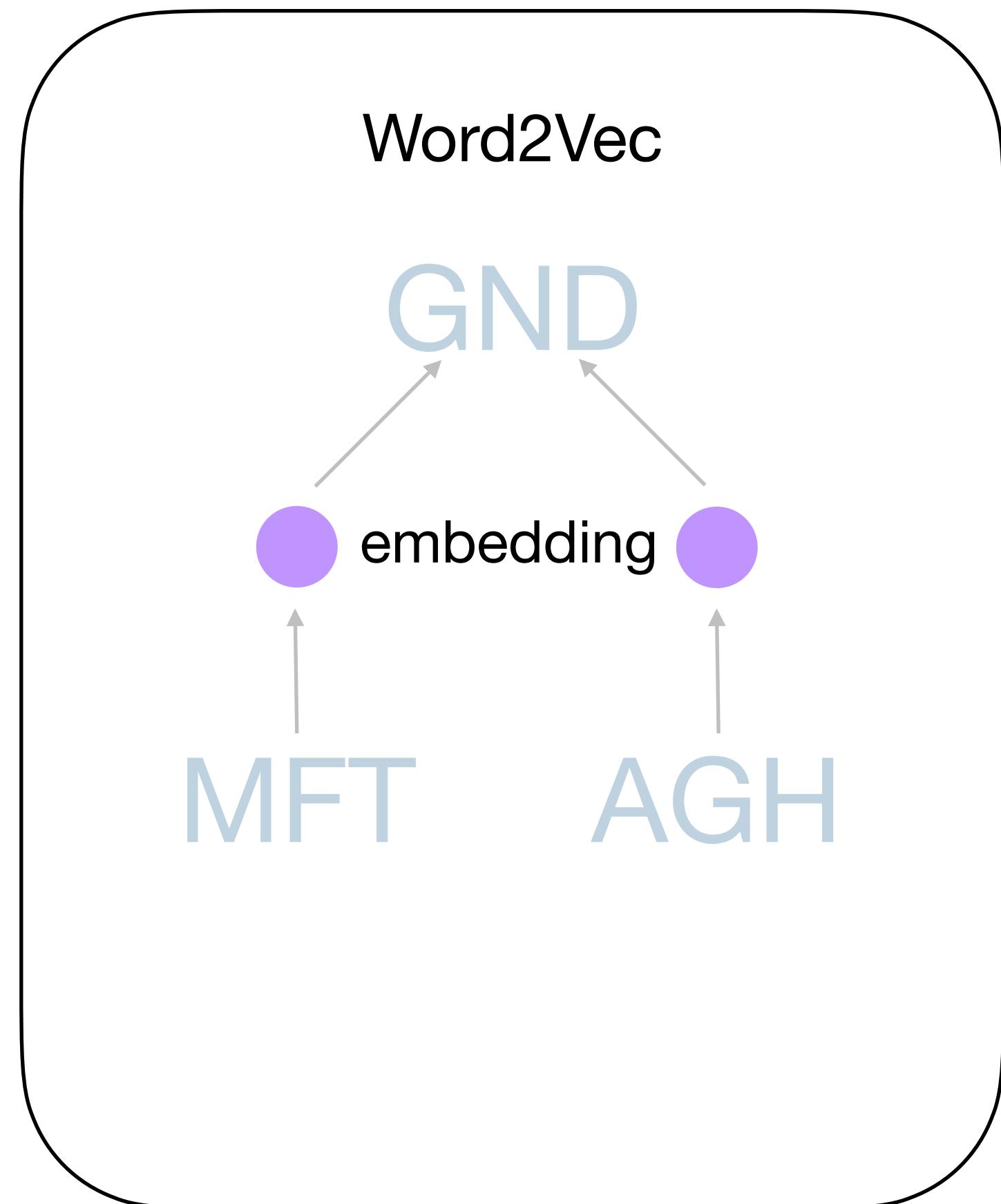
MFTGNDAGH

# Many methods pretend proteins are language



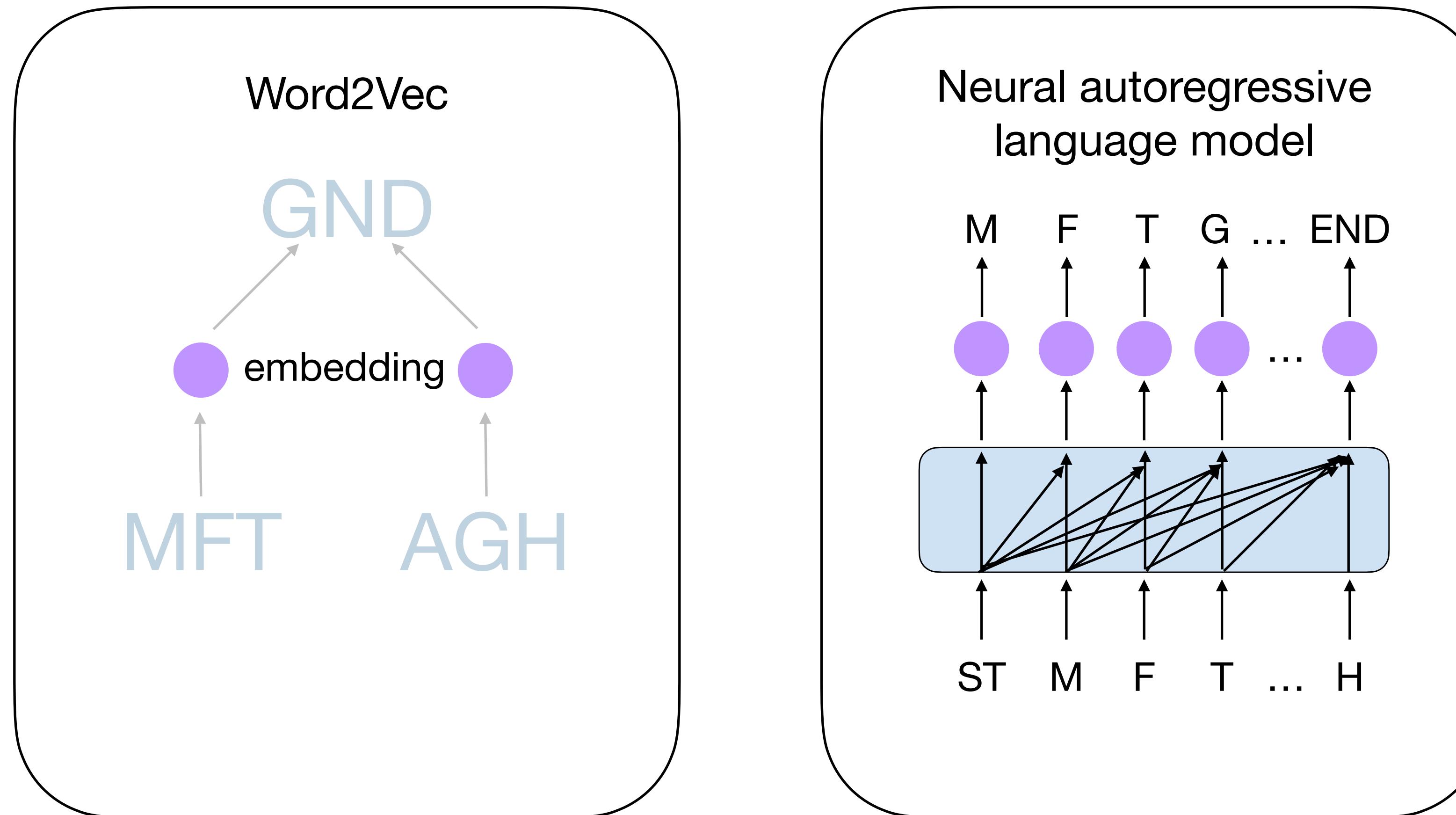
MFTGNDAGH

# Many methods pretend proteins are language



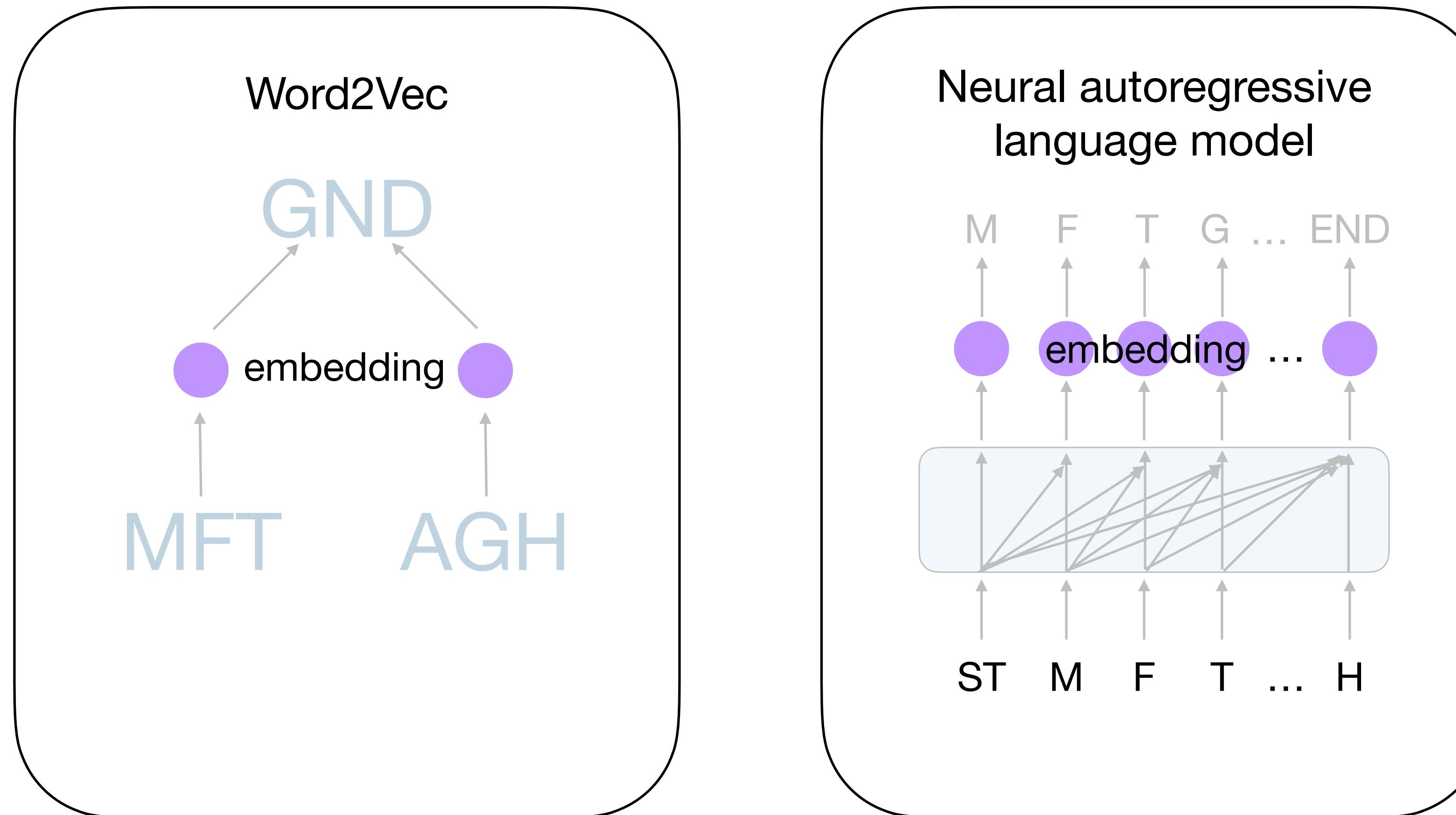
MFTGNDAGH

# Many methods pretend proteins are language



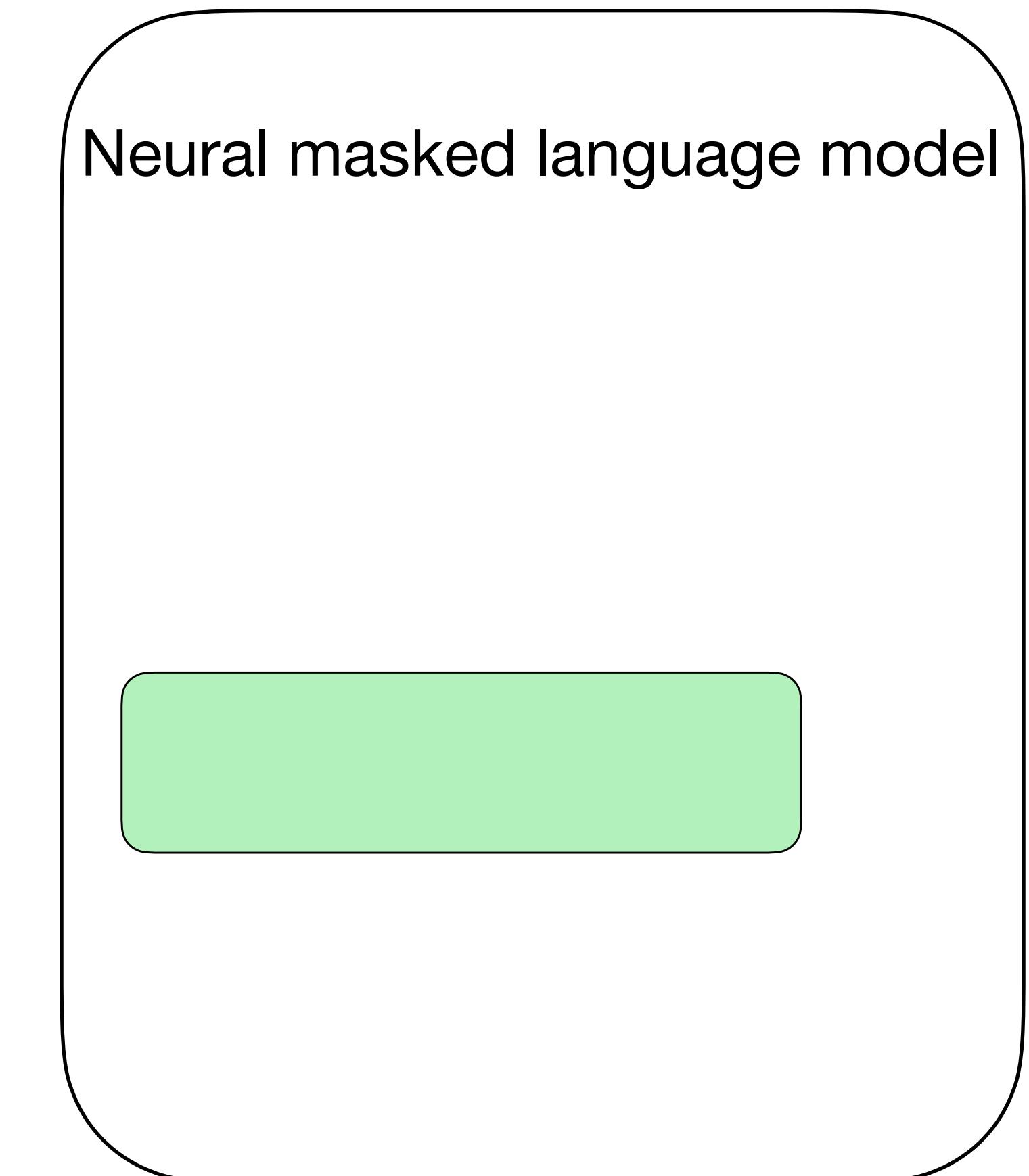
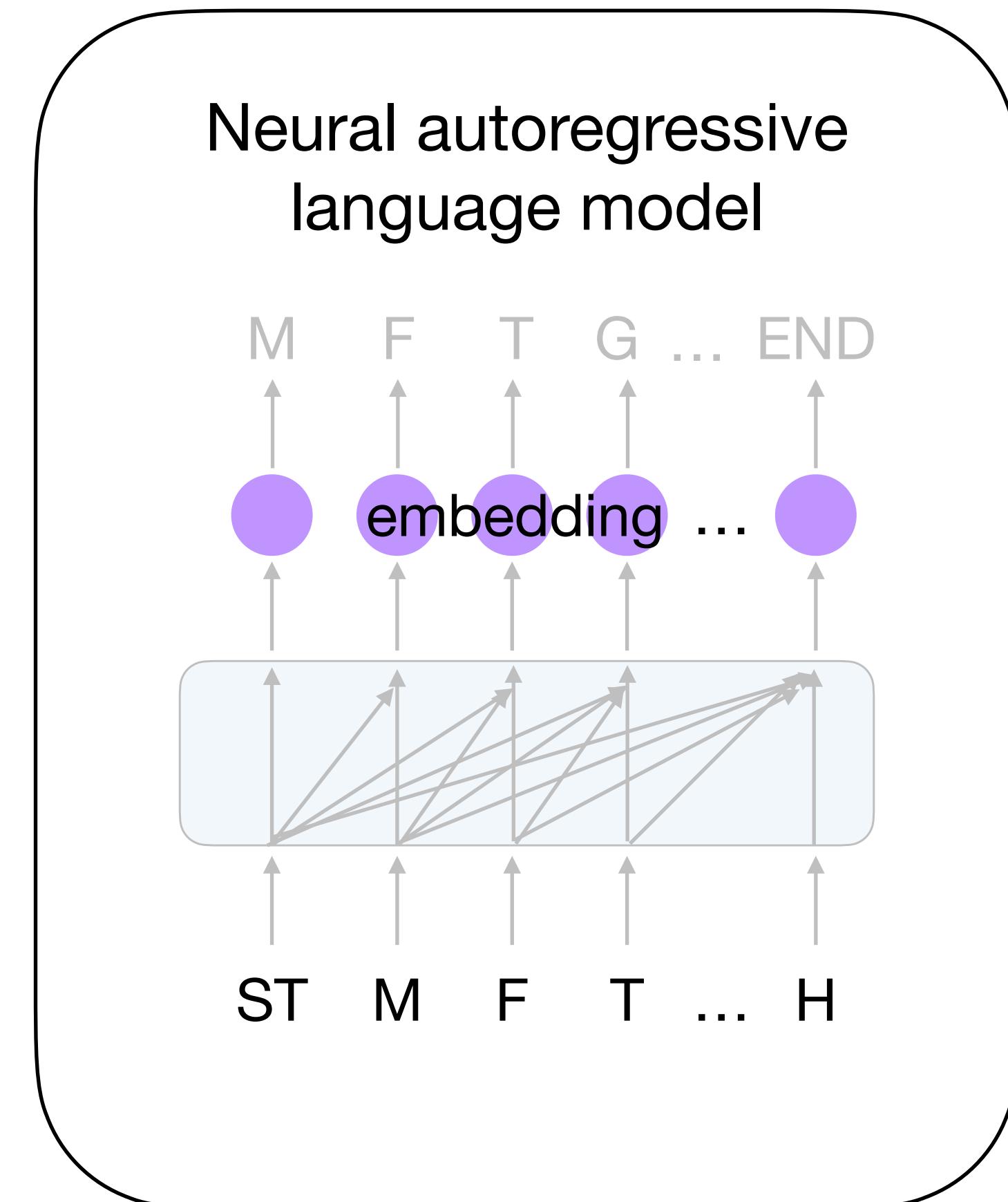
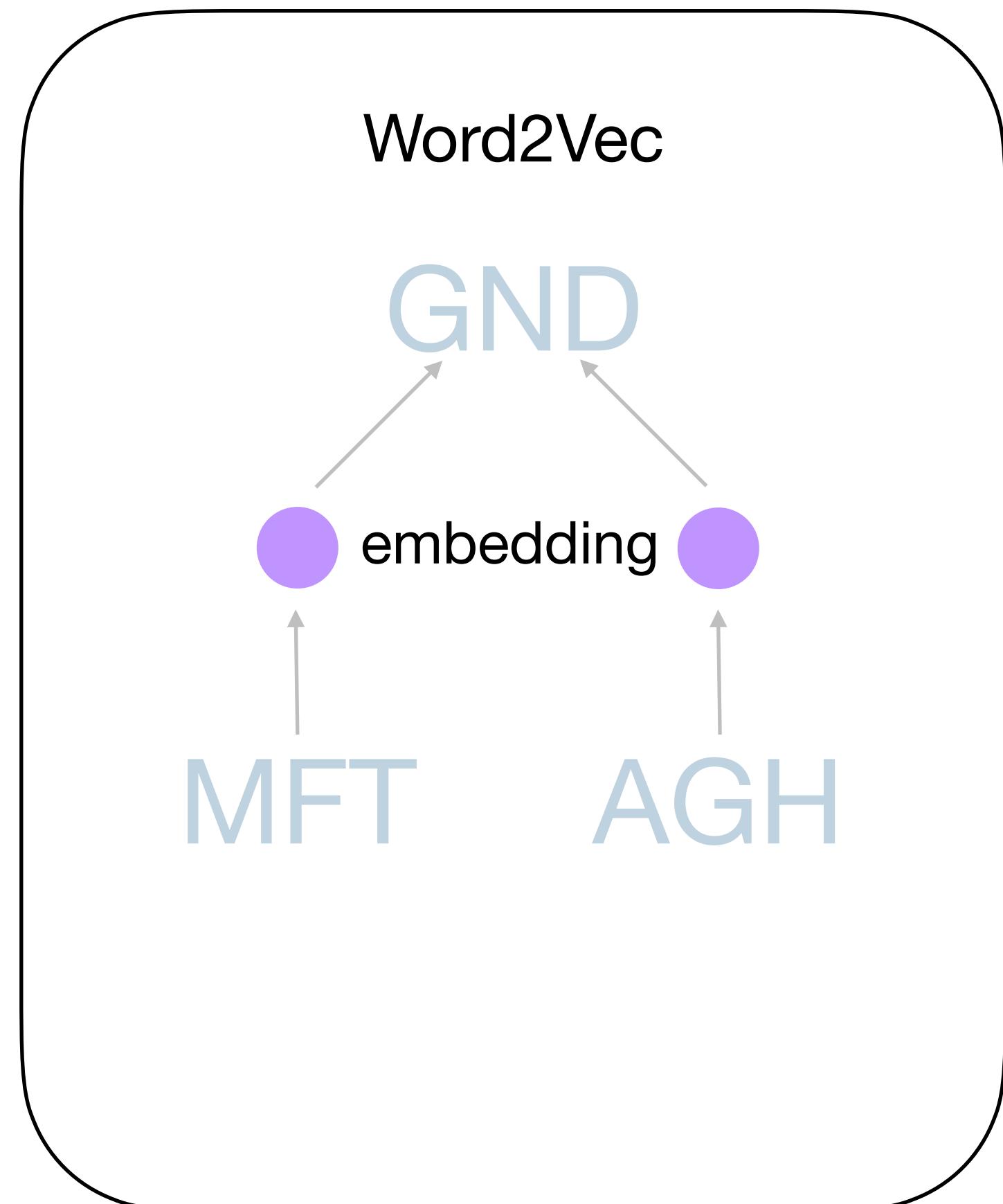
MFTGNDAGH

# Many methods pretend proteins are language



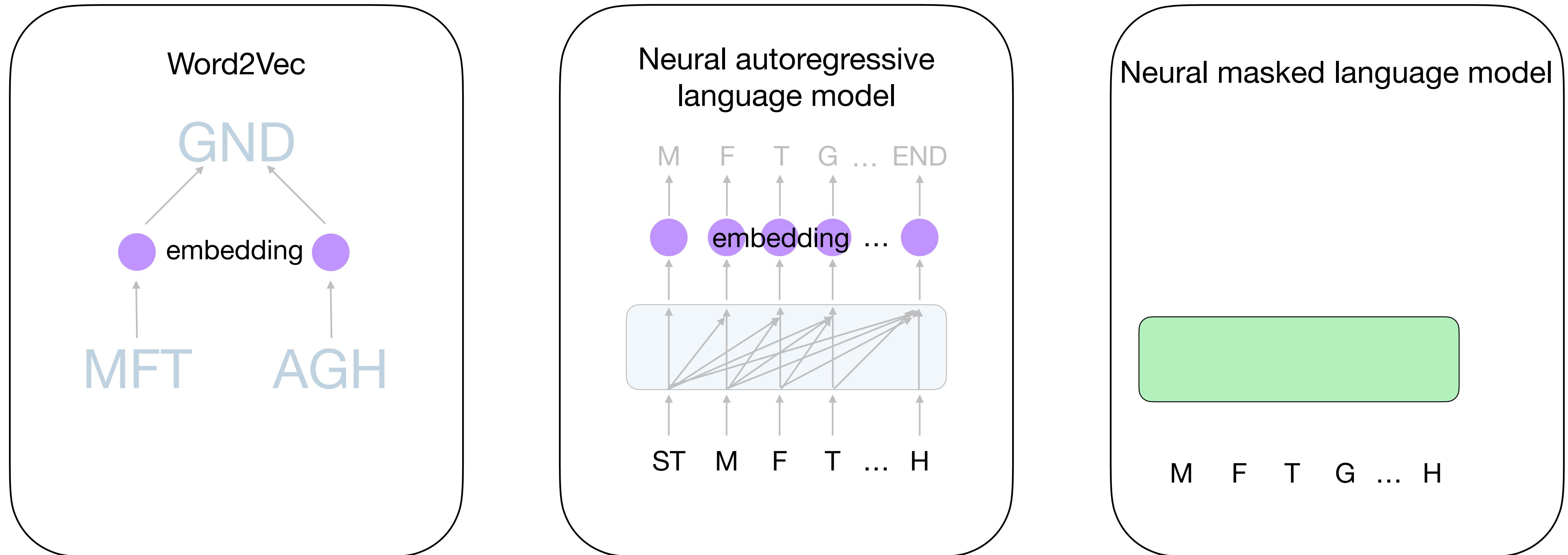
MFTGNDAGH

# Many methods pretend proteins are language



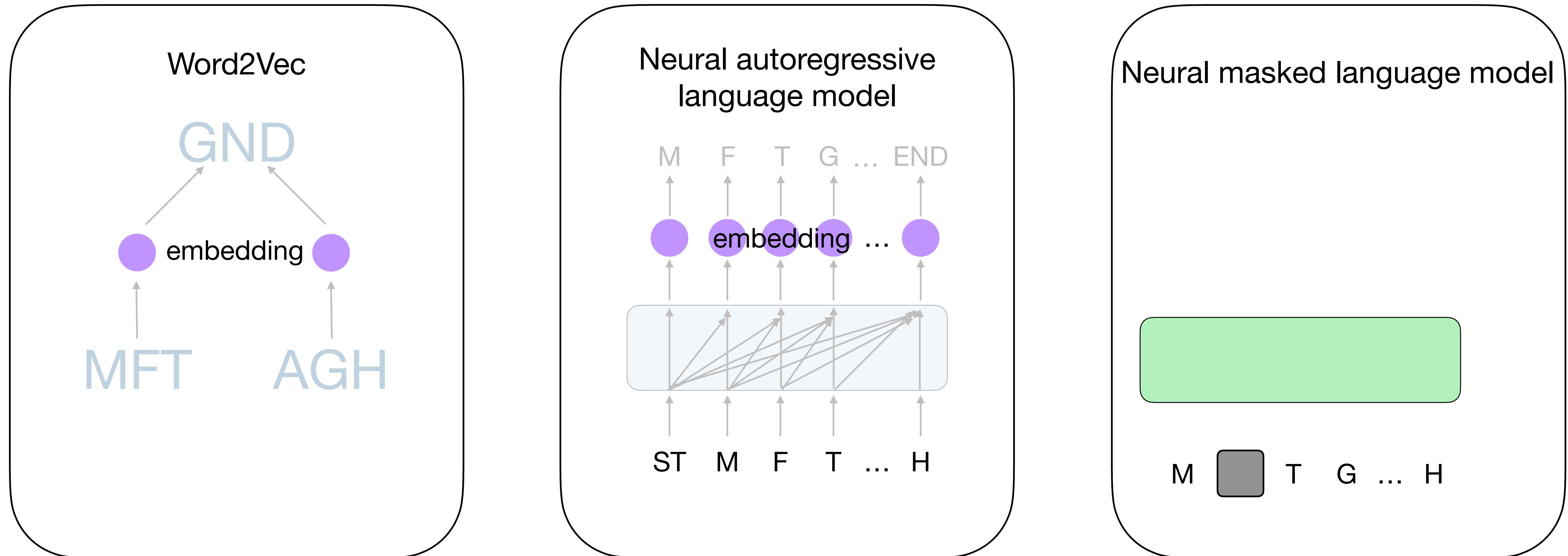
MFTGNDAGH

# Many methods pretend proteins are language

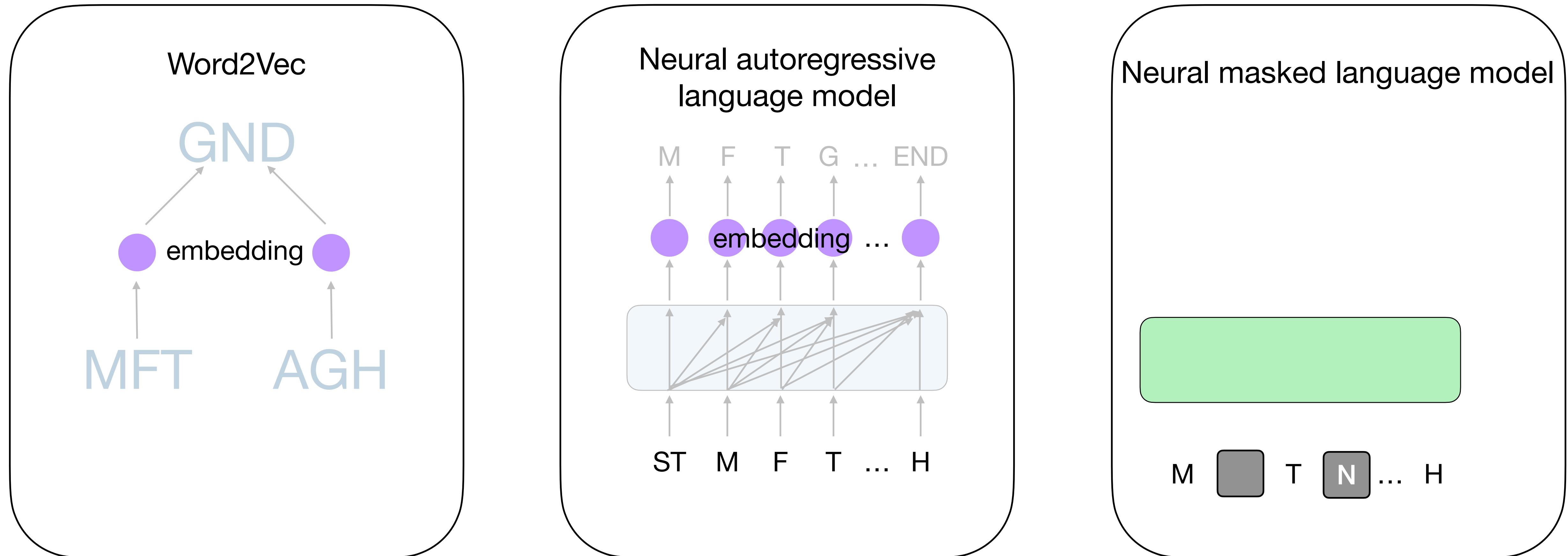


MFTGNDAGH

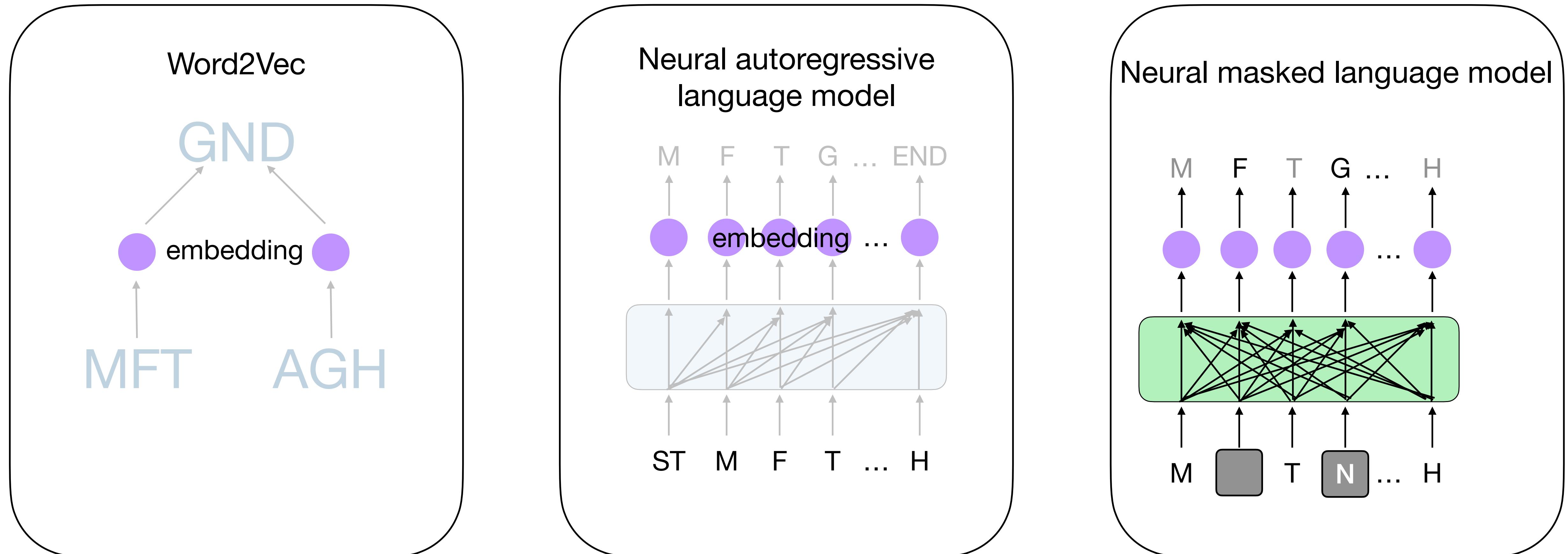
# Many methods pretend proteins are language



# Many methods pretend proteins are language

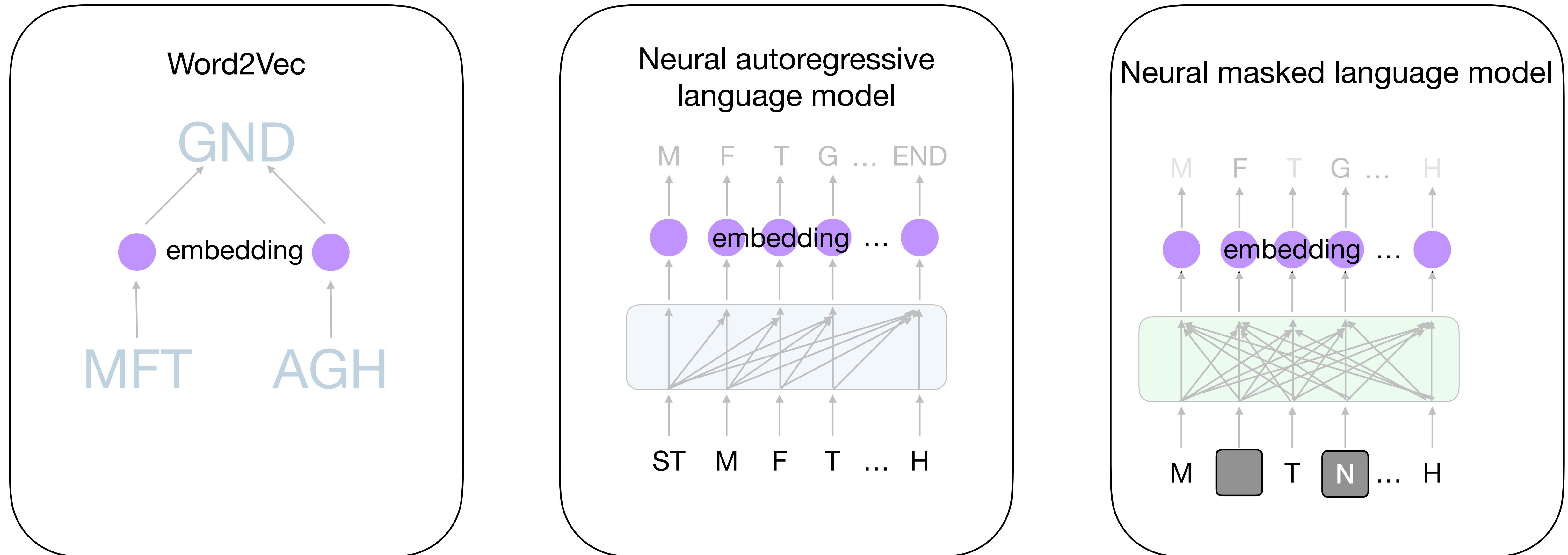


# Many methods pretend proteins are language



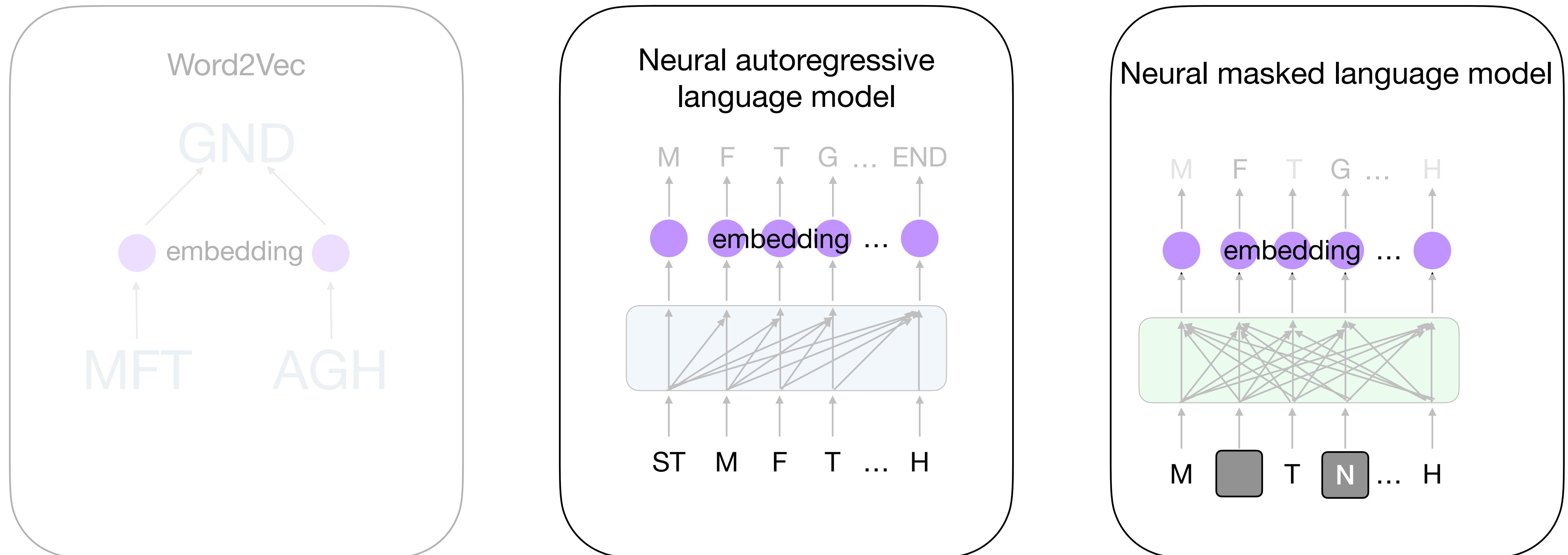
MFTGNDAGH

# Many methods pretend proteins are language



MFTGNDAGH

# Many methods pretend proteins are language

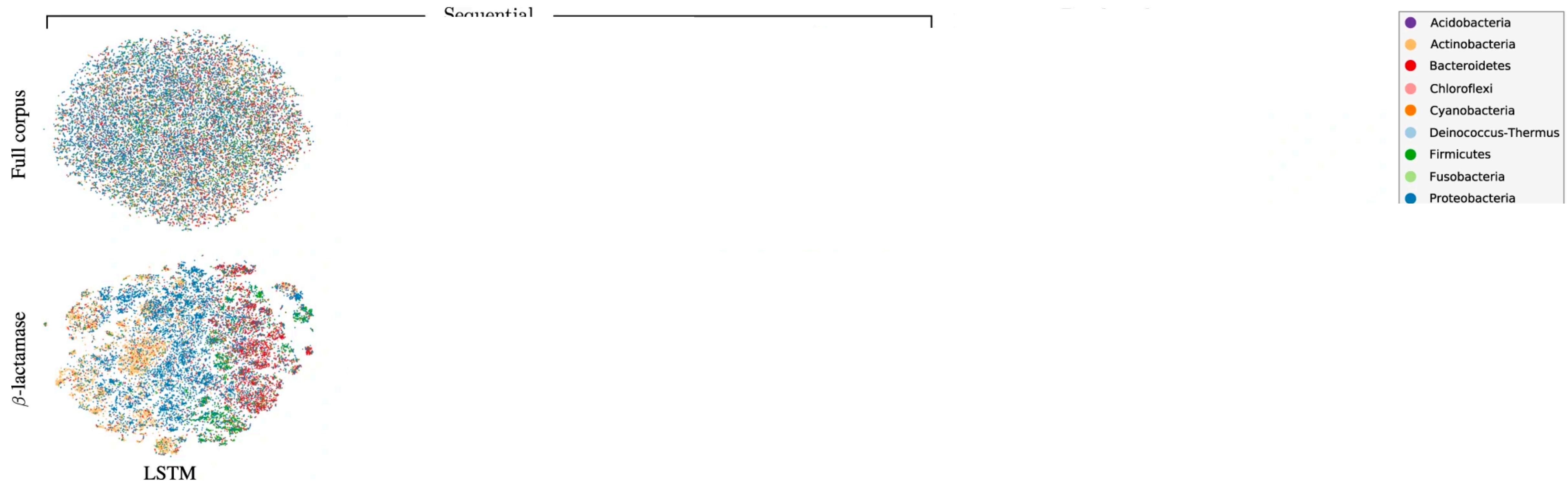


MFTGNDAGH

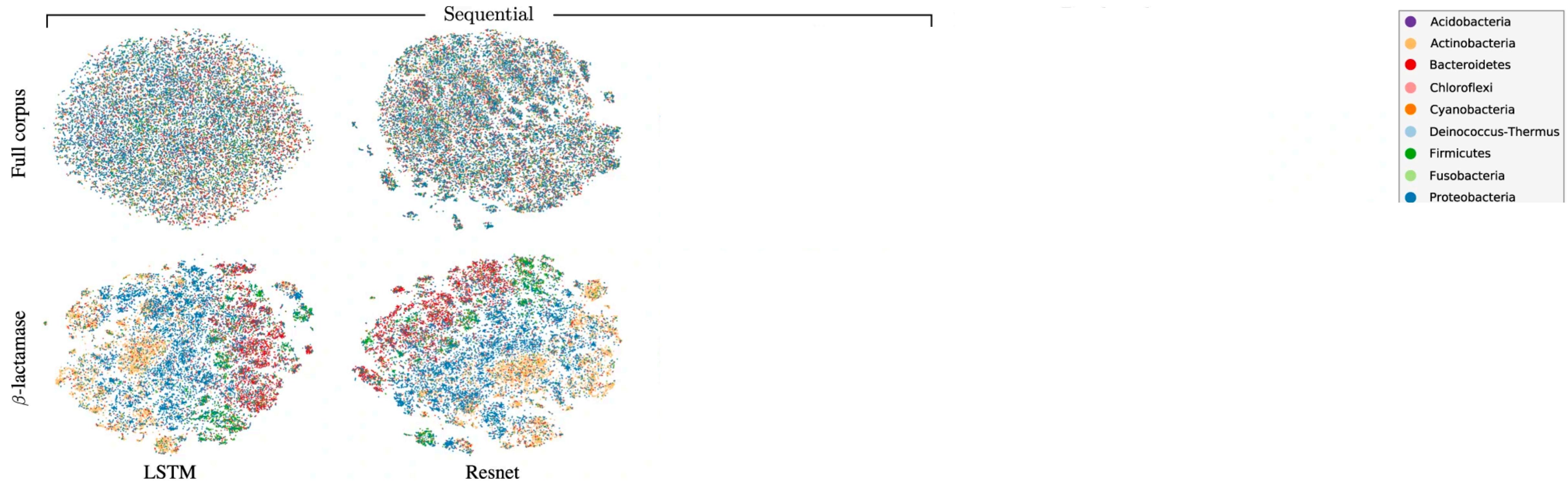
# Architecture, dataset, and task all influence representation quality



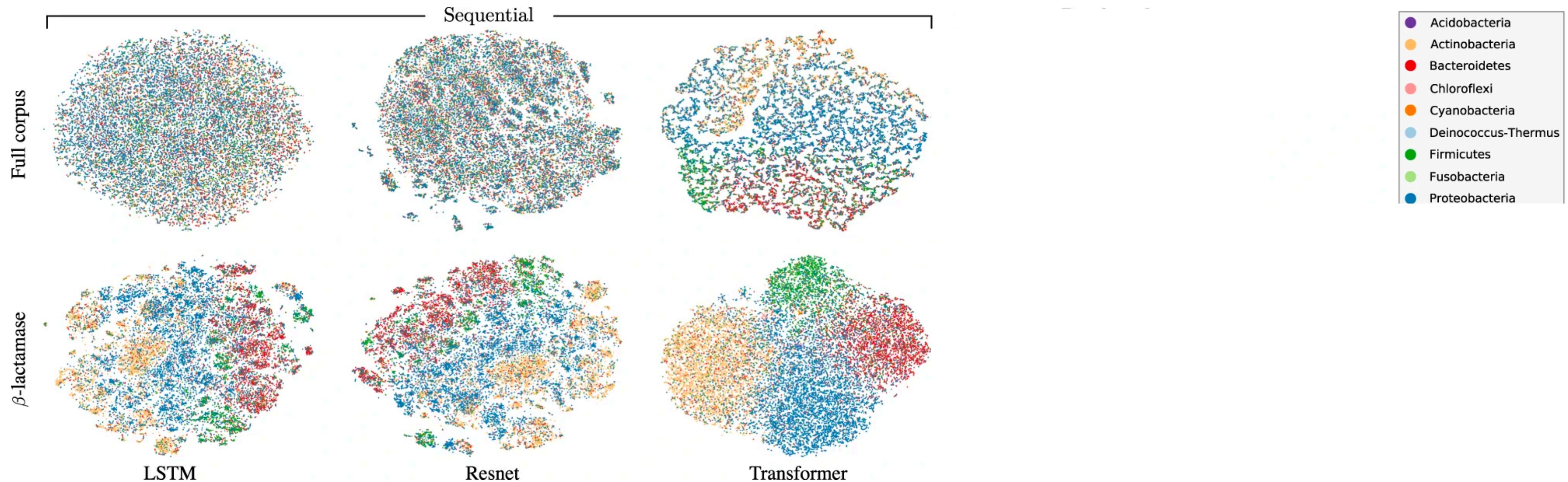
# Architecture, dataset, and task all influence representation quality



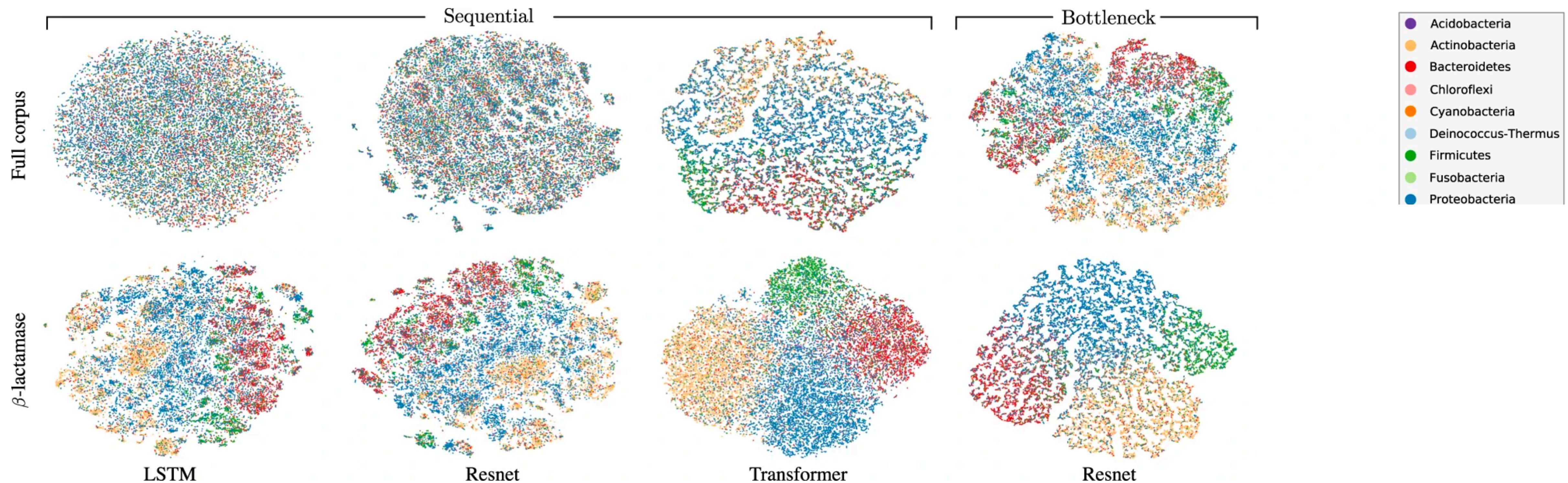
# Architecture, dataset, and task all influence representation quality



# Architecture, dataset, and task all influence representation quality

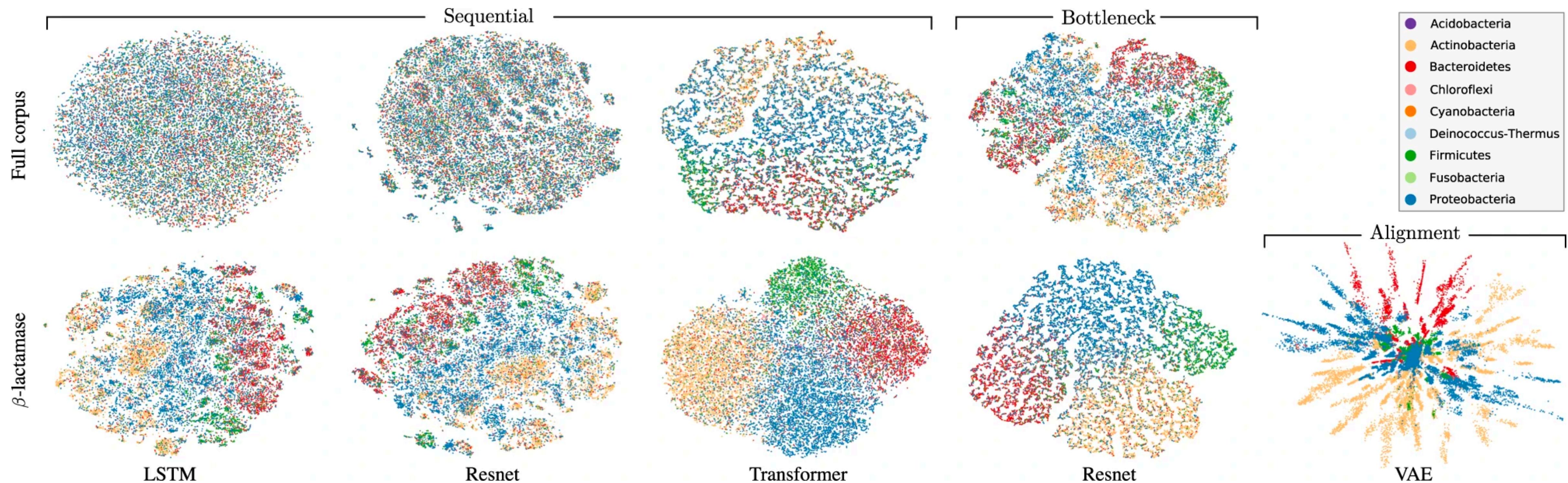


# Architecture, dataset, and task all influence representation quality



Detlefsen et al. 2022

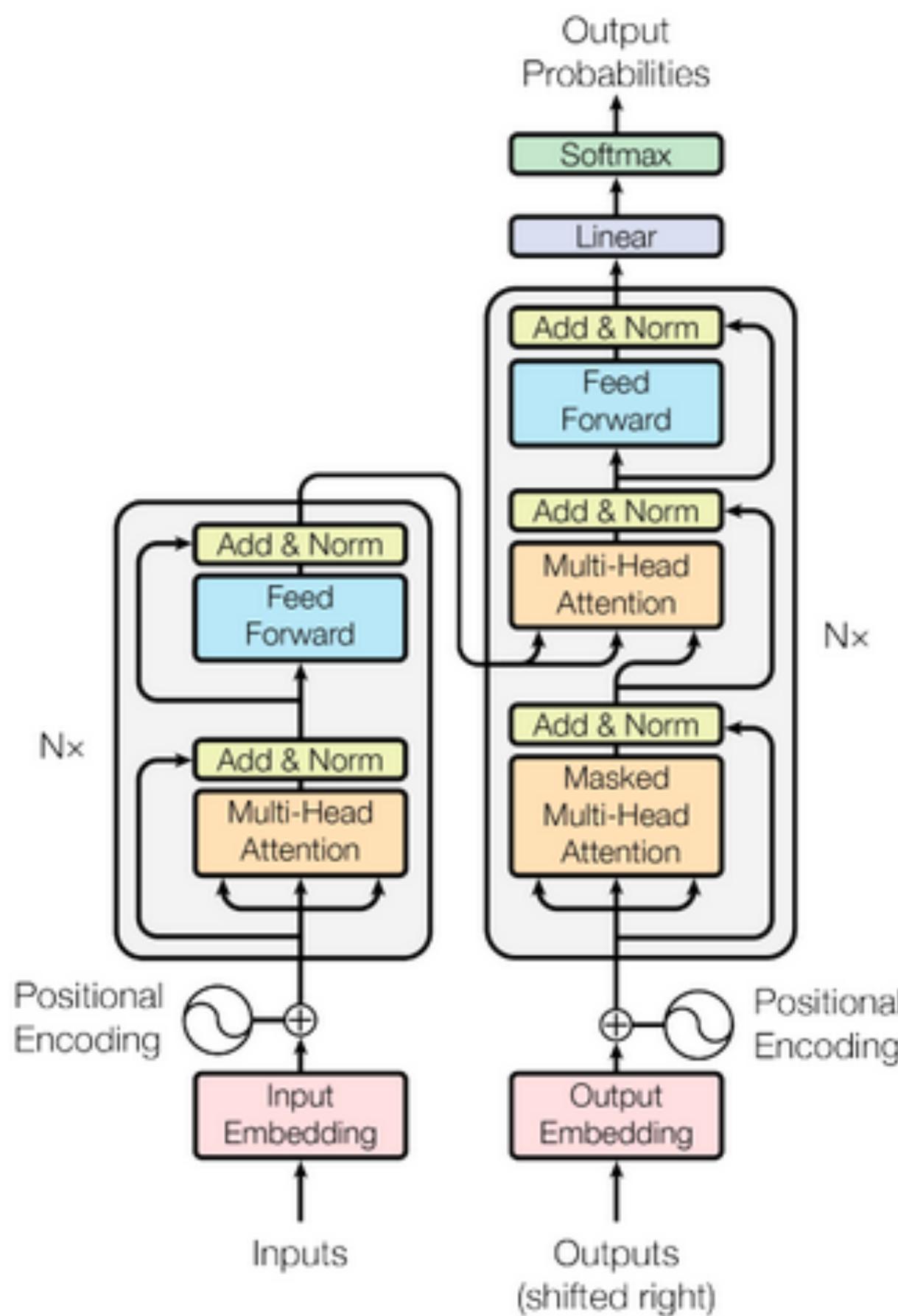
# Architecture, dataset, and task all influence representation quality



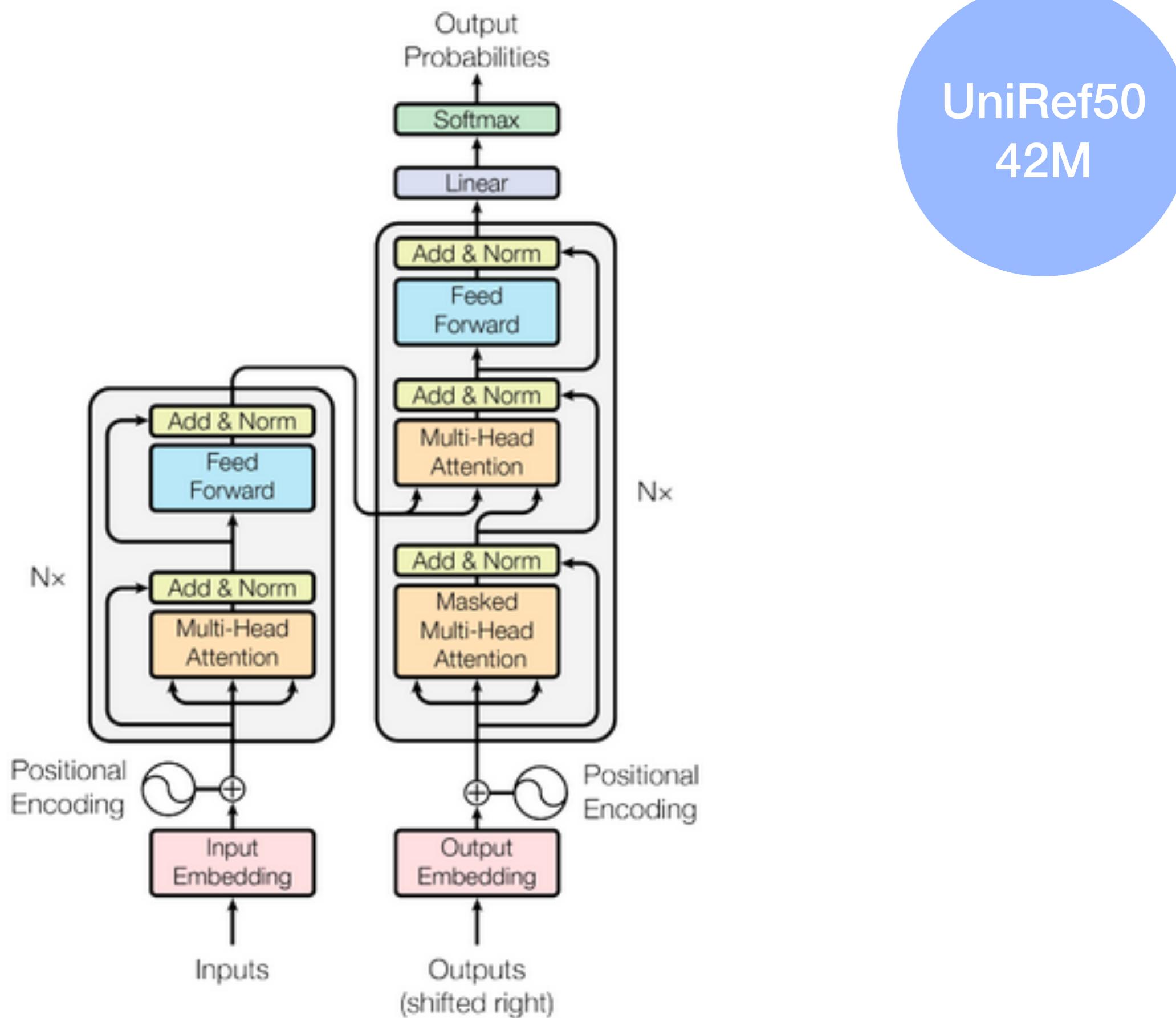
Detlefsen et al. 2022

# Transformer masked language models are currently popular

# Transformer masked language models are currently popular

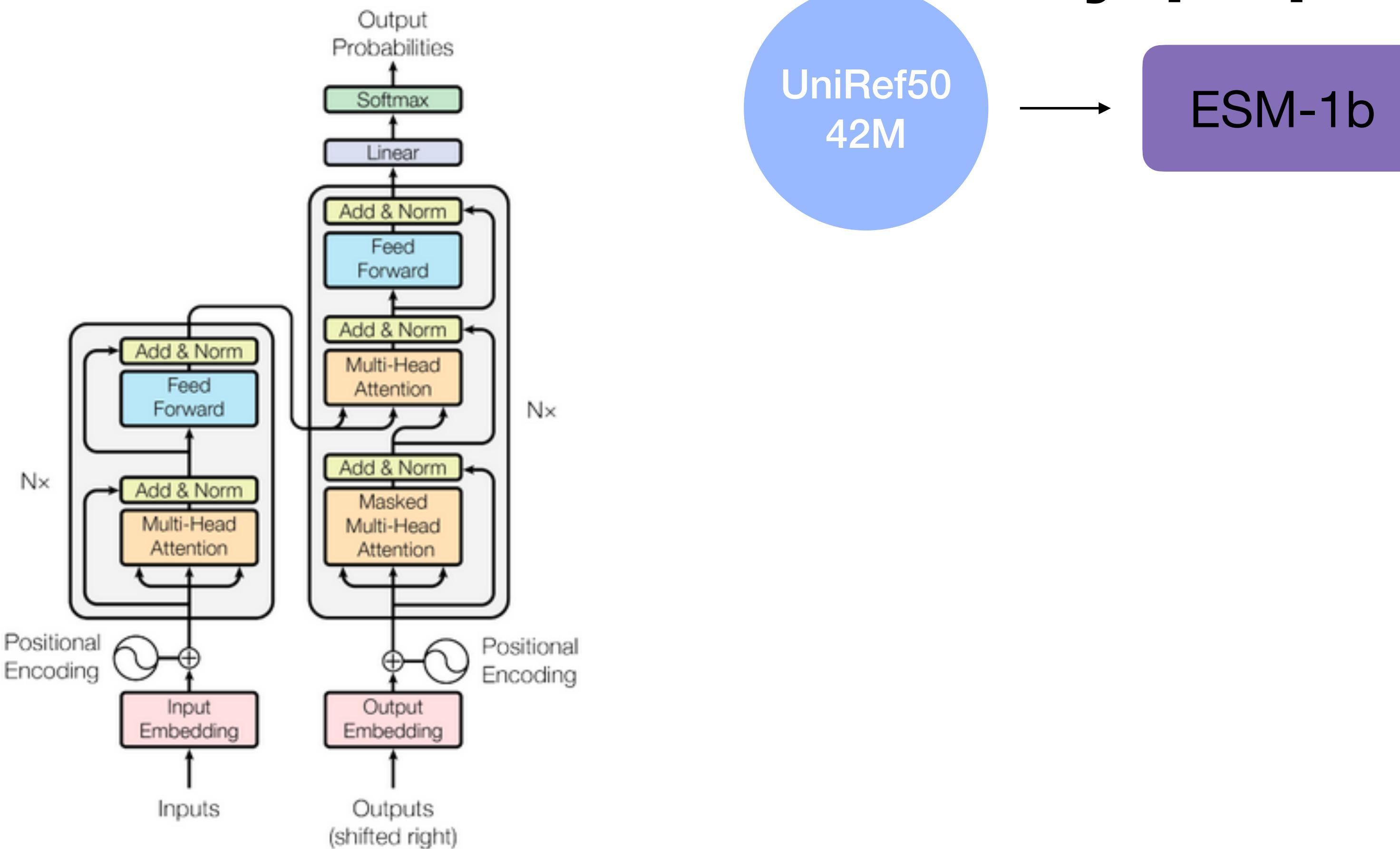


# Transformer masked language models are currently popular

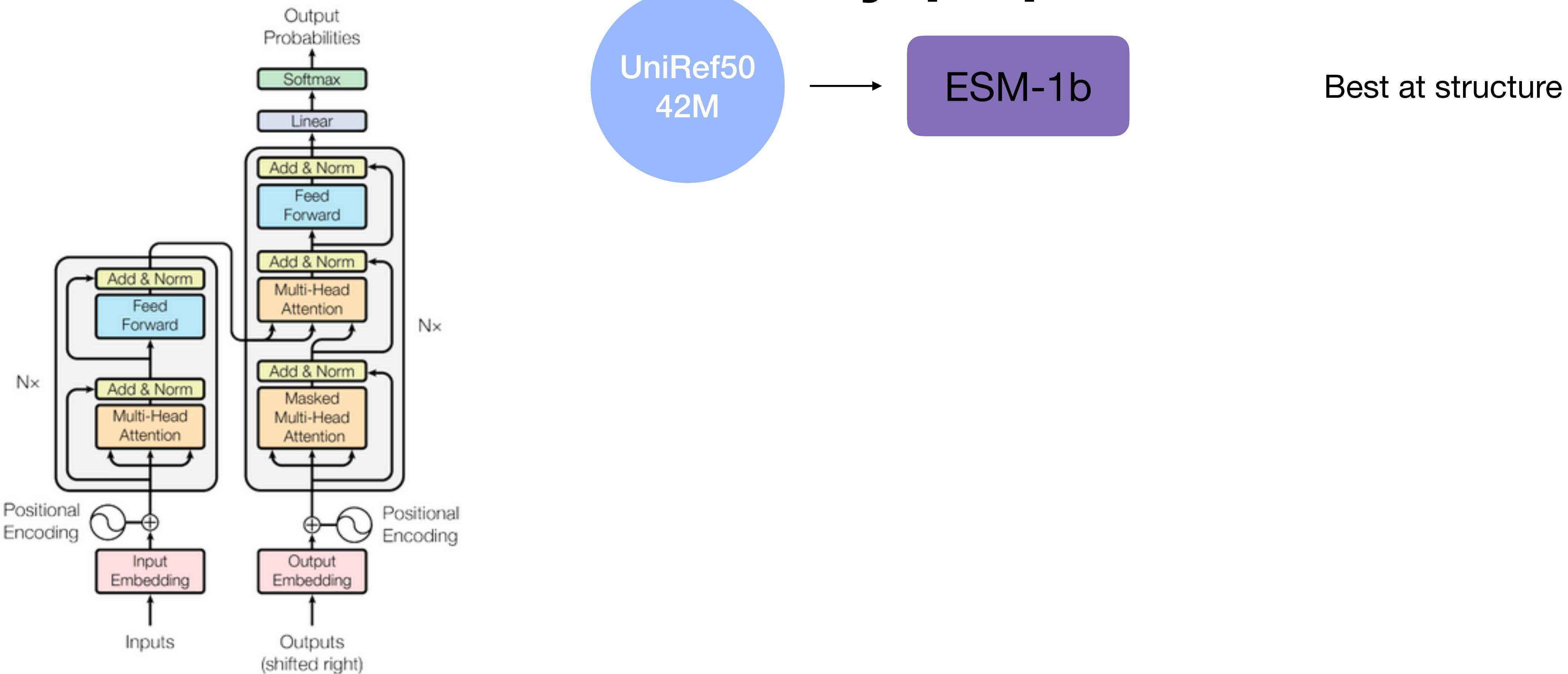


UniRef50  
42M

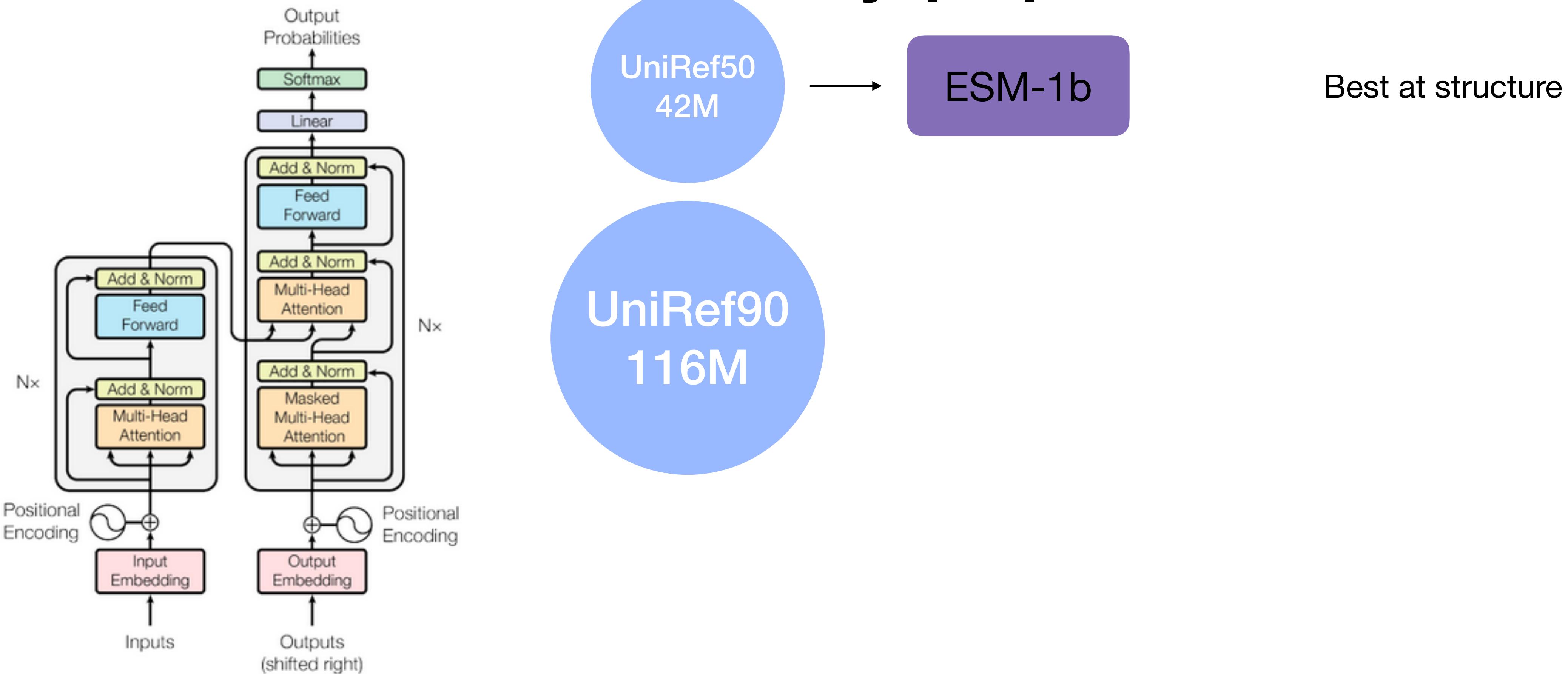
# Transformer masked language models are currently popular



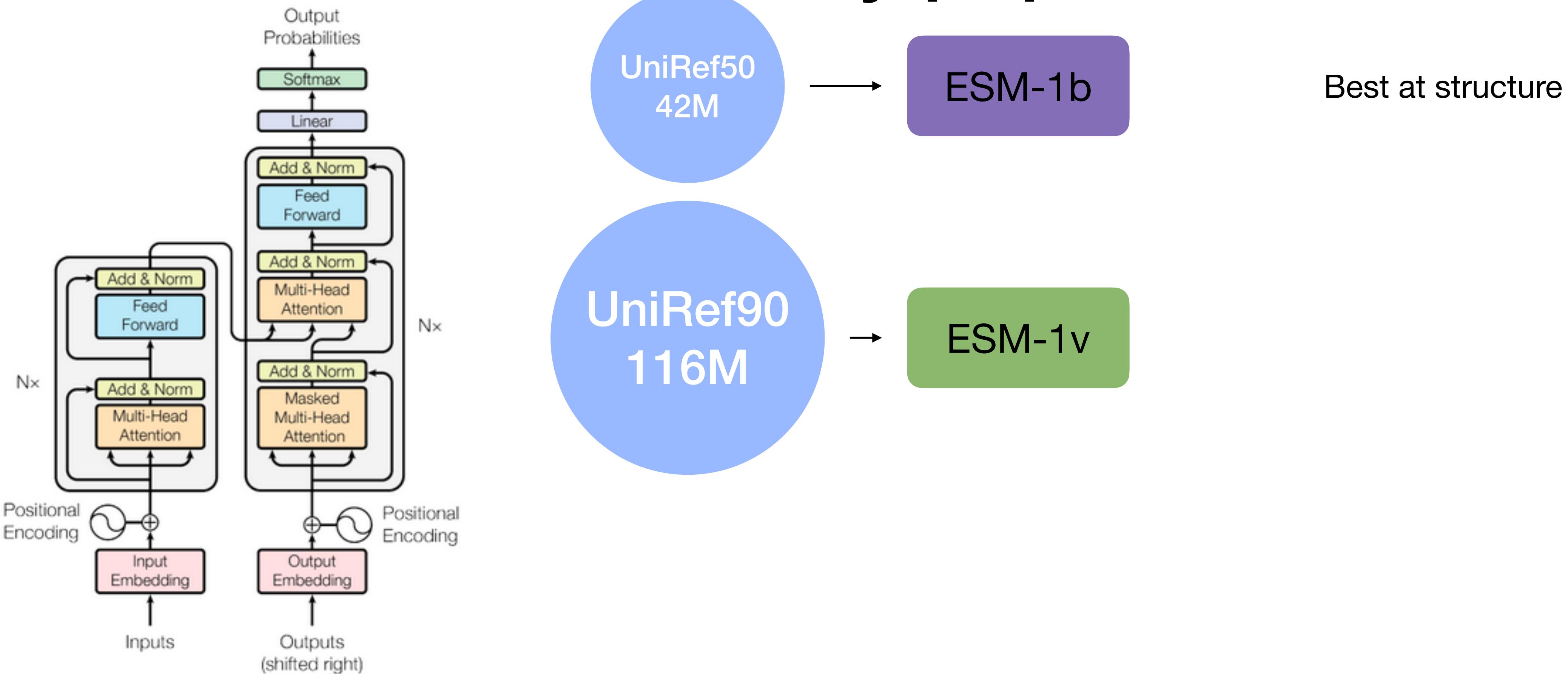
# Transformer masked language models are currently popular



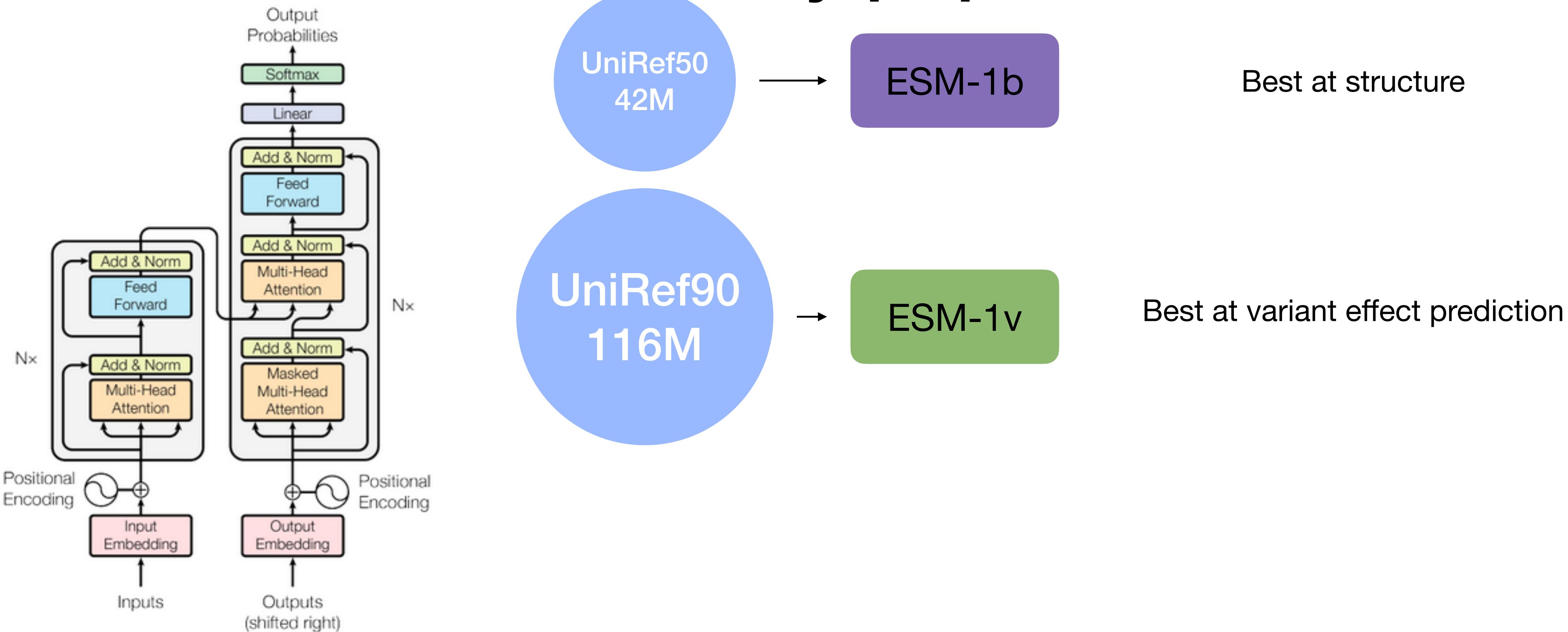
# Transformer masked language models are currently popular



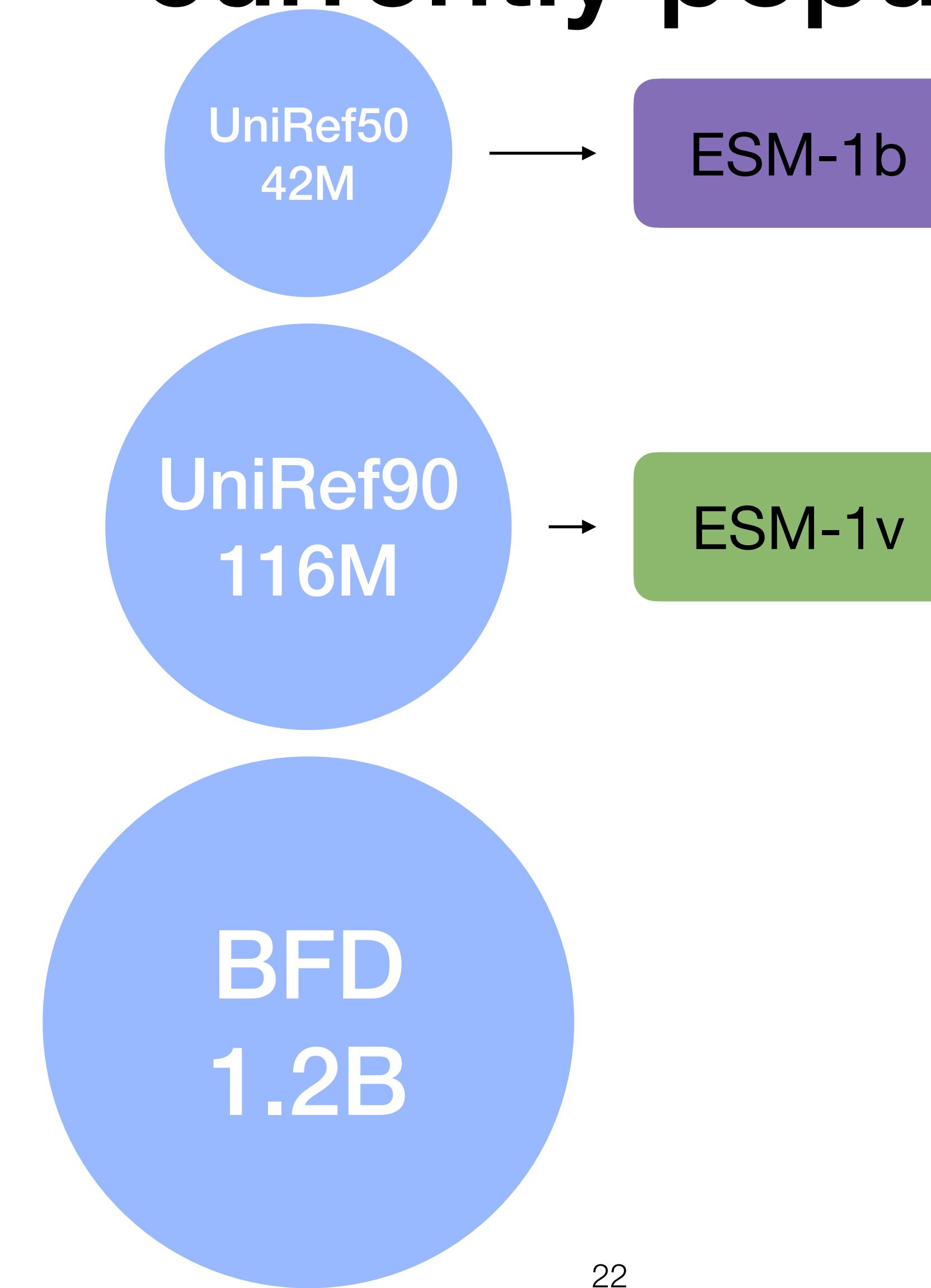
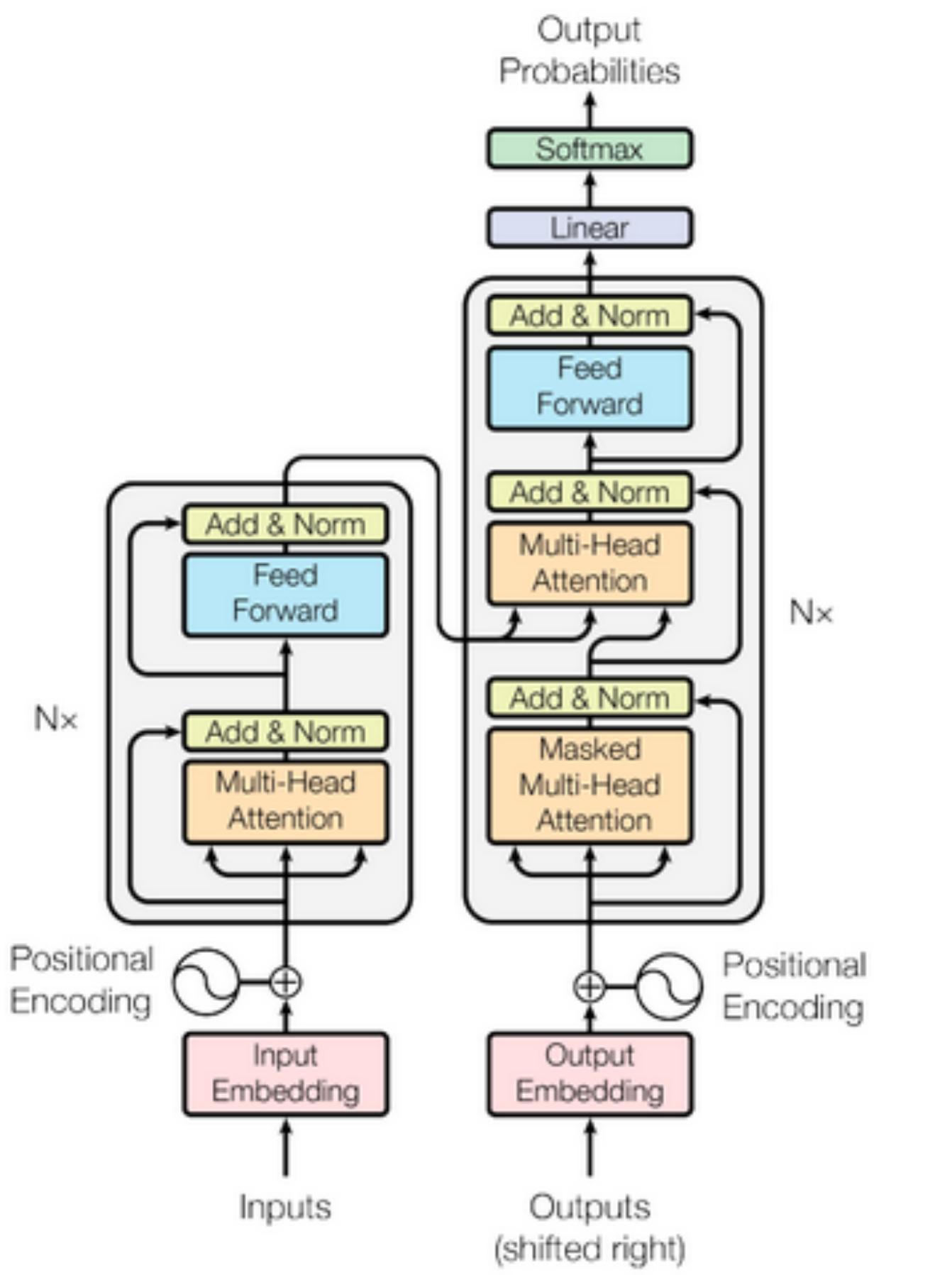
# Transformer masked language models are currently popular



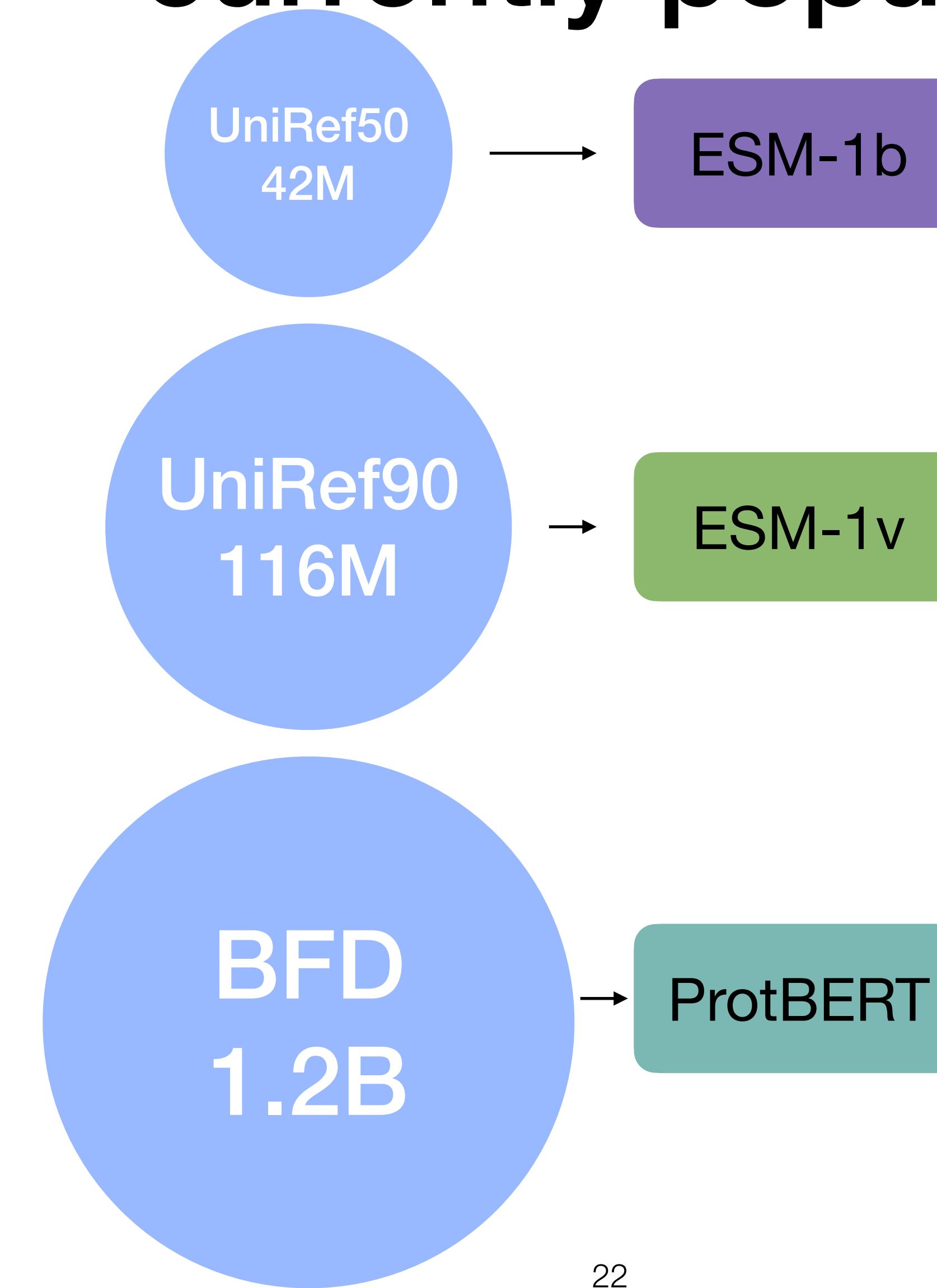
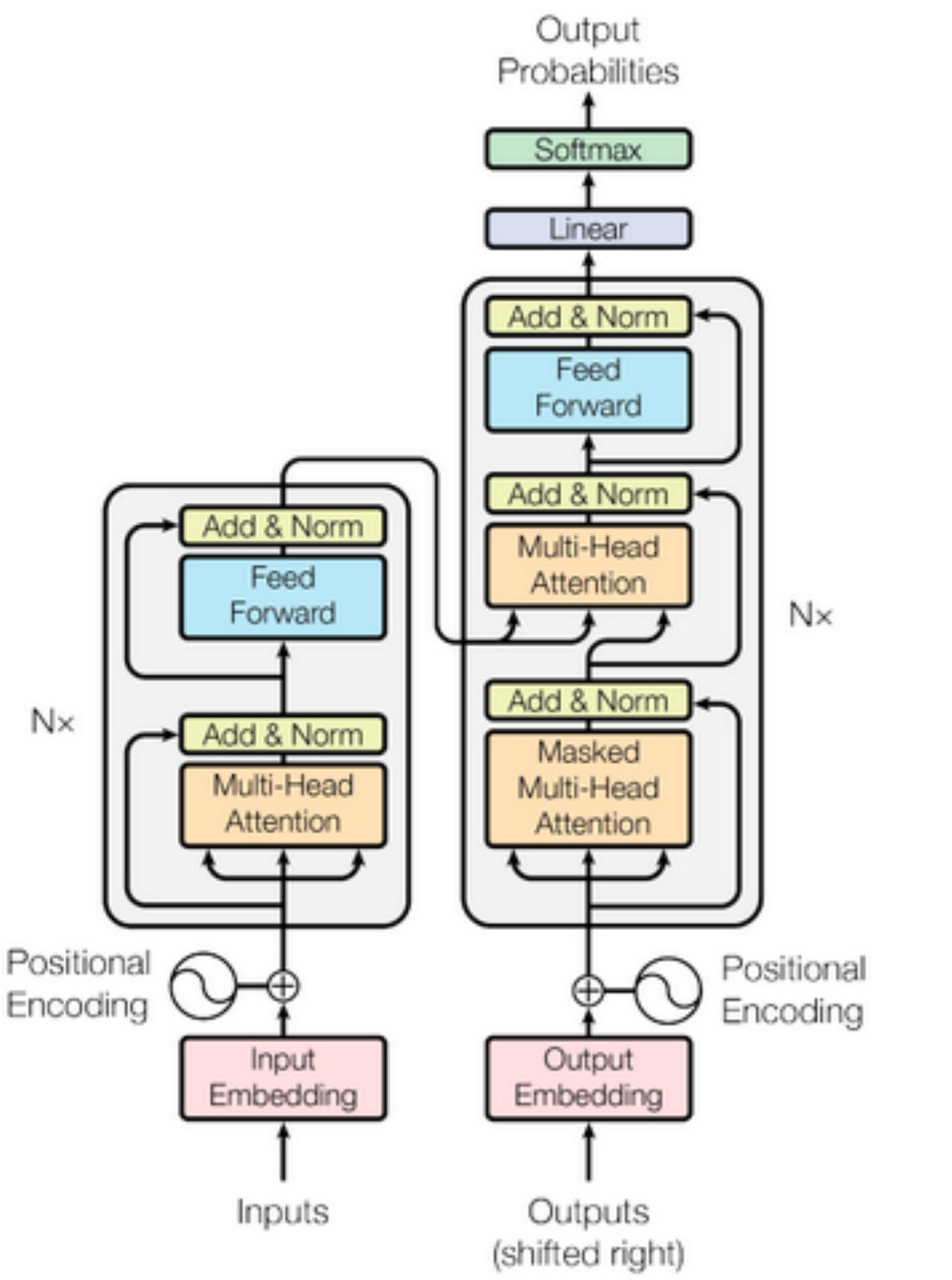
# Transformer masked language models are currently popular



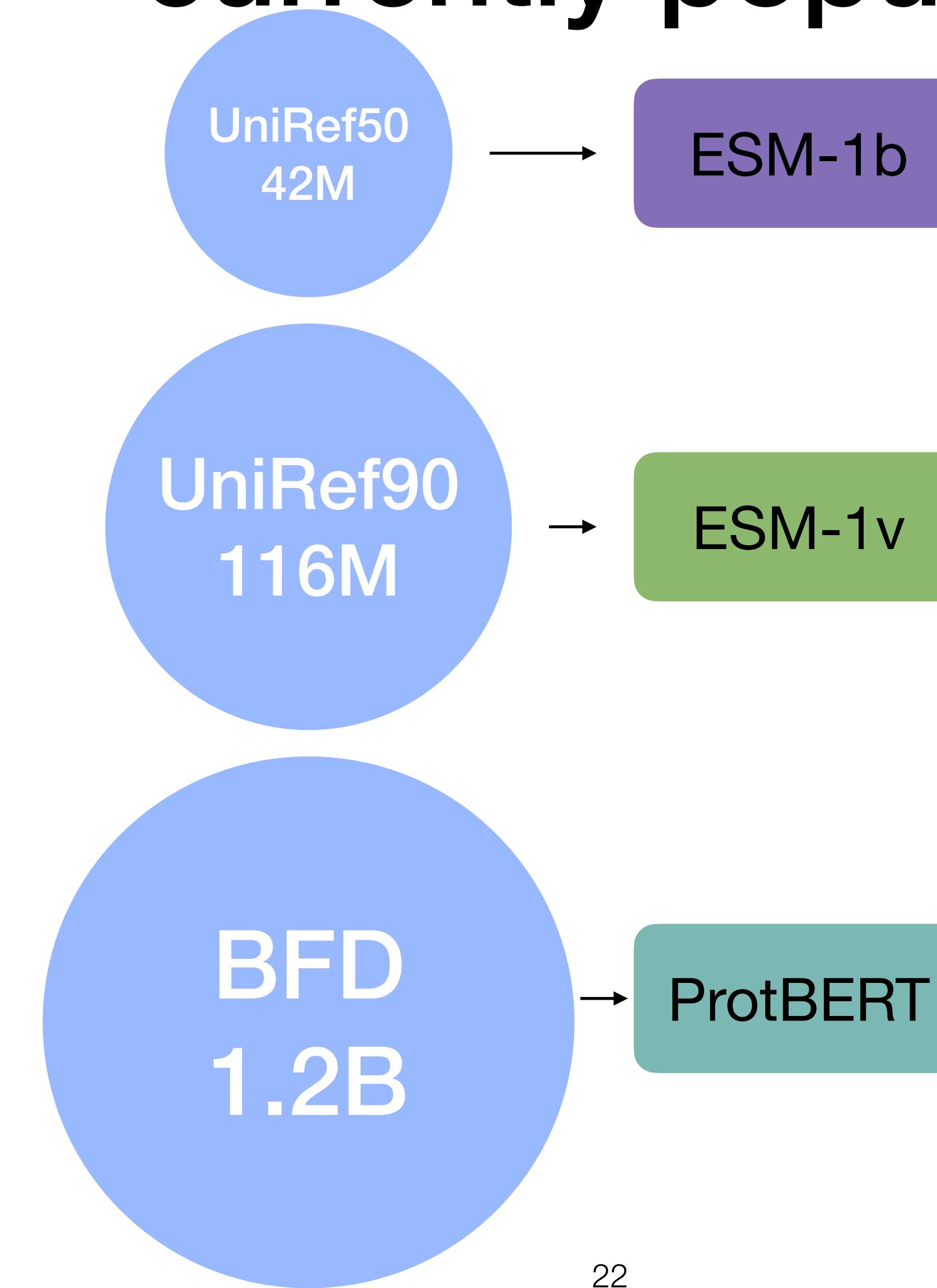
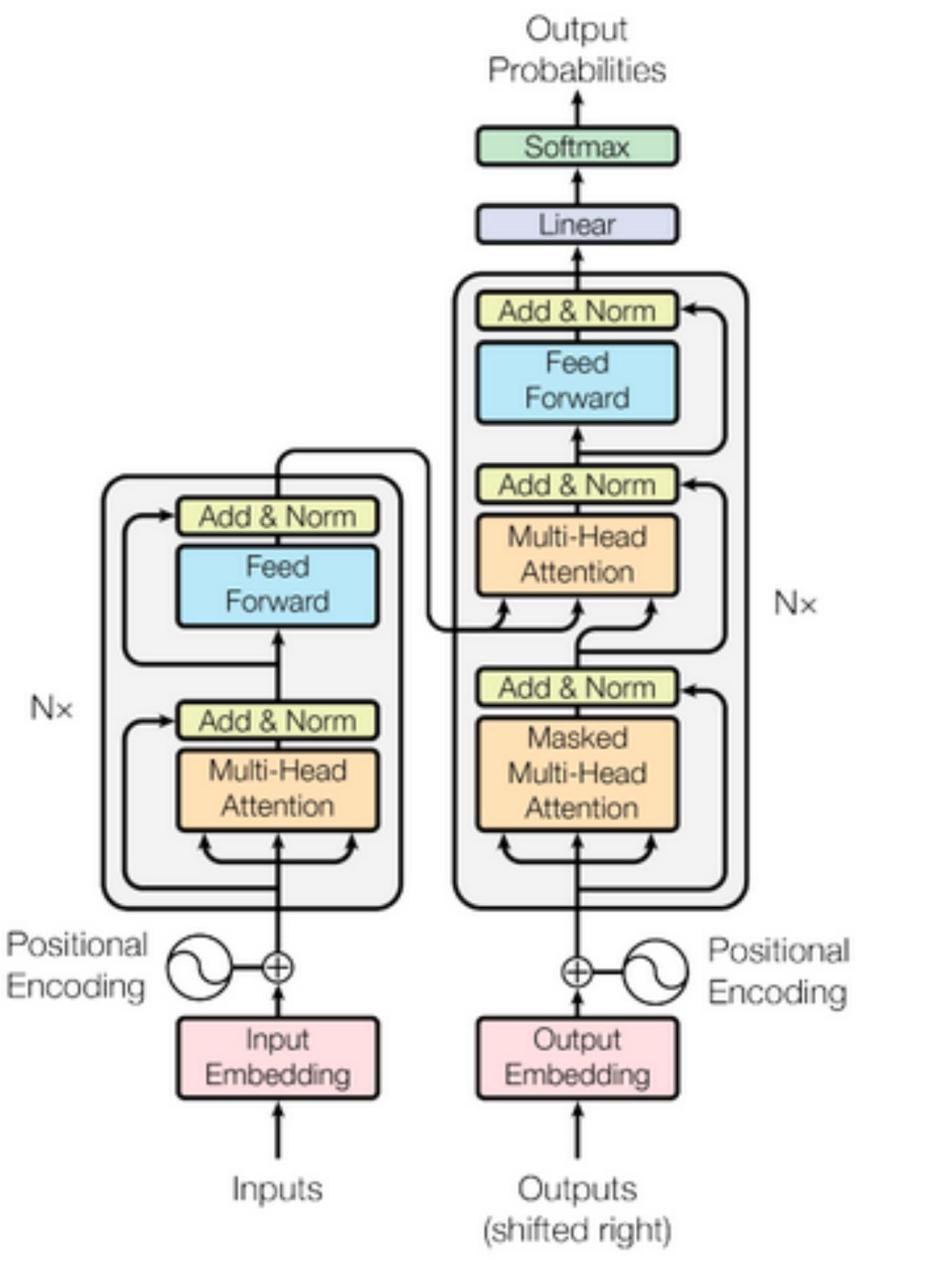
# Transformer masked language models are currently popular



# Transformer masked language models are currently popular

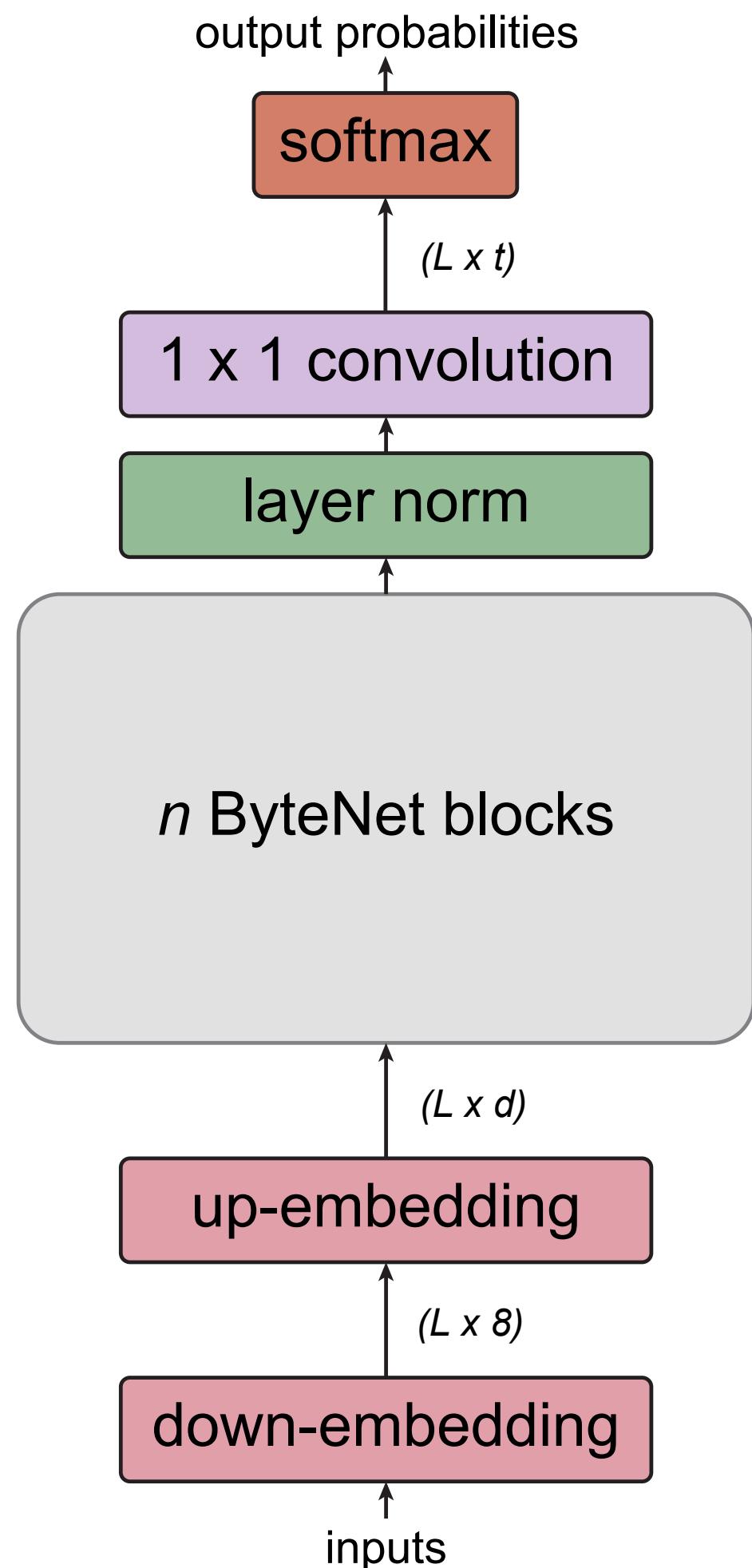


# Transformer masked language models are currently popular

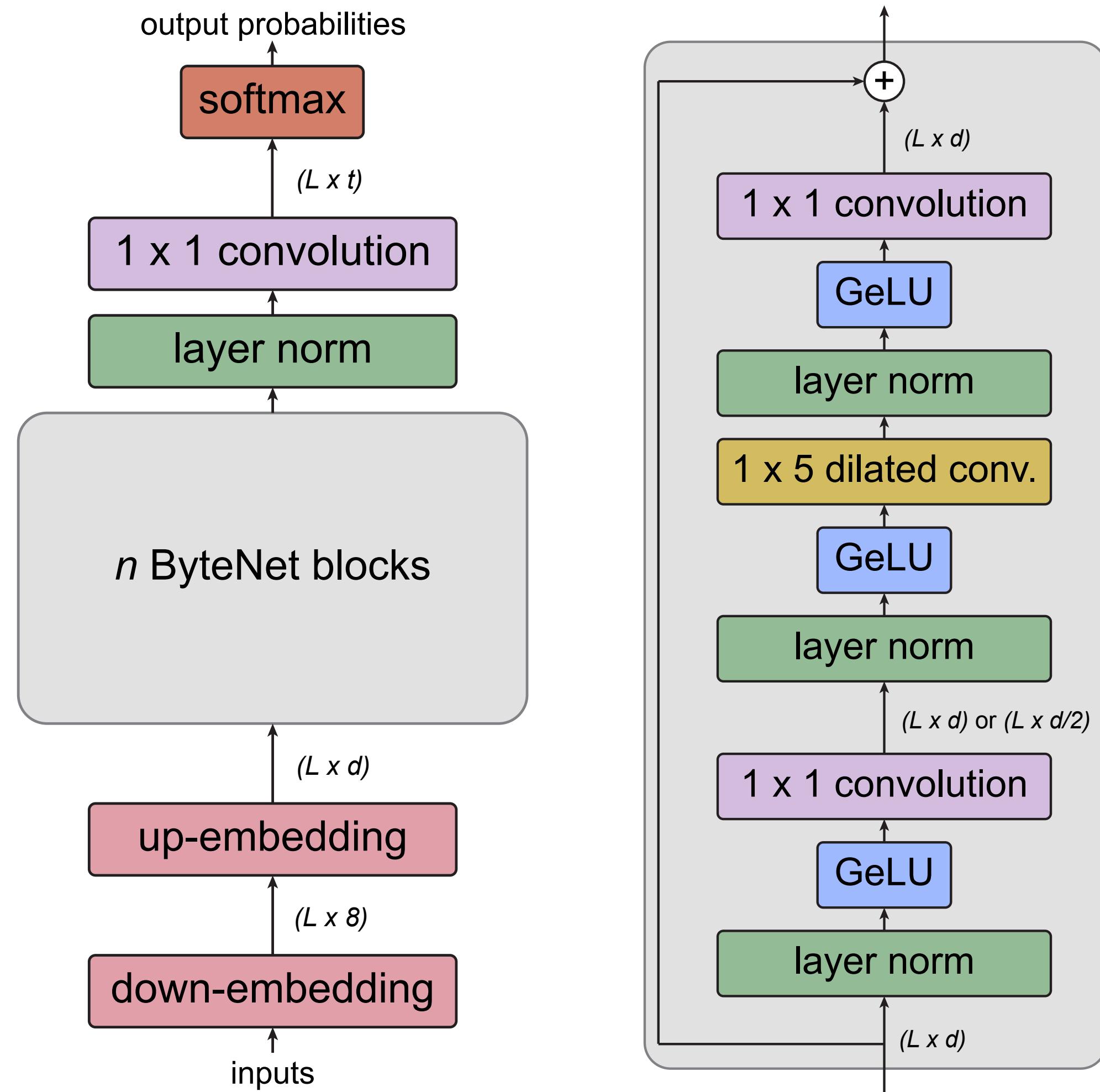


# CNNs are competitive with transformers for pretraining

# CNNs are competitive with transformers for pretraining

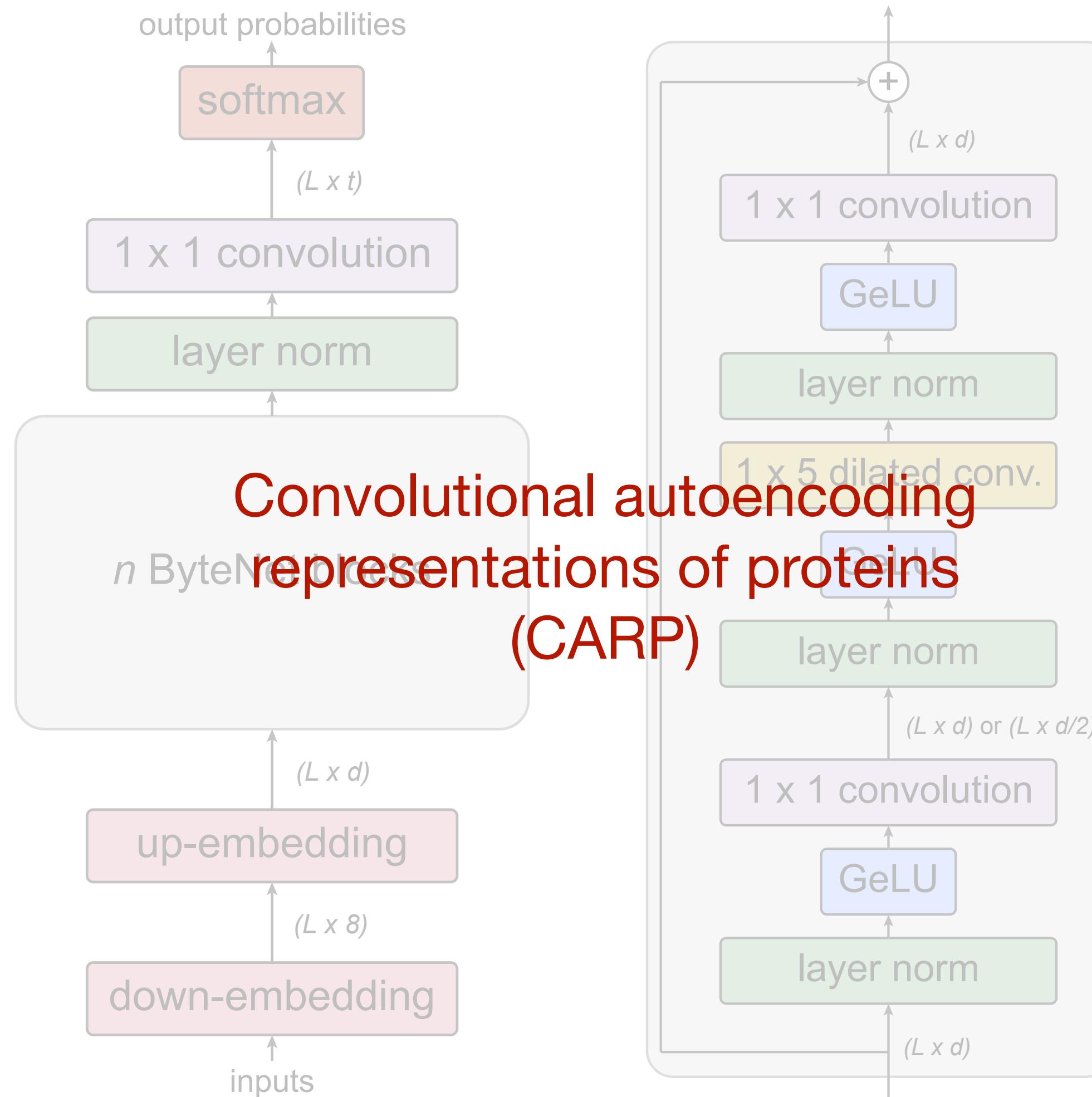


# CNNs are competitive with transformers for pretraining

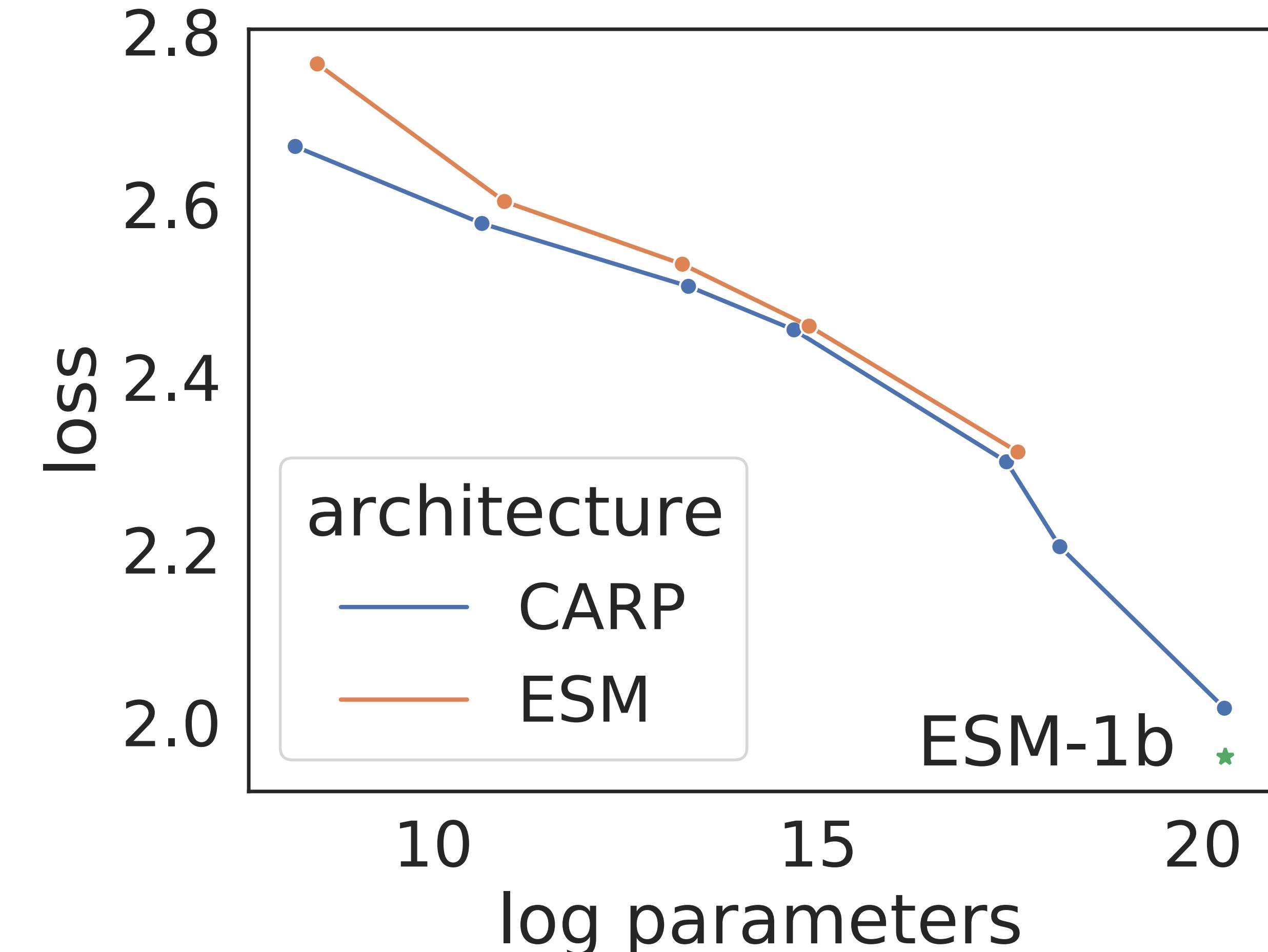
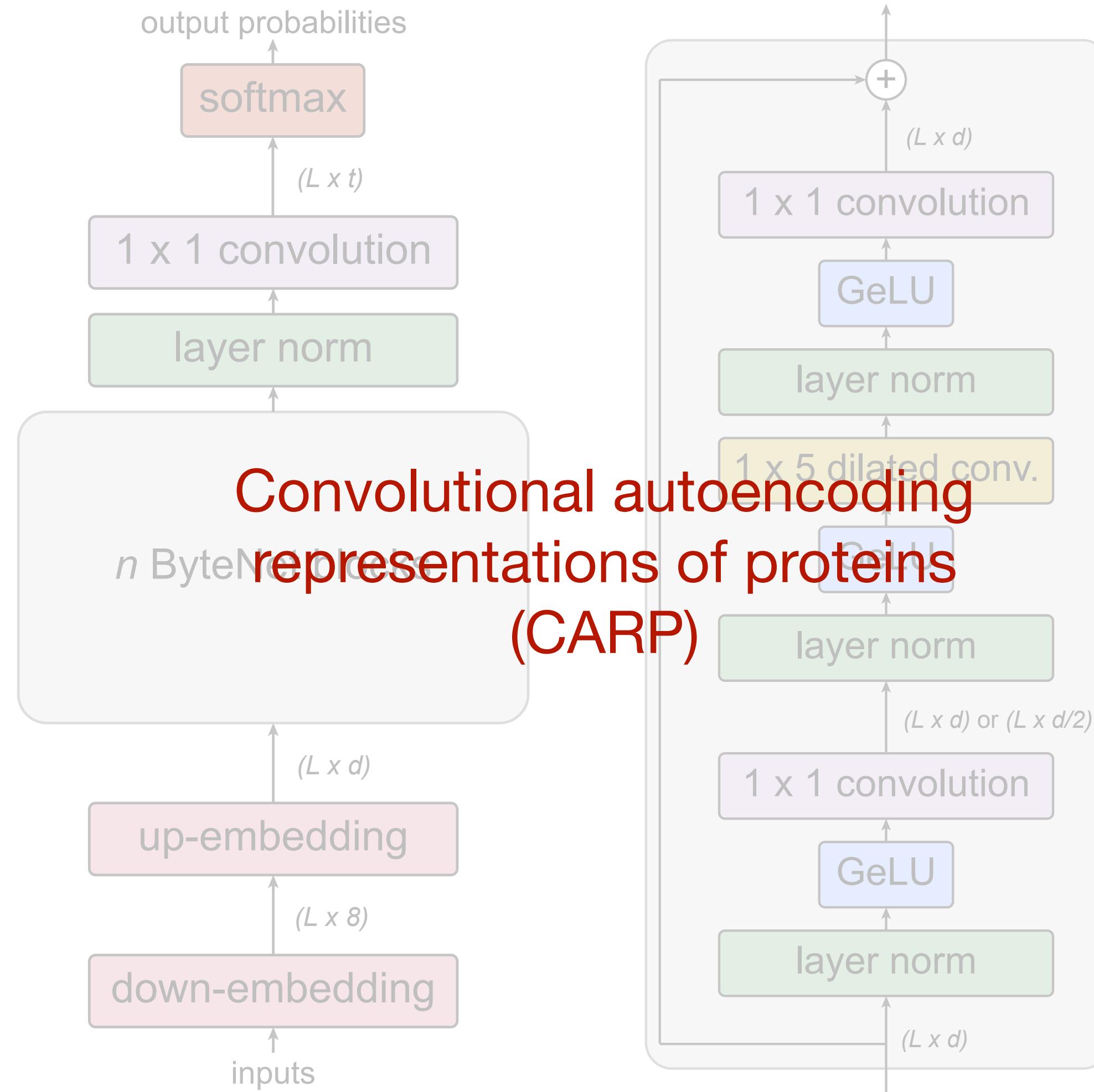


# CNNs are competitive with transformers for pretraining

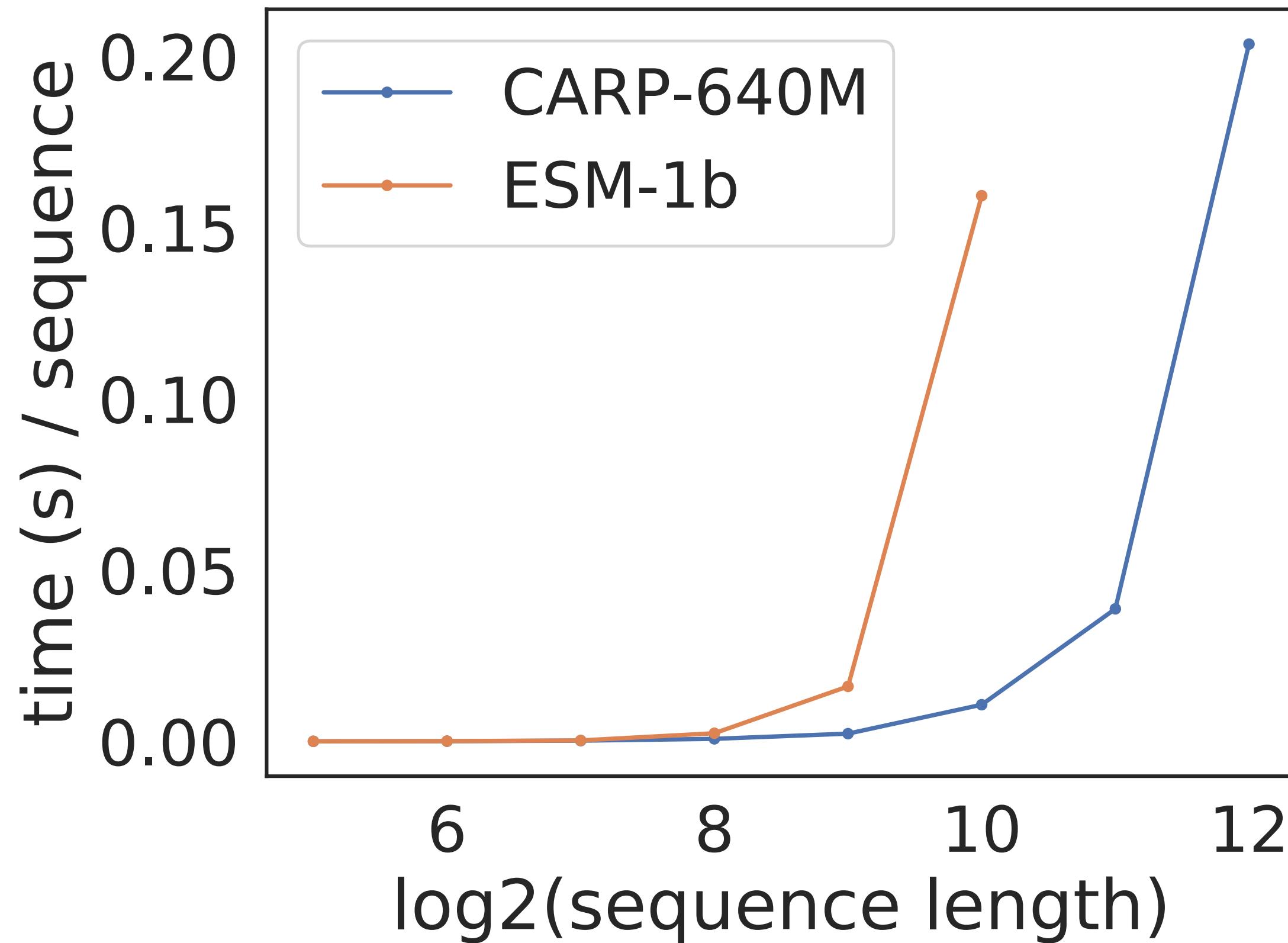
# CNNs are competitive with transformers for pretraining



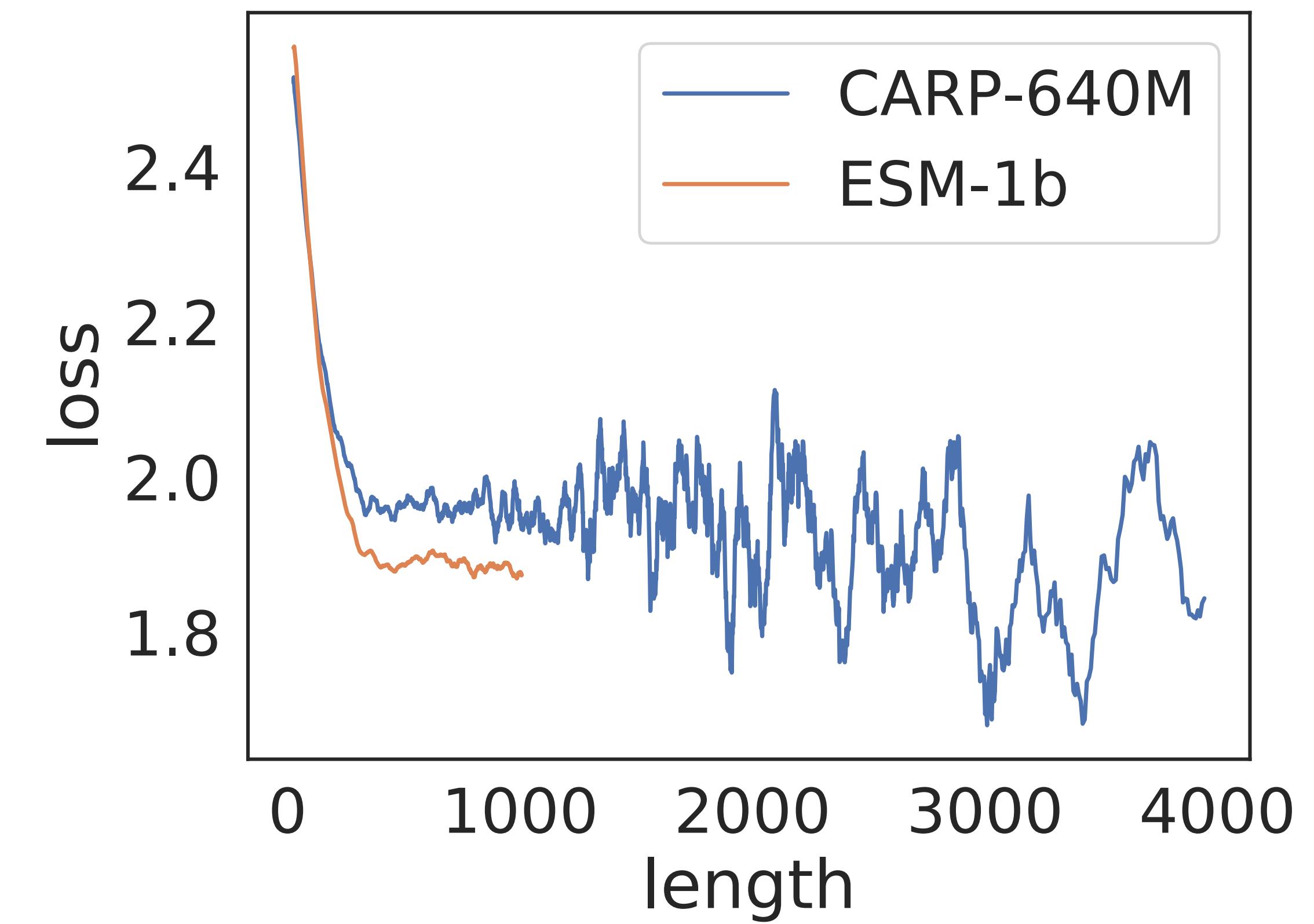
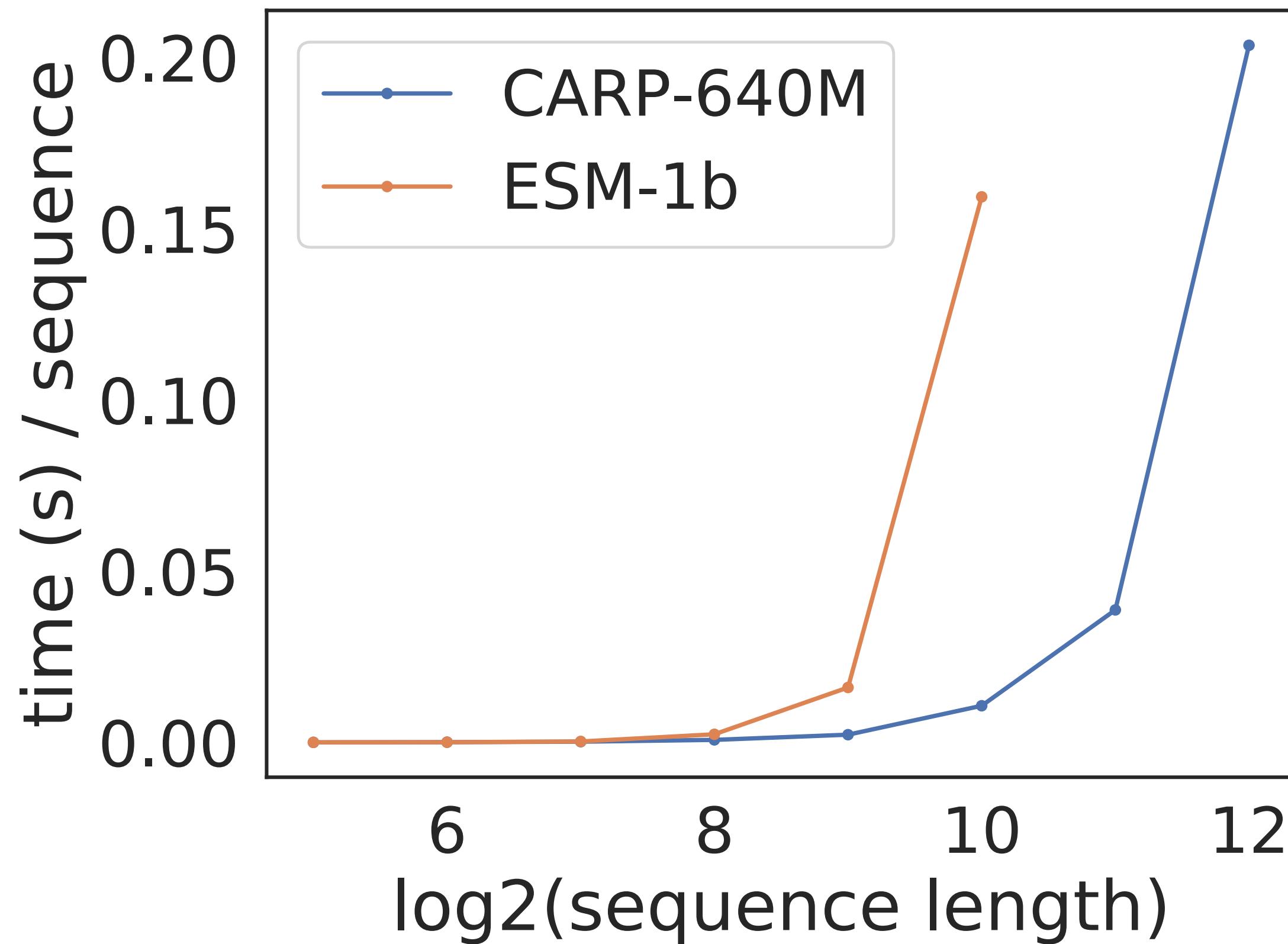
# CNNs are competitive with transformers for pretraining



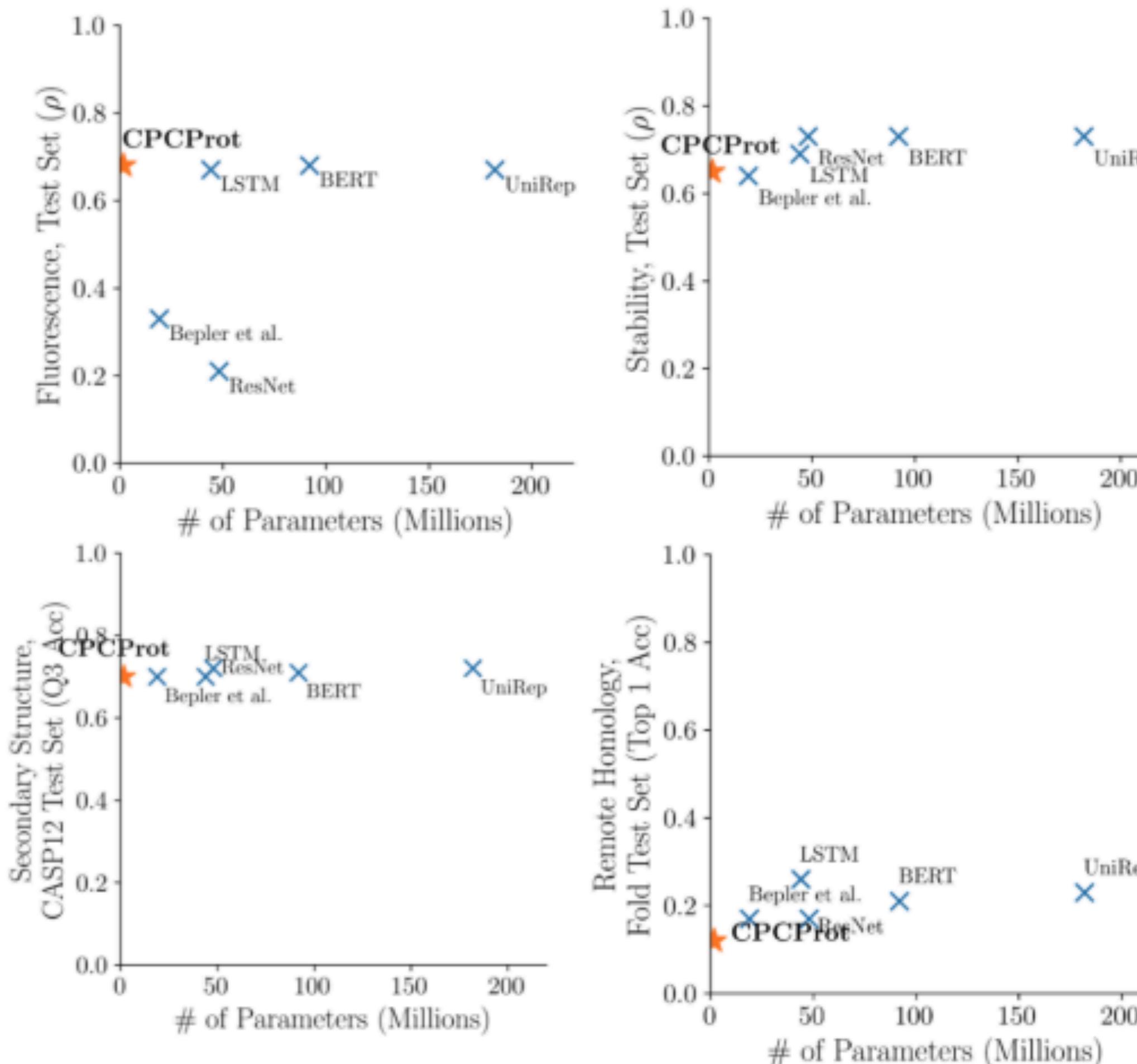
# CNNs scale better with length



# CNNs scale better with length

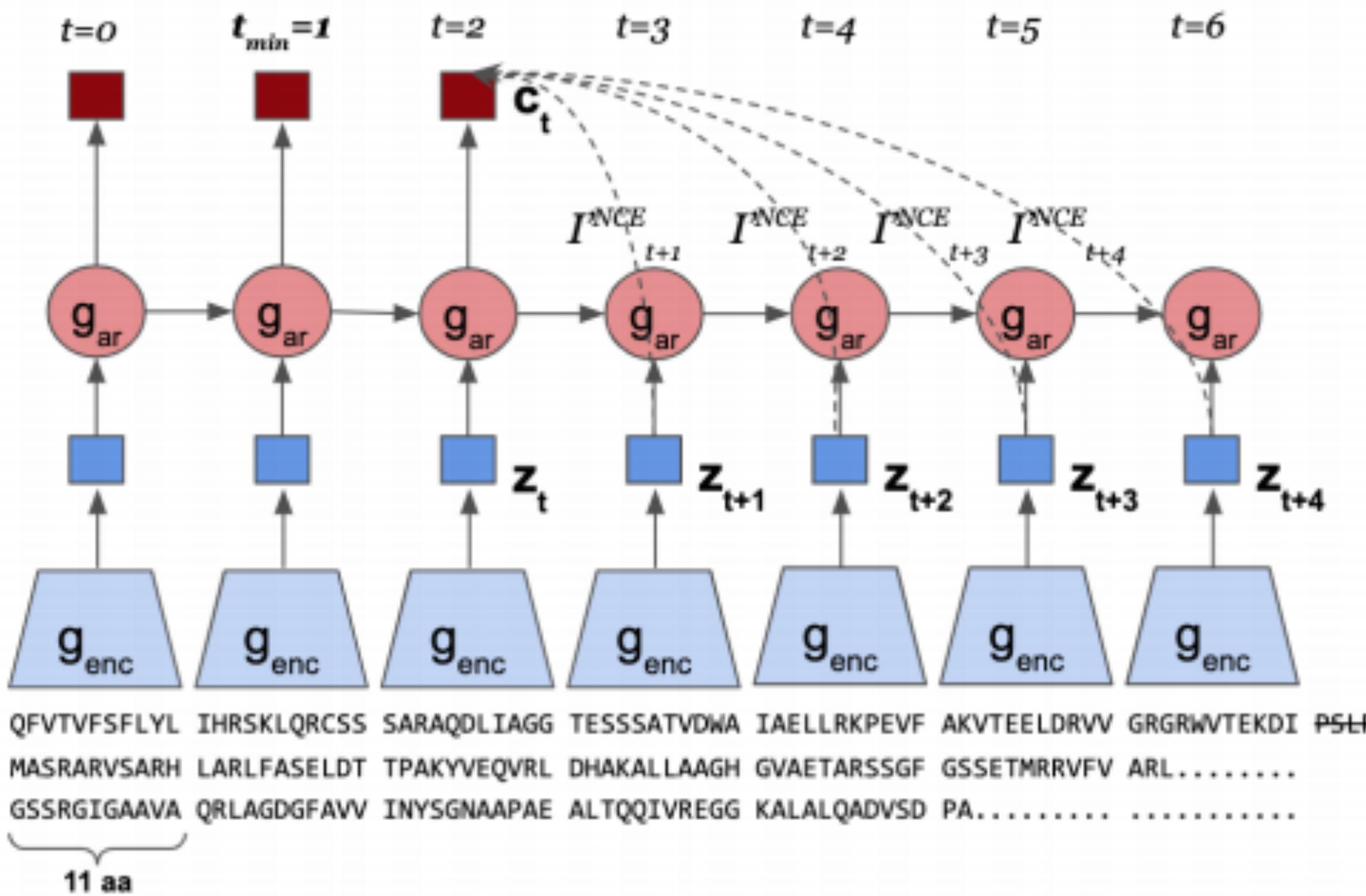


# Contrastive learning may require fewer parameters



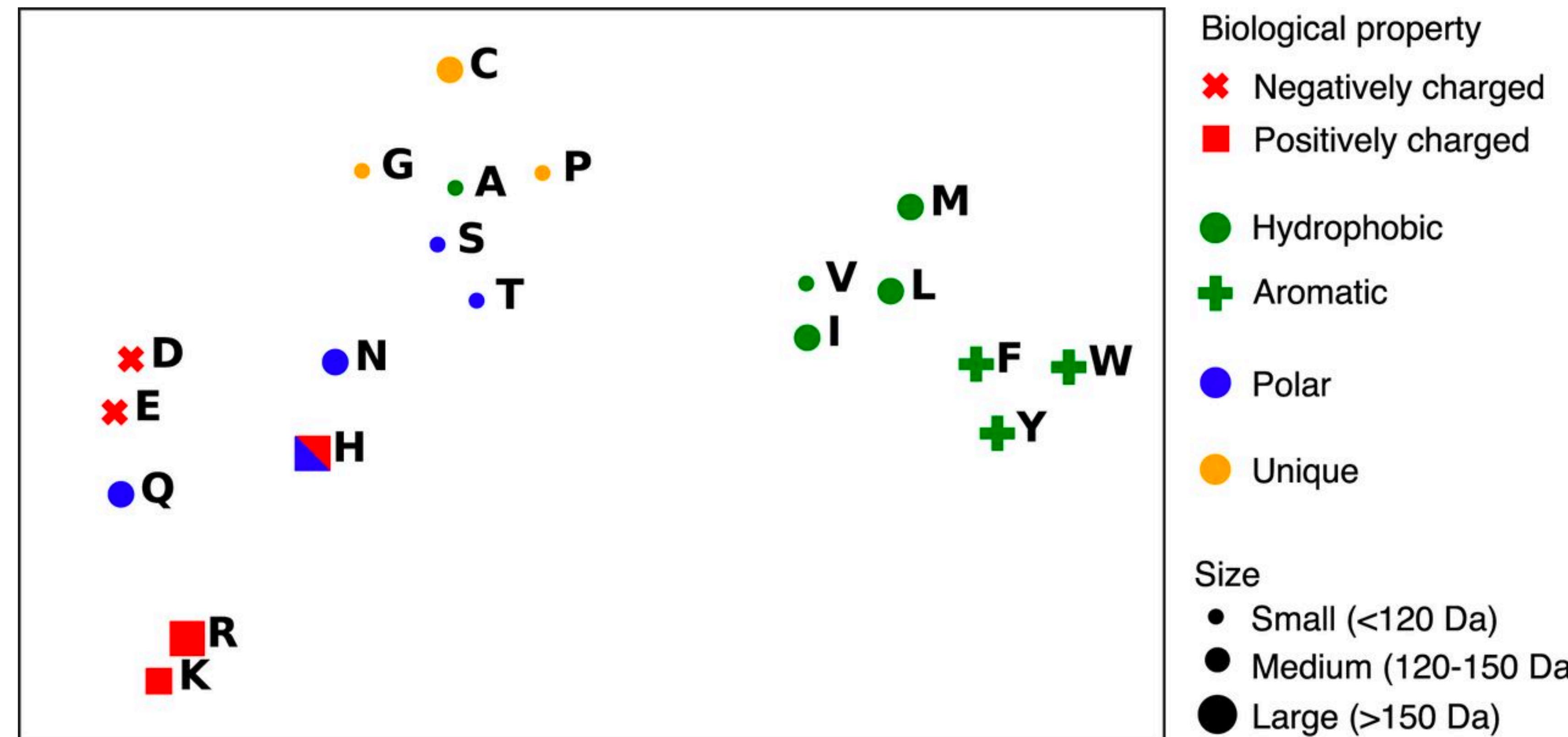
# of Embedding Parameters	Remote Homology			Secondary Structure			Stability	Fluorescence
	Fold	Superfamily	Family	CB513	CASP12	TS115		
Unirep	182M	0.23	0.38	0.87	0.73	<b>0.72</b>	0.77	<b>0.73</b>
BERT	92M	0.21	0.34	0.88	0.73	0.71	0.77	<b>0.73</b>
ResNet	48M	0.17	0.31	0.77	<b>0.75</b>	<b>0.72</b>	<b>0.78</b>	<b>0.73</b>
LSTM	44M	<b>0.26</b>	<b>0.43</b>	<b>0.92</b>	<b>0.75</b>	0.70	<b>0.78</b>	0.69
Bepler et al.	19M	0.17	0.20	0.79	0.73	0.70	0.76	0.64
One Hot	0	0.09	0.08	0.39	0.69	0.68	0.72	0.19
CPCProt	1.7M	0.12	0.12	0.48	0.69	0.70	0.73	0.65
CPCProt <sub>GRU</sub> <sub>large</sub>	8.4M	0.13	0.14	0.52	0.70	0.70	0.73	0.65
CPCProt <sub>LSTM</sub>	71M	0.11	0.11	0.47	0.68	0.66	0.70	0.68

# But is it really the contrastive learning that helps?

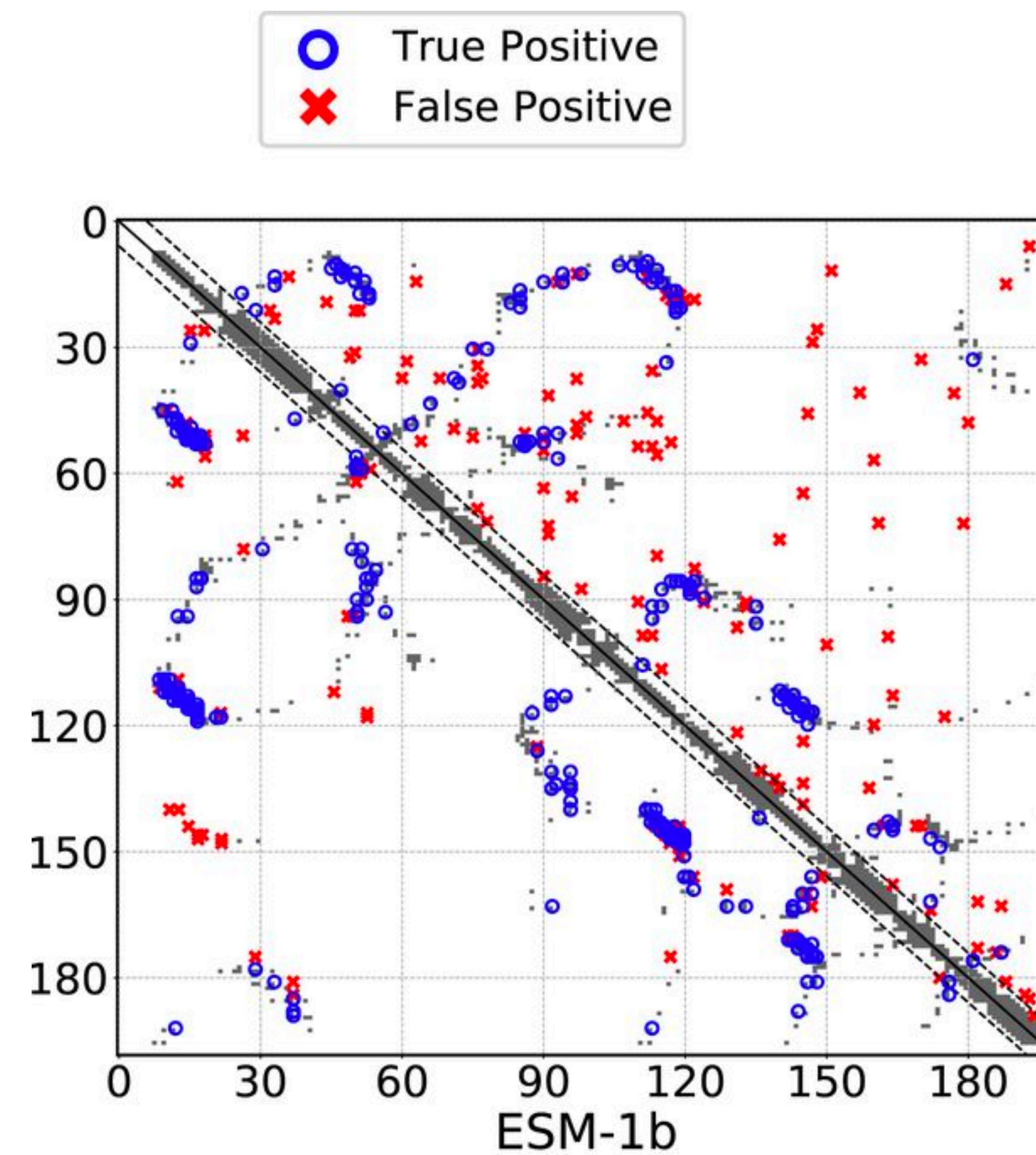


- In images, the augmentation pipeline is most important
- Maybe we just need to be predicting larger blocks?
- Or have better negatives

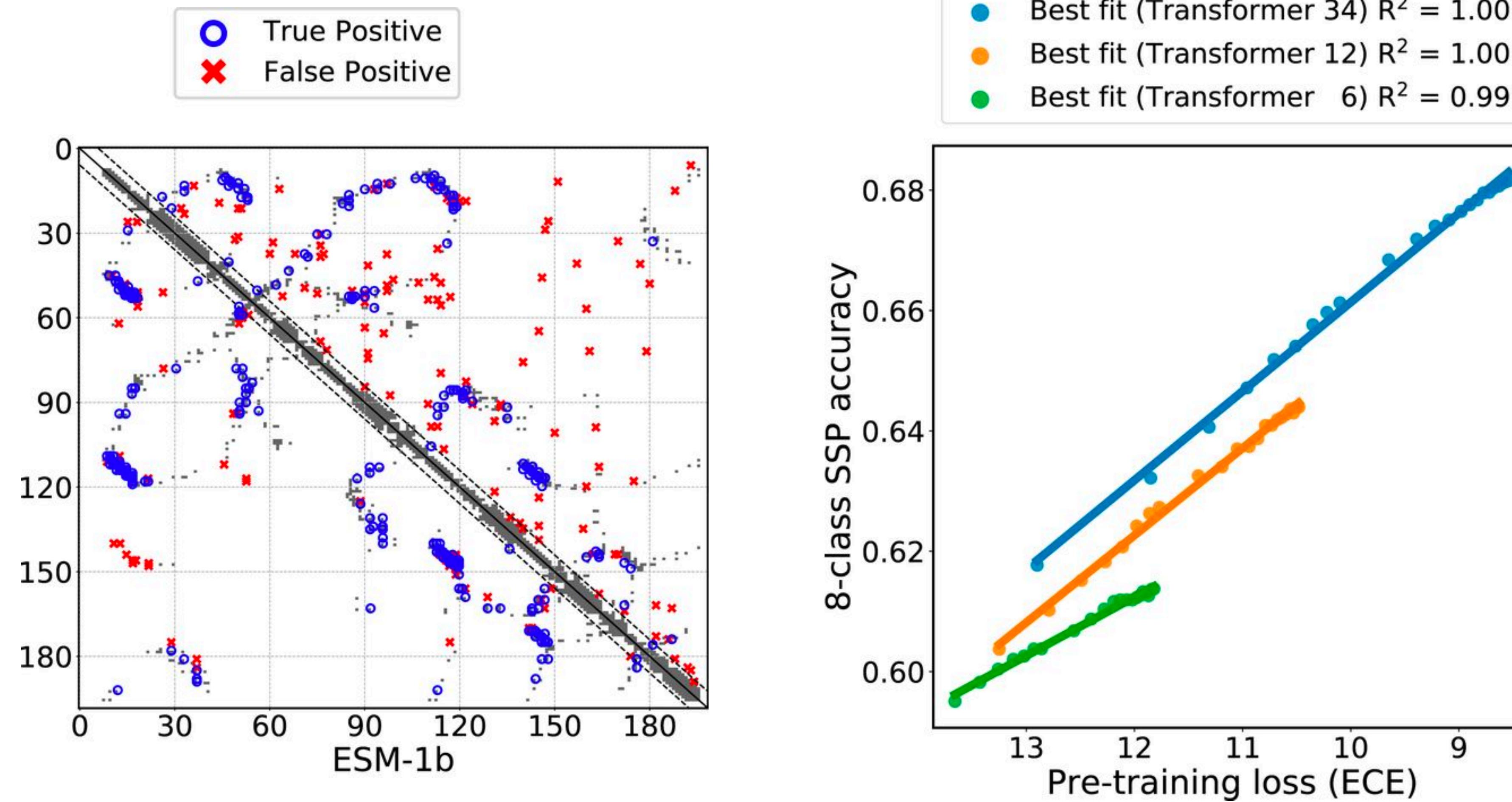
# Deep sequence embeddings recapitulate biophysical properties



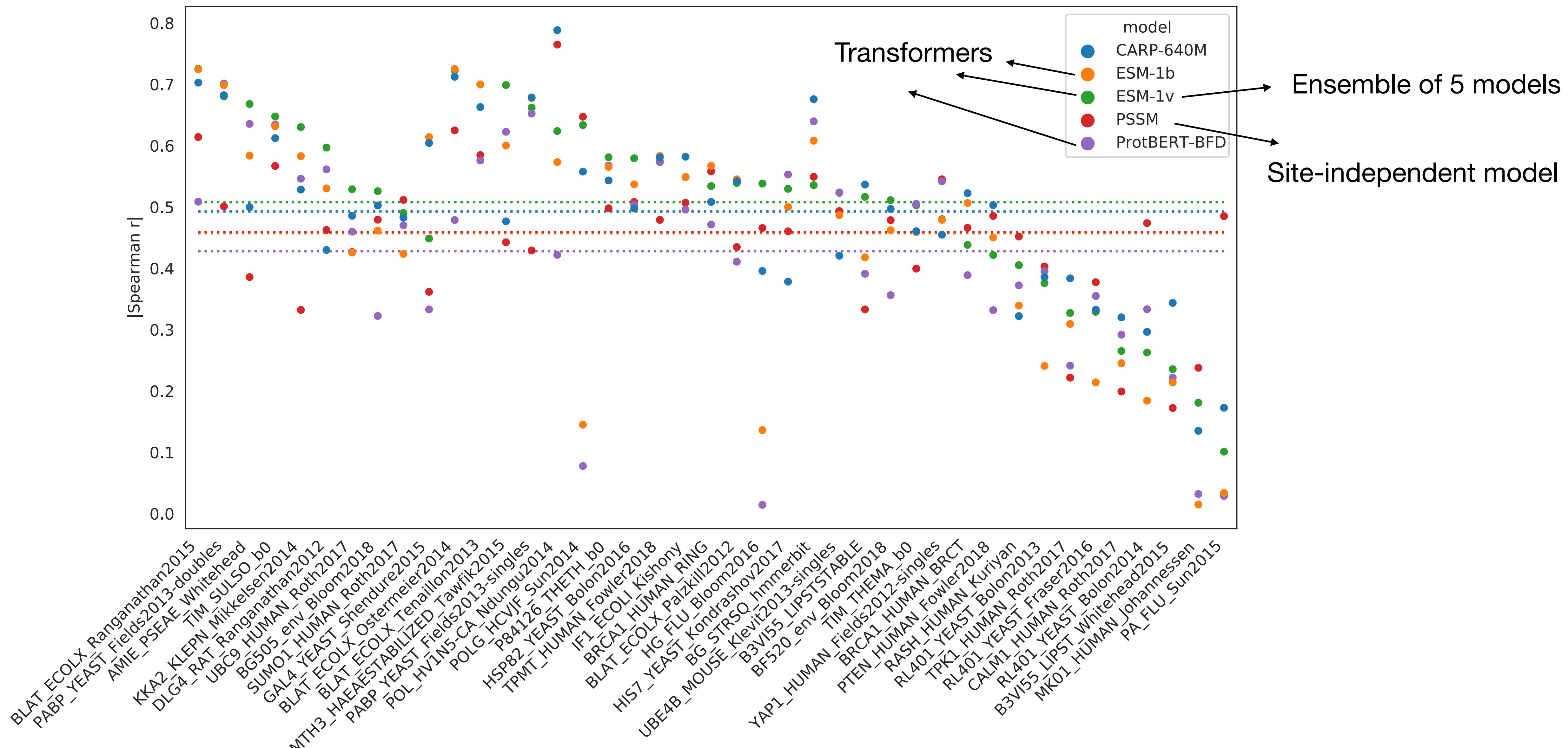
# Deep sequence embeddings contain structural information



# Deep sequence embeddings contain structural information



# MLM pseudo likelihood predicts mutant fitness



# Pretrained models contain evolutionary information

# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance

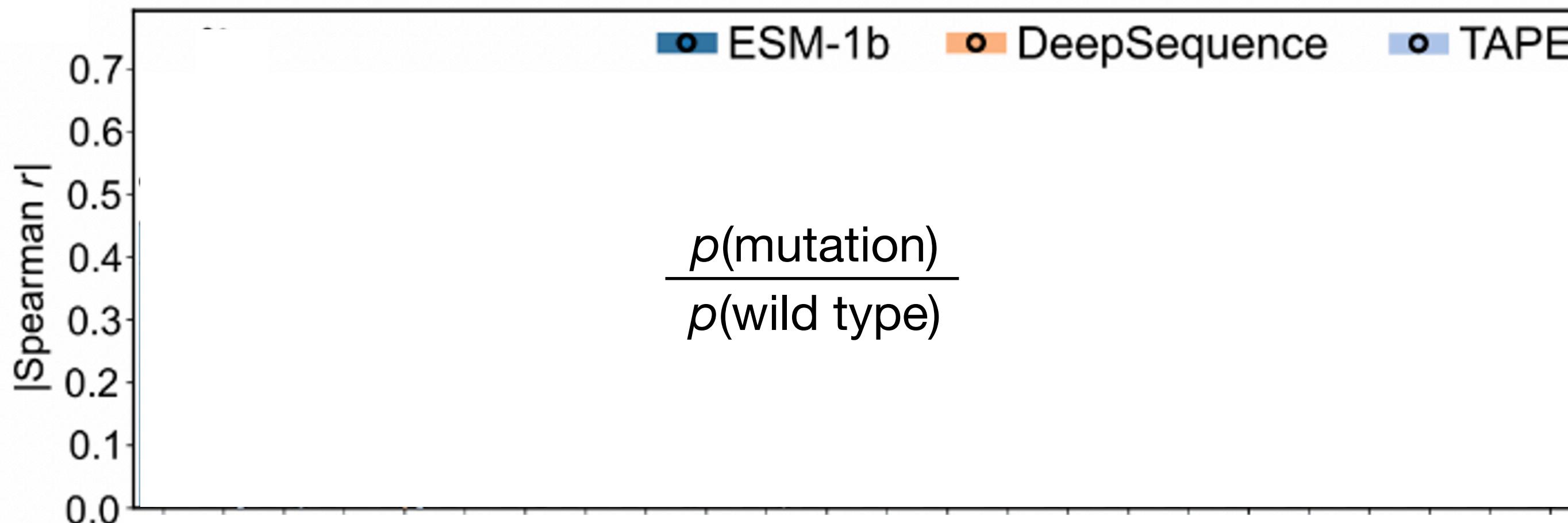
# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance

$$\frac{p(\text{mutation})}{p(\text{wild type})}$$

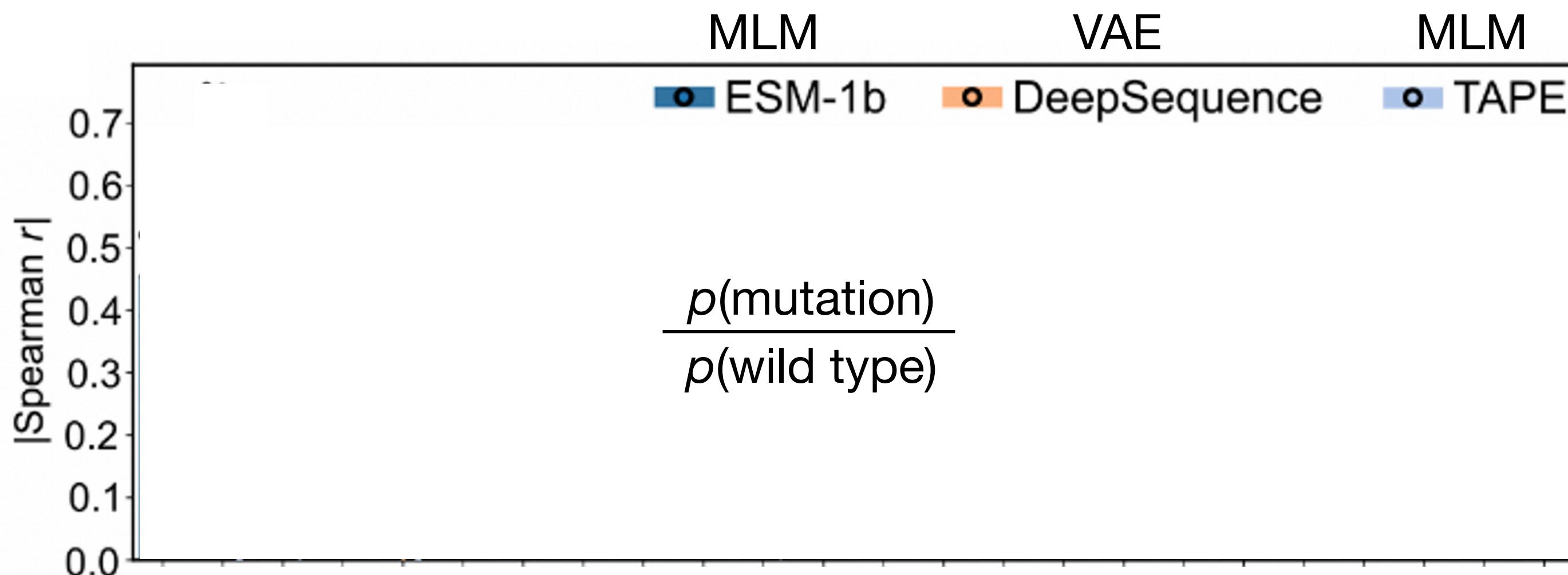
# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance



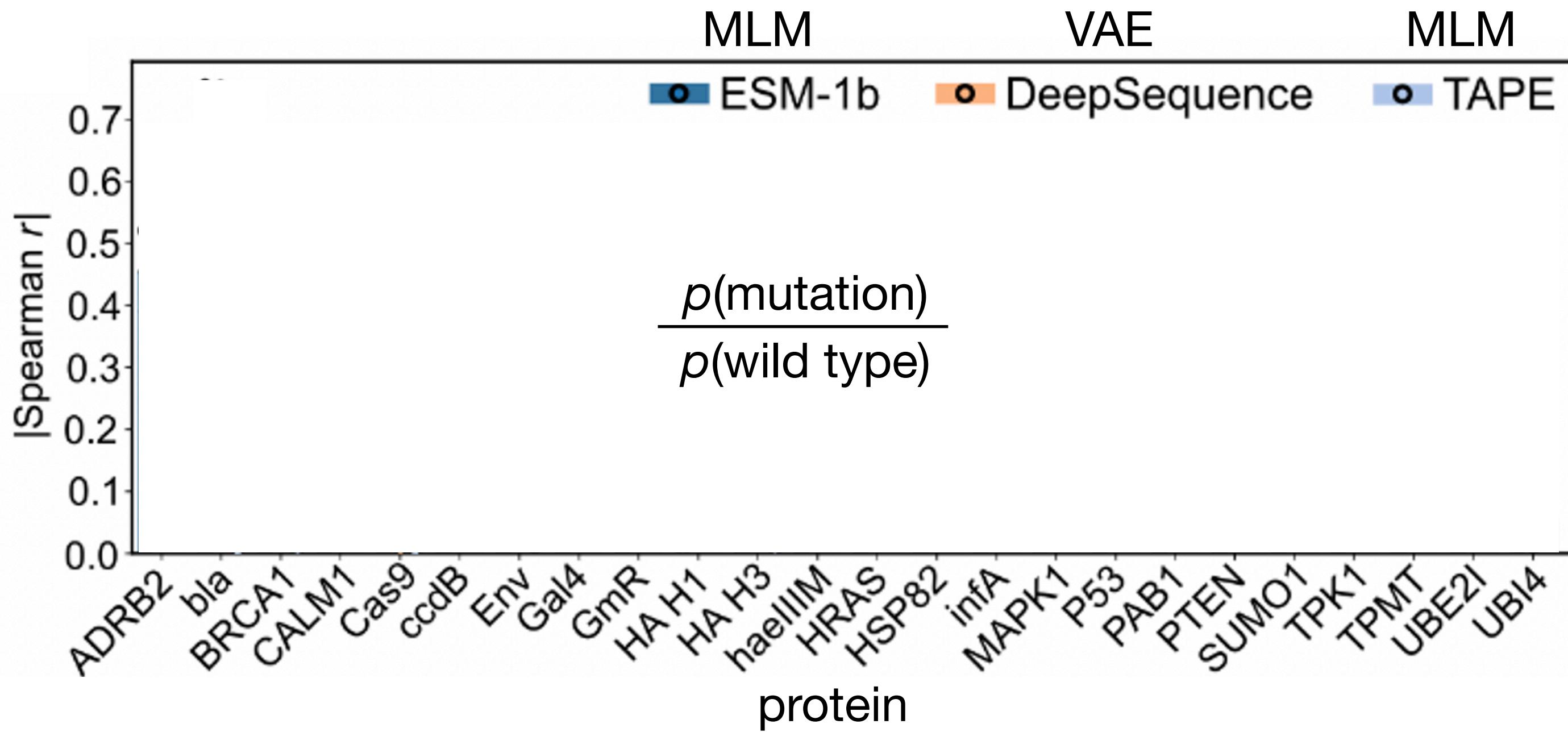
# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance



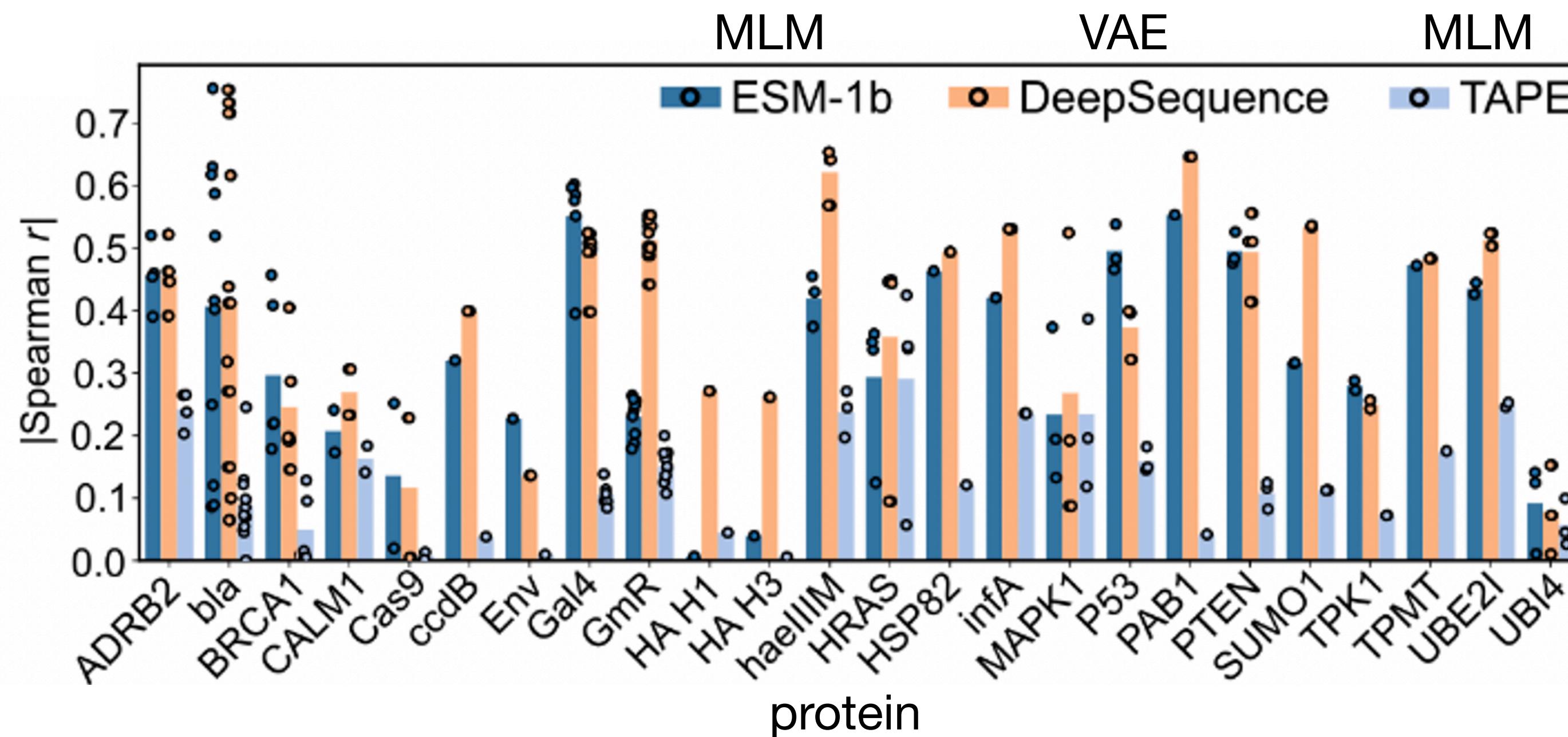
# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance



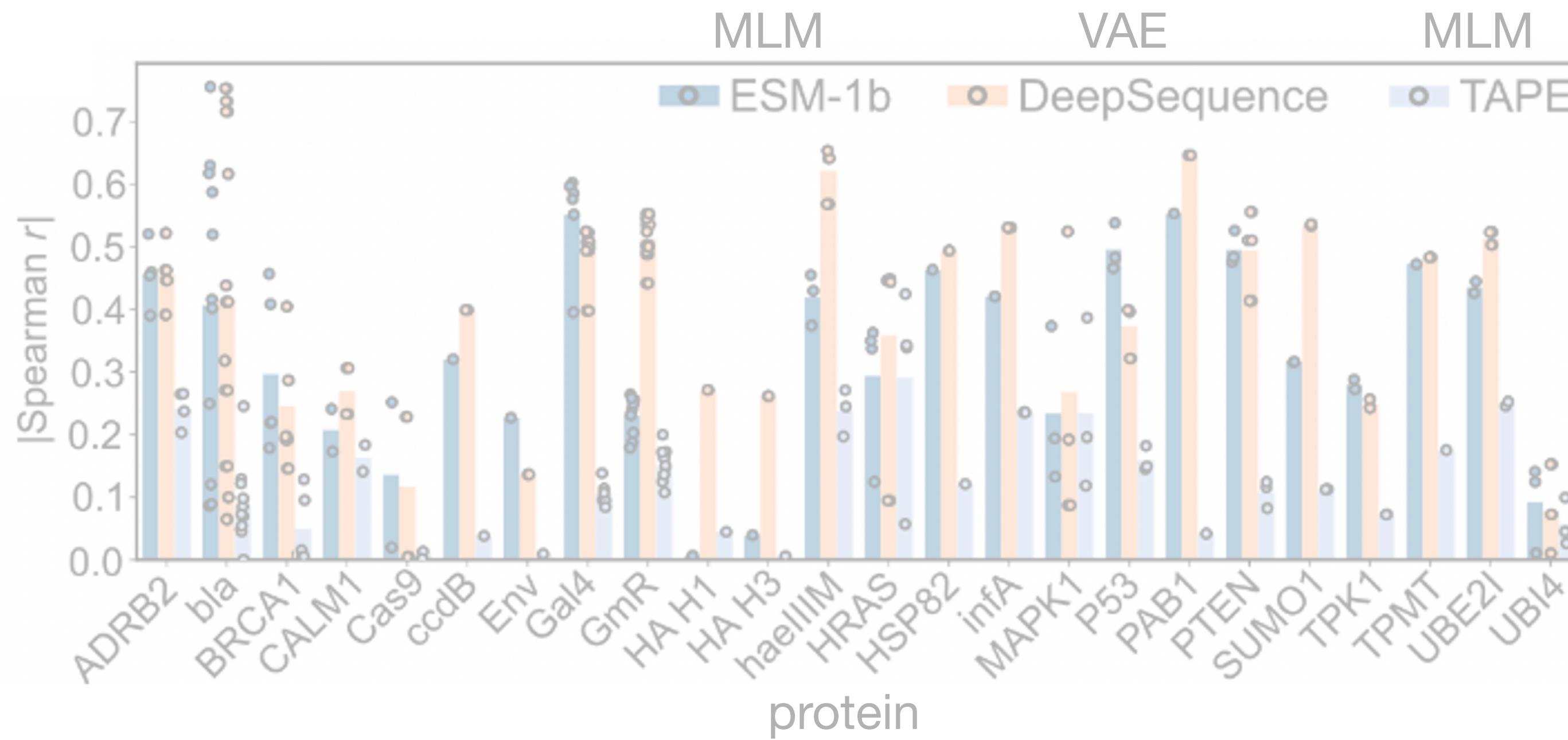
# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance



# Pretrained models contain evolutionary information

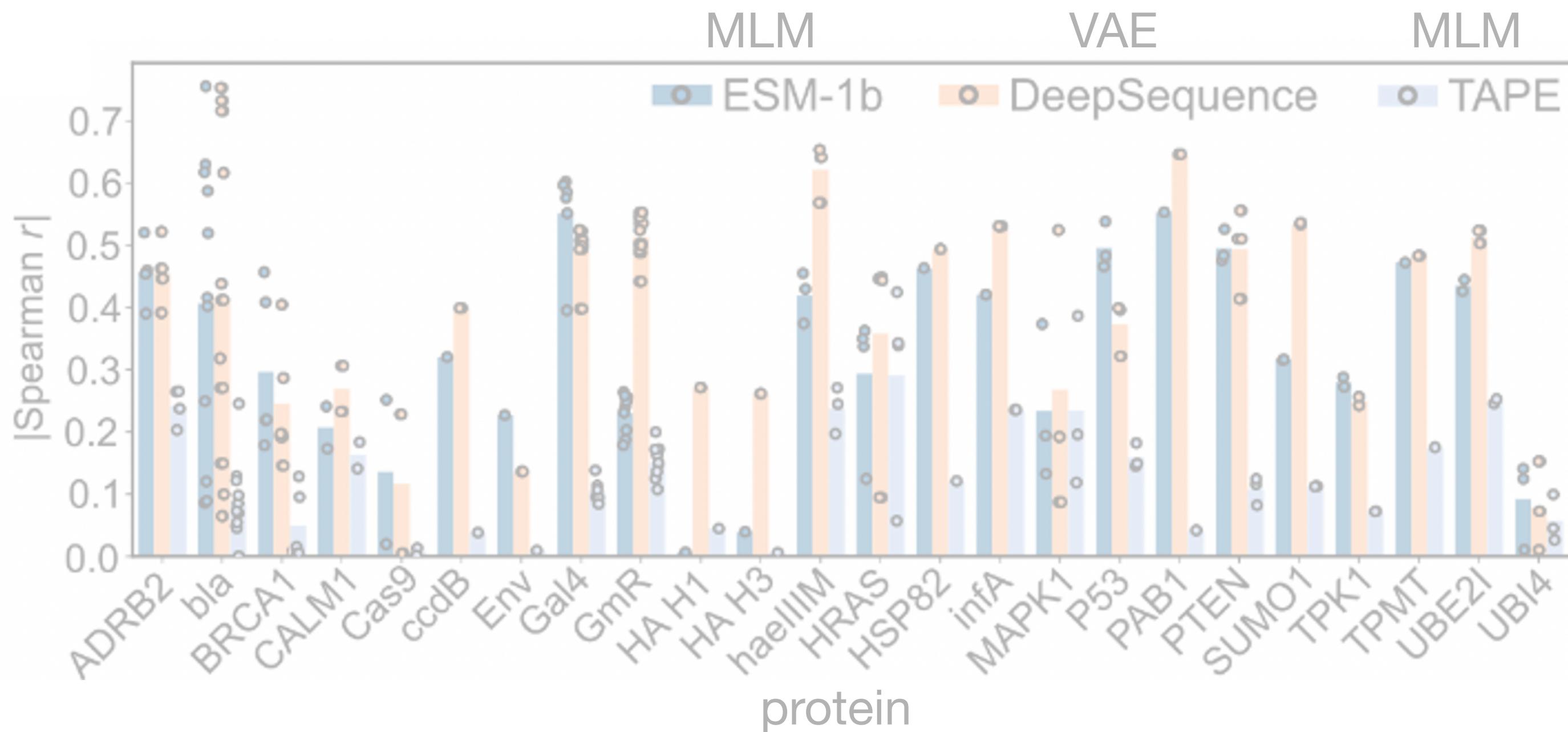
Pseudolikelihood correlates with mutational tolerance



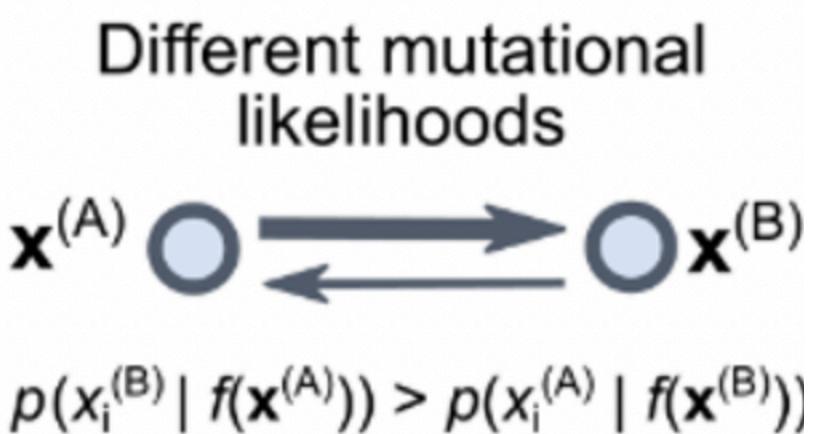
Use pseudolikelihood to predict evolutionary vector field

# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance



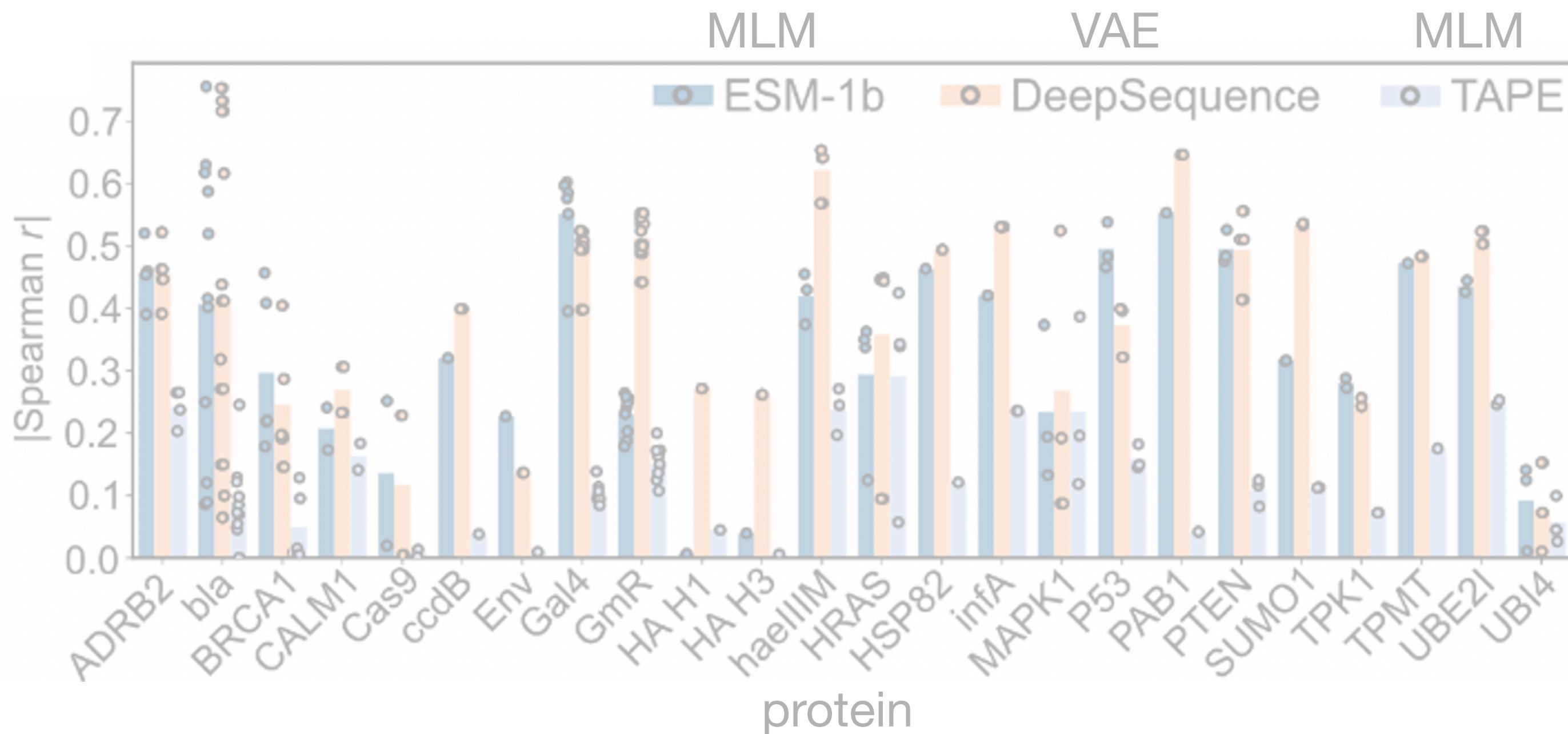
Use pseudolikelihood to predict evolutionary vector field



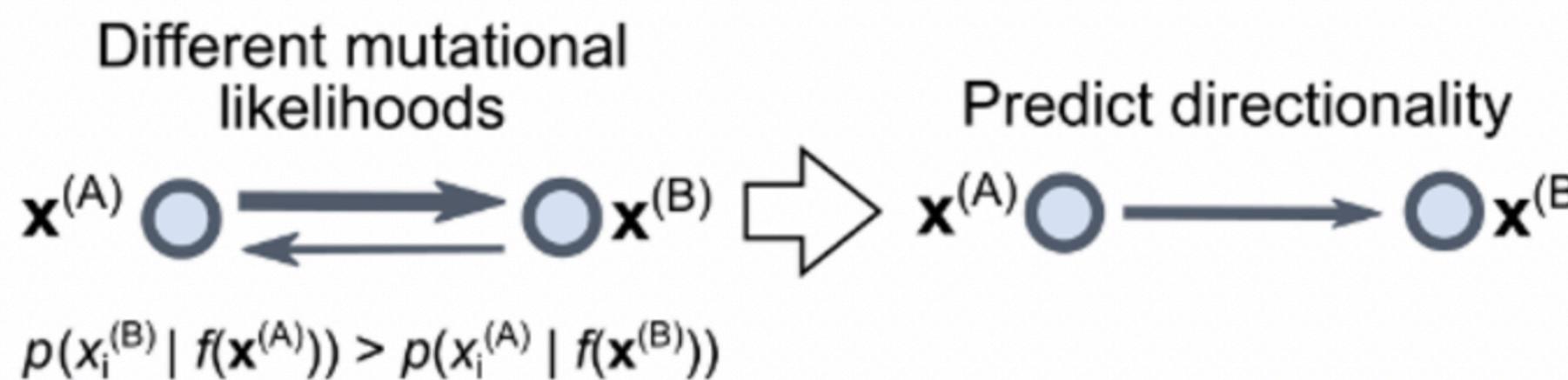
Hie et al. 2022

# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance



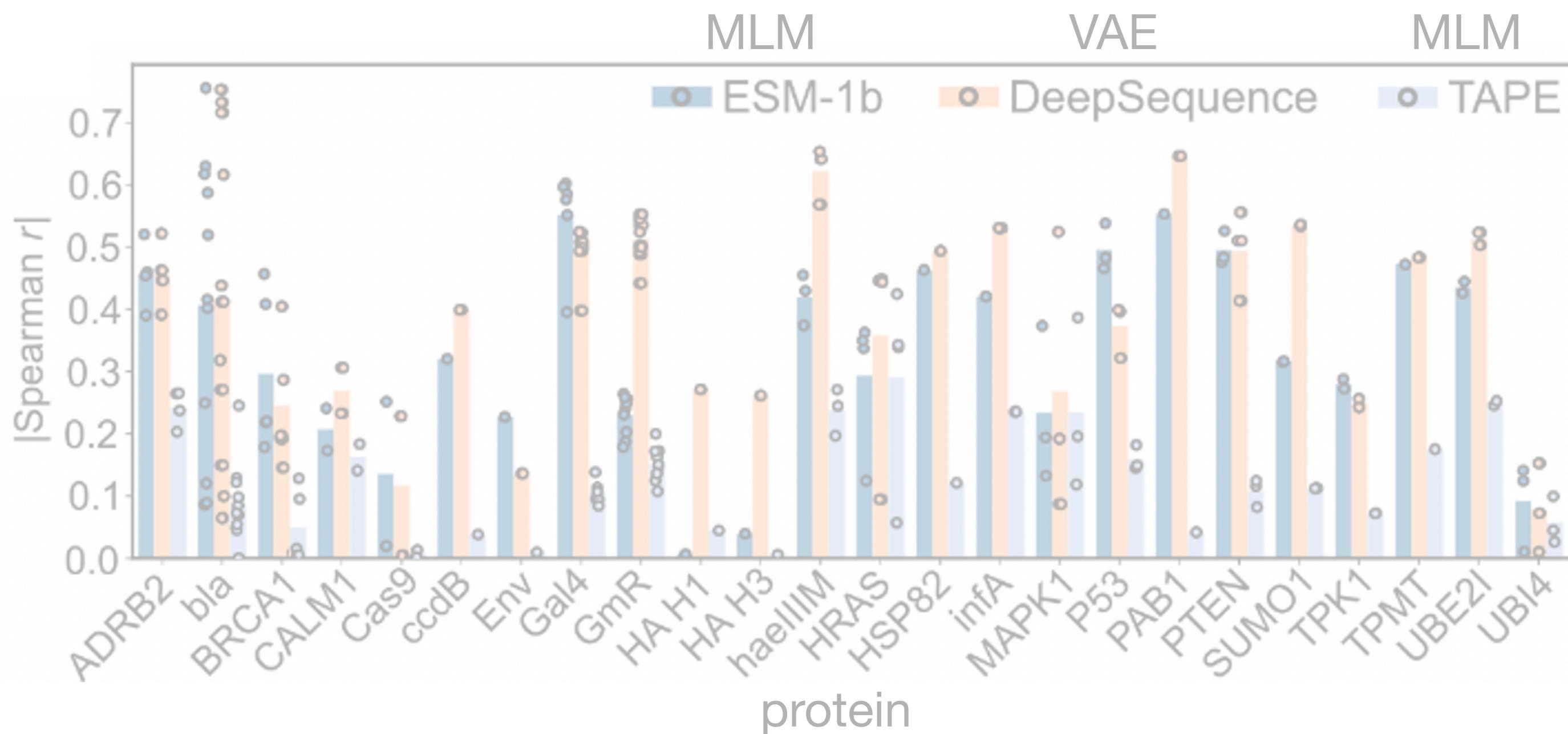
Use pseudolikelihood to predict evolutionary vector field



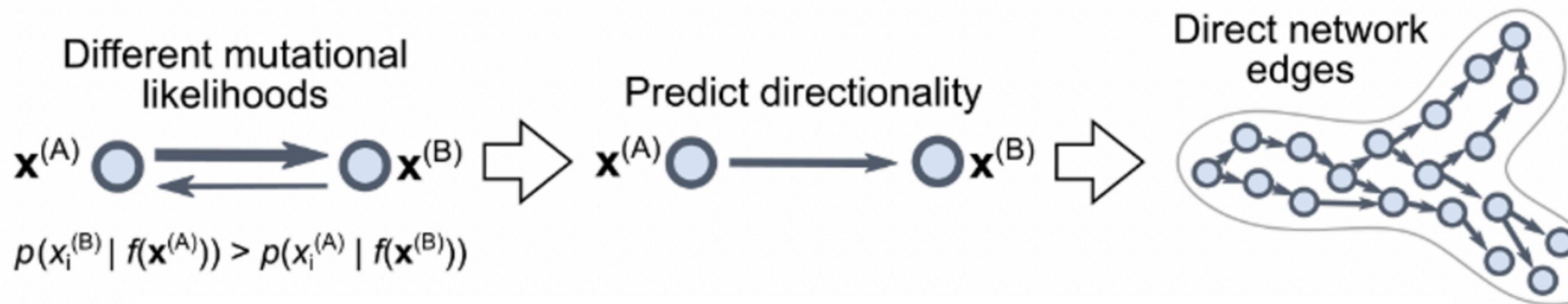
Hie et al. 2022

# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance



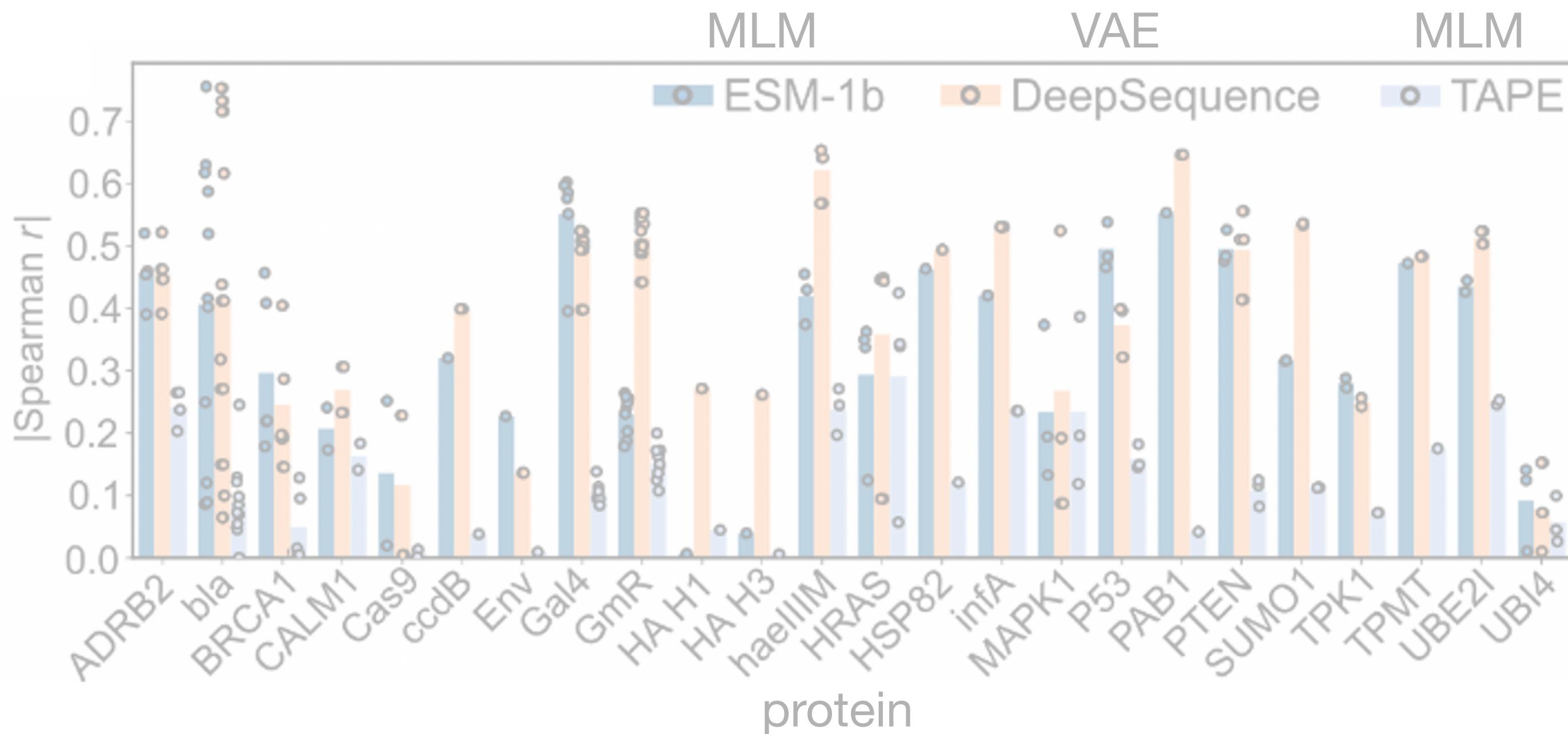
Use pseudolikelihood to predict evolutionary vector field



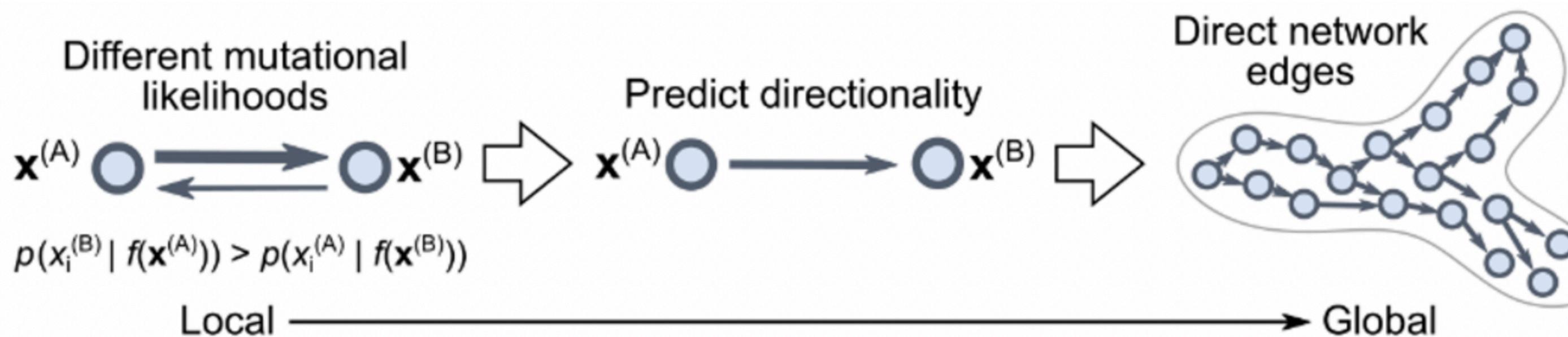
Hie et al. 2022

# Pretrained models contain evolutionary information

Pseudolikelihood correlates with mutational tolerance



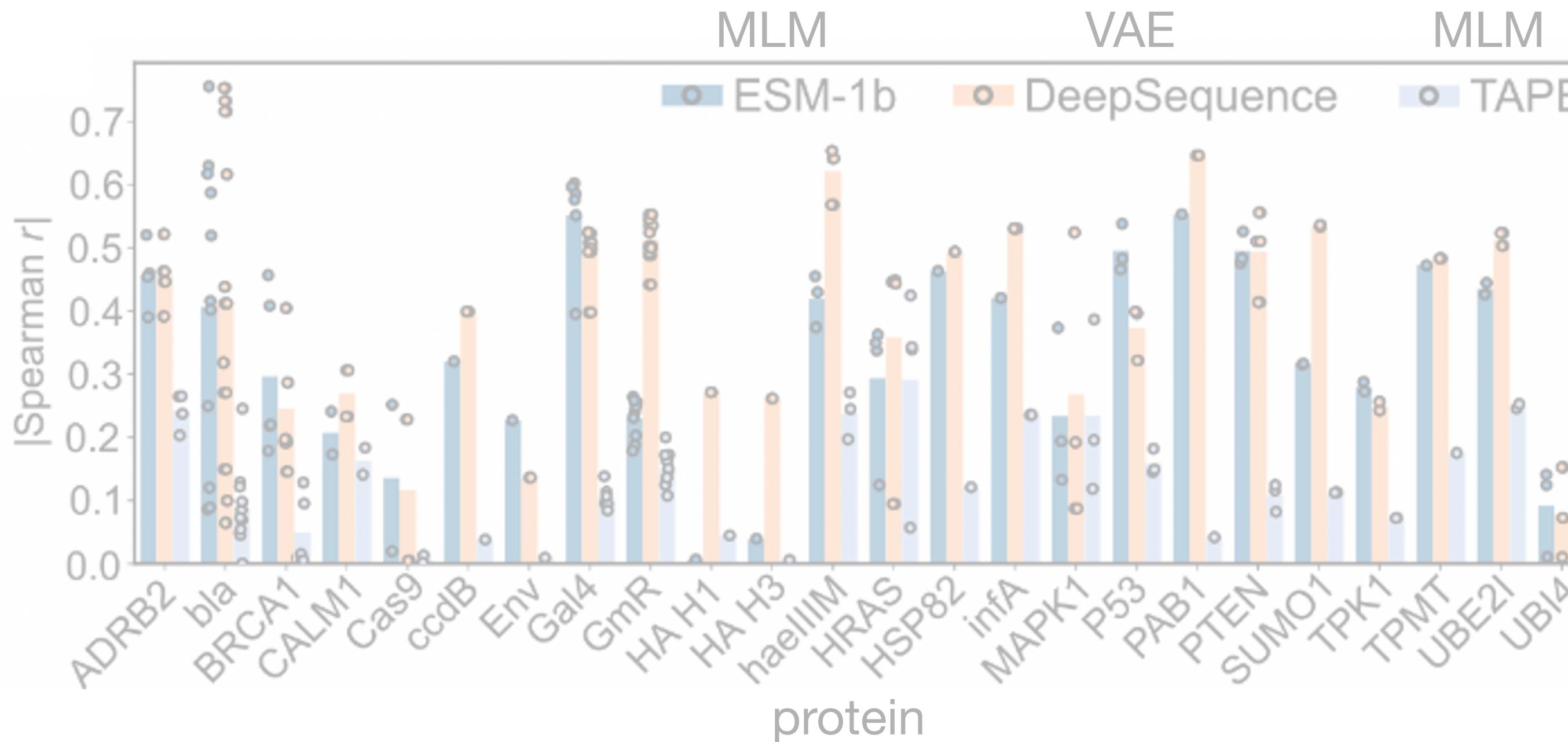
Use pseudolikelihood to predict evolutionary vector field



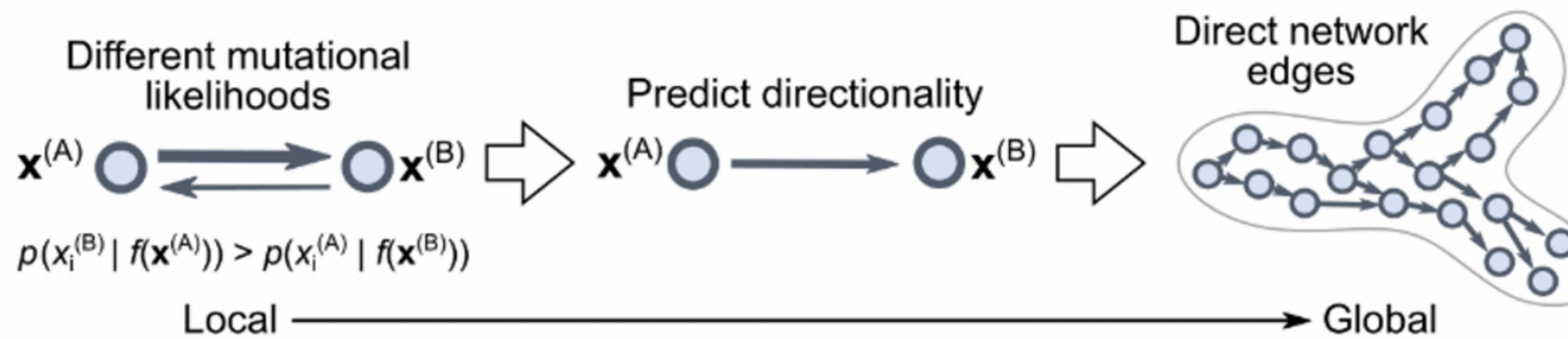
Hie et al. 2022

# Pretrained models contain evolutionary information

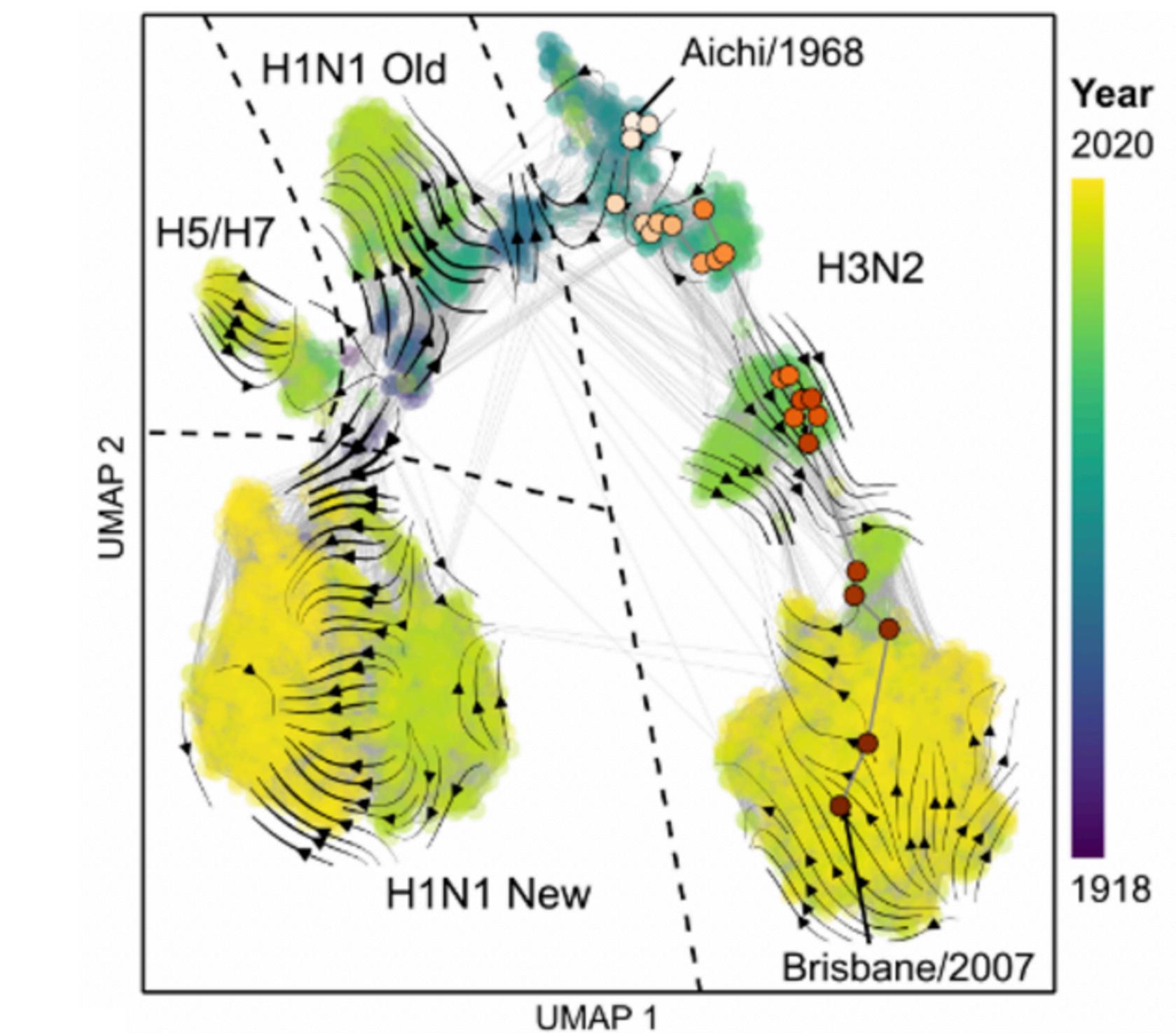
Pseudolikelihood correlates with mutational tolerance



Use pseudolikelihood to predict evolutionary vector field



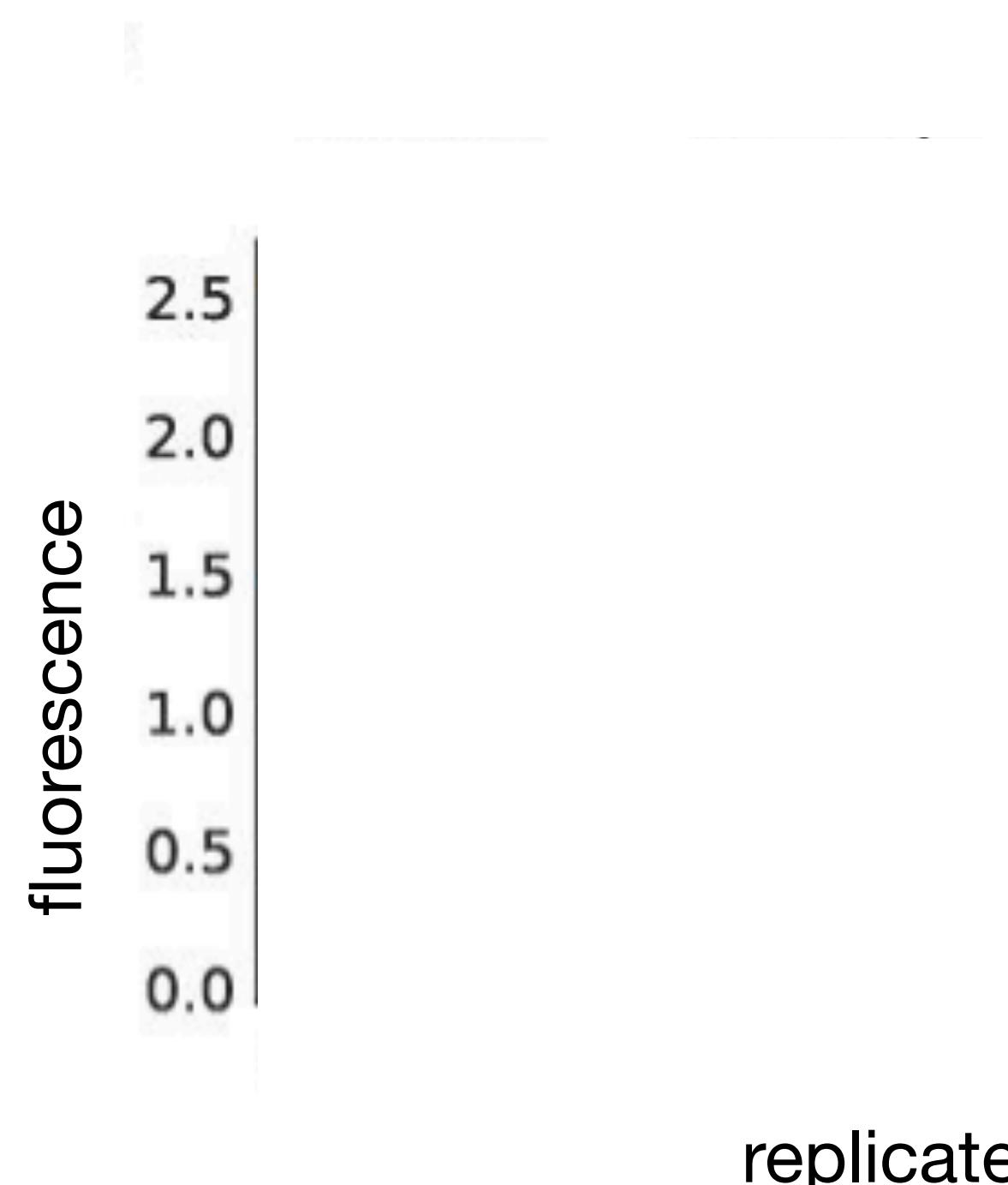
Influenza A nucleoprotein



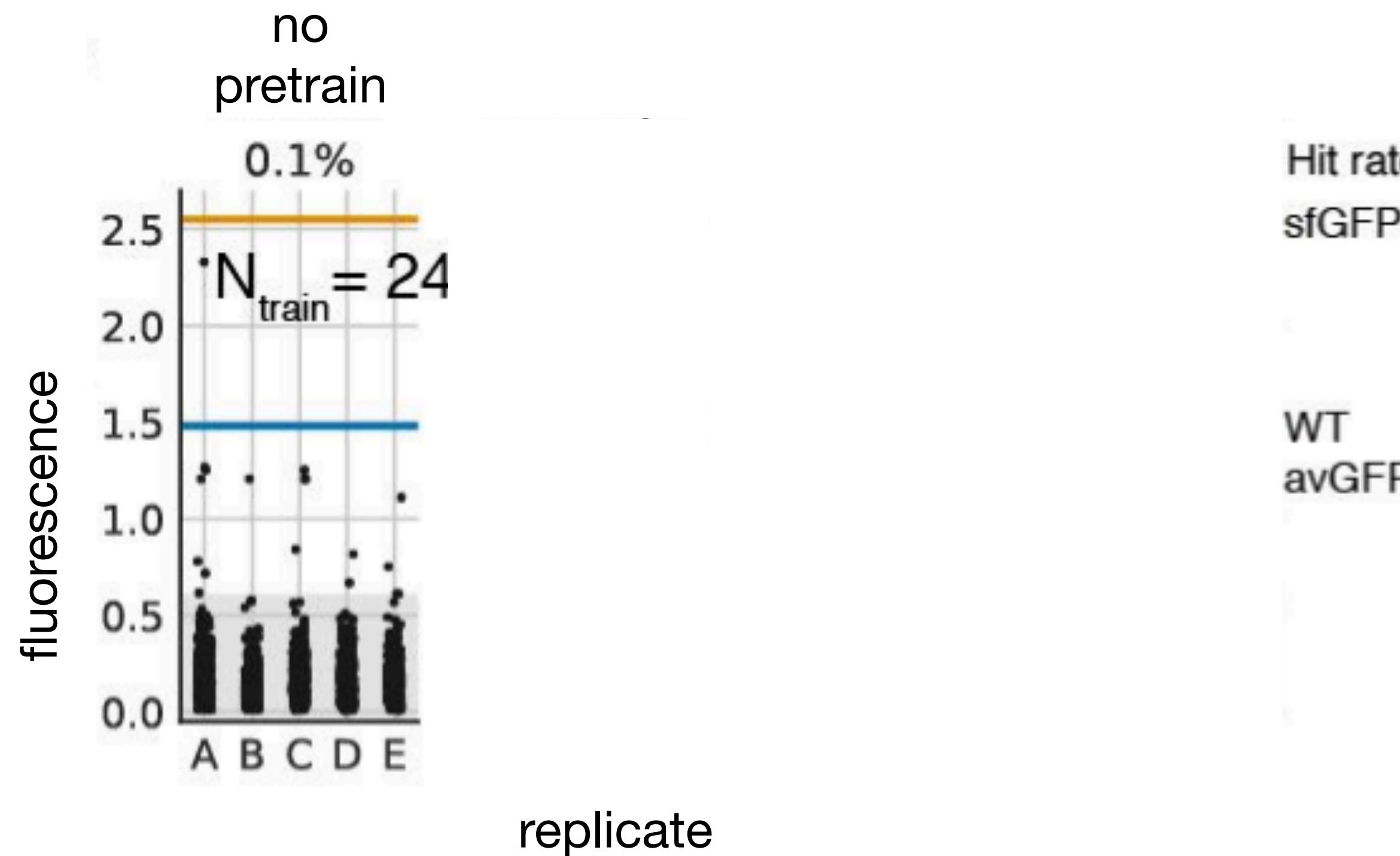
Hie et al. 2022

# Pretraining guides search away from loss-of-function sequences

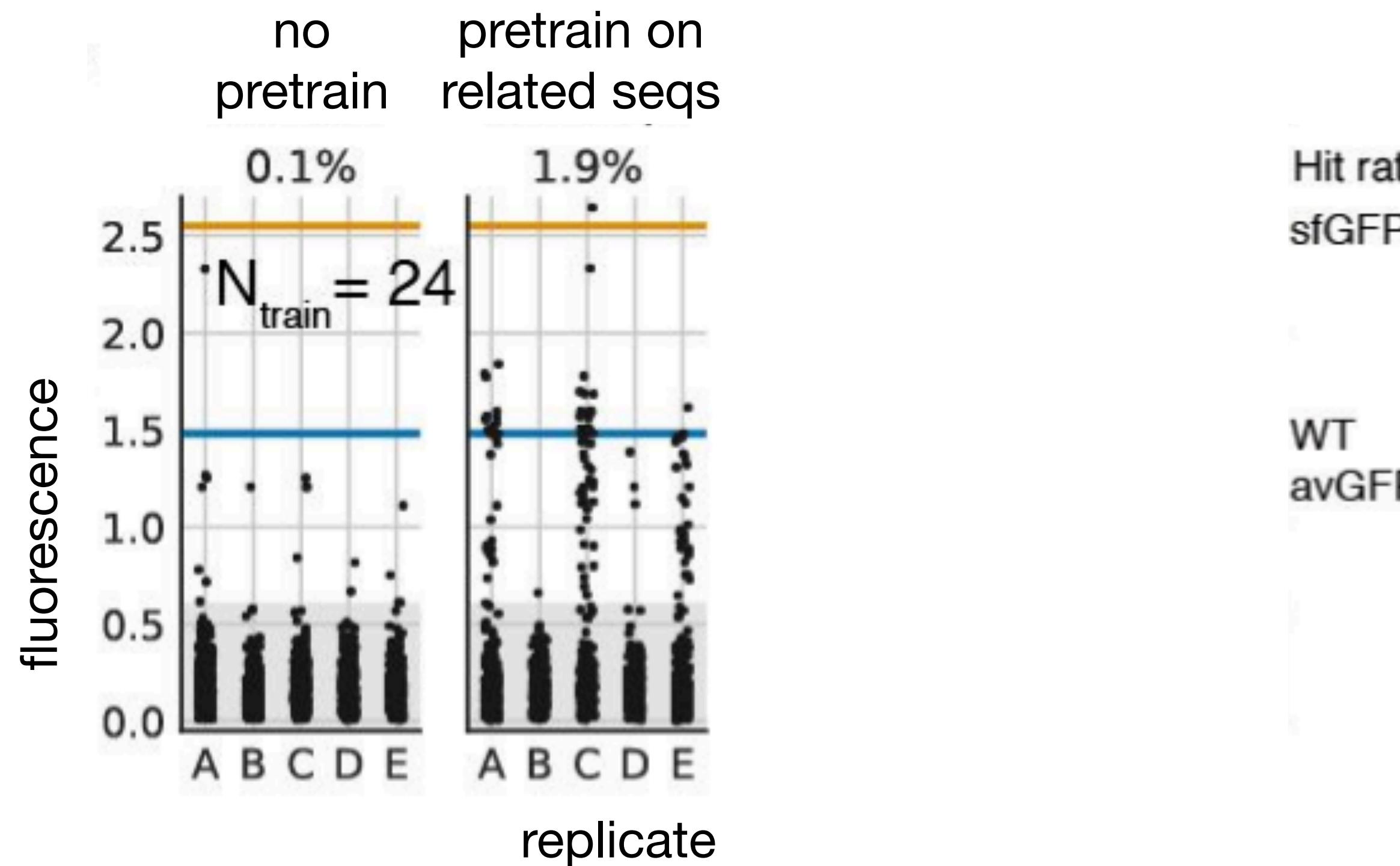
# Pretraining guides search away from loss-of-function sequences



# Pretraining guides search away from loss-of-function sequences

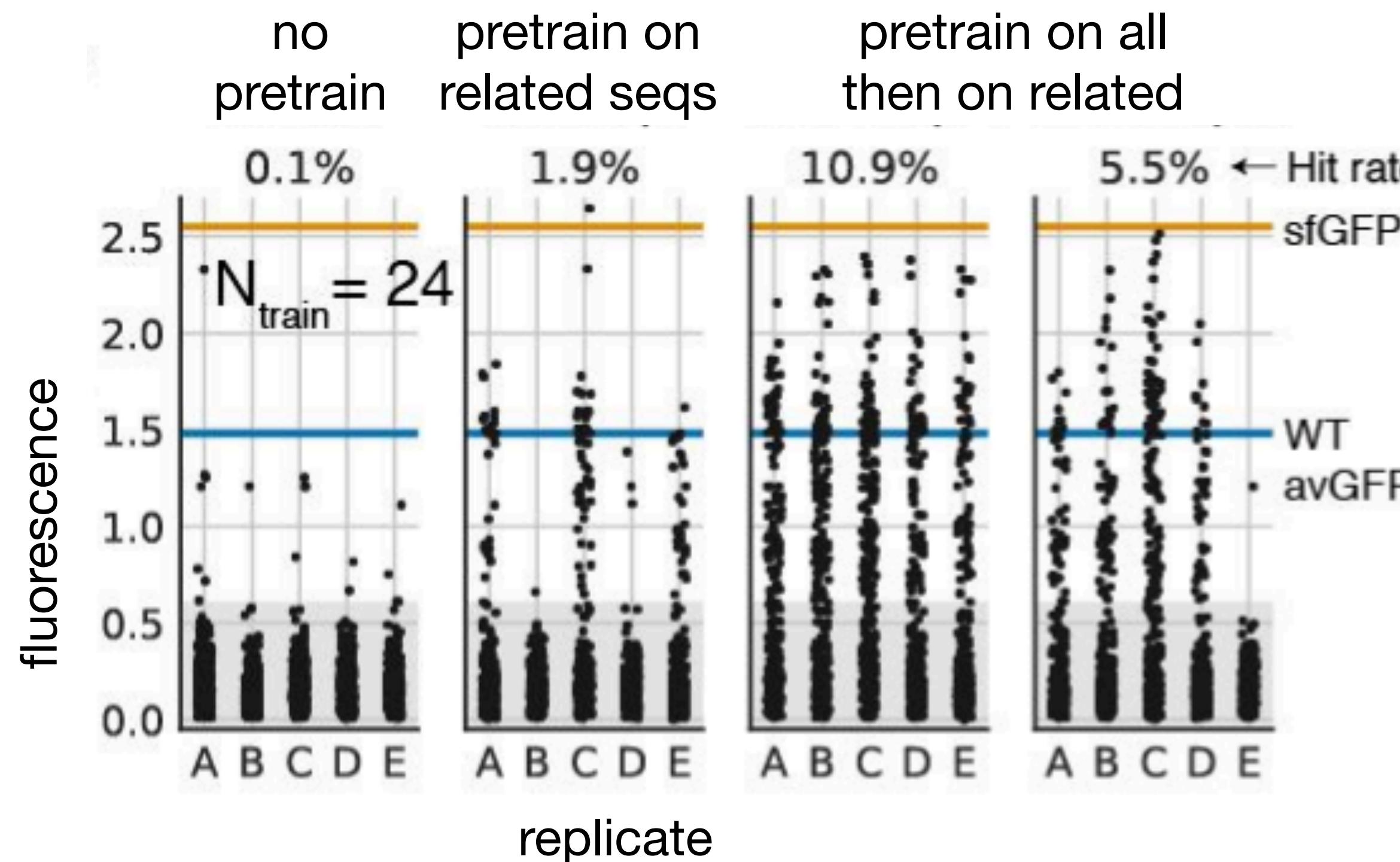


# Pretraining guides search away from loss-of-function sequences

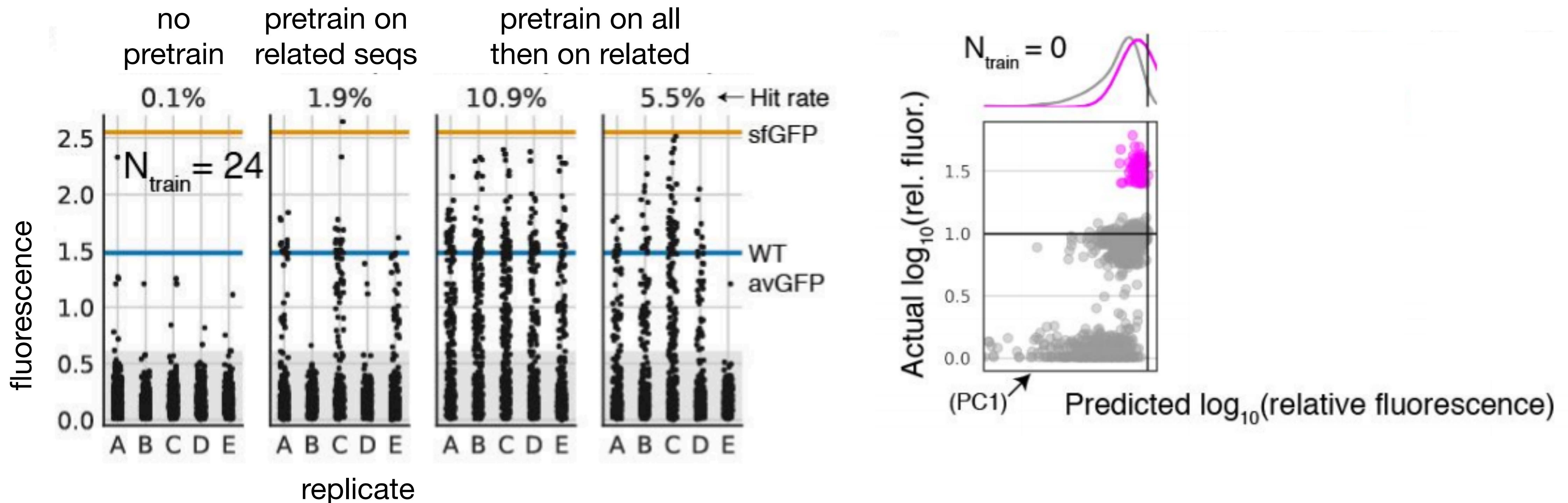


Hit rate  
sfGFP  
WT  
avGFP

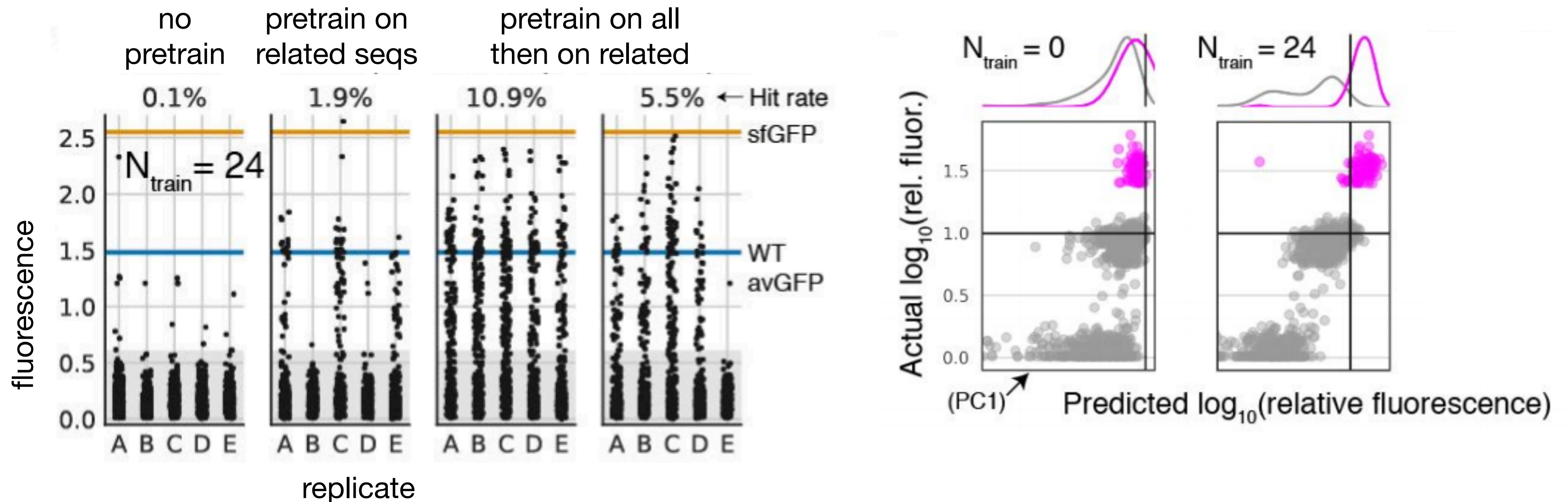
# Pretraining guides search away from loss-of-function sequences



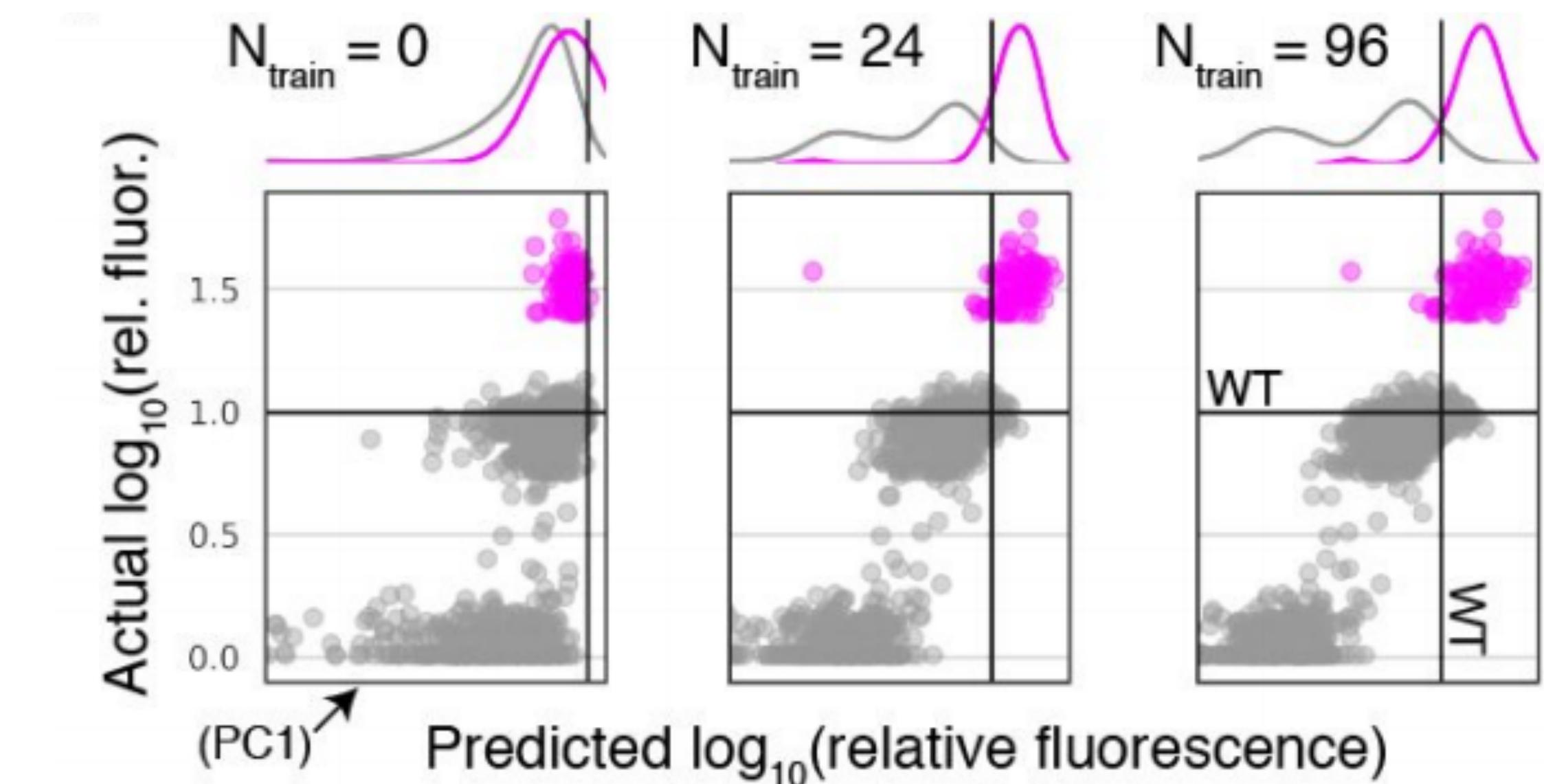
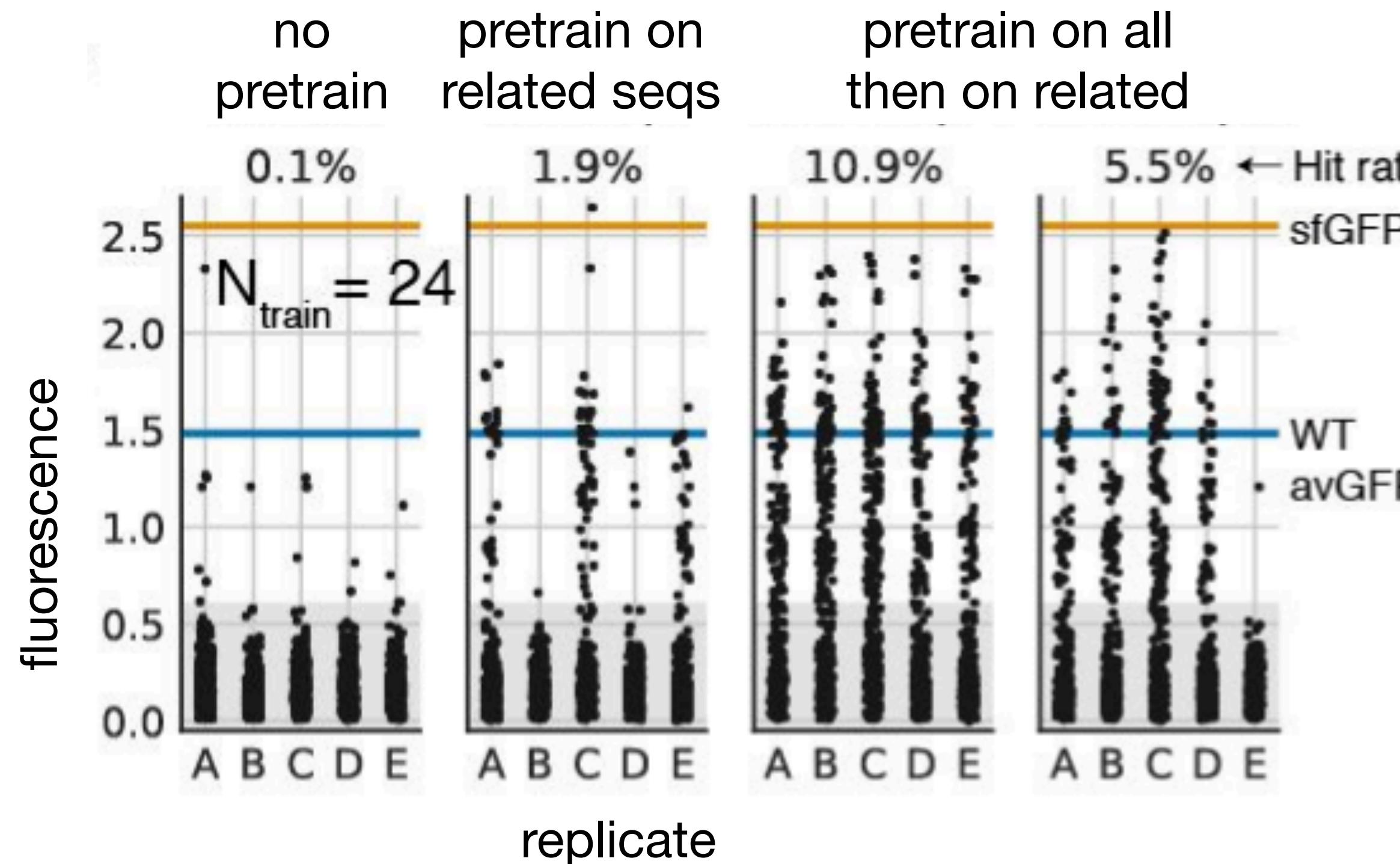
# Pretraining guides search away from loss-of-function sequences



# Pretraining guides search away from loss-of-function sequences



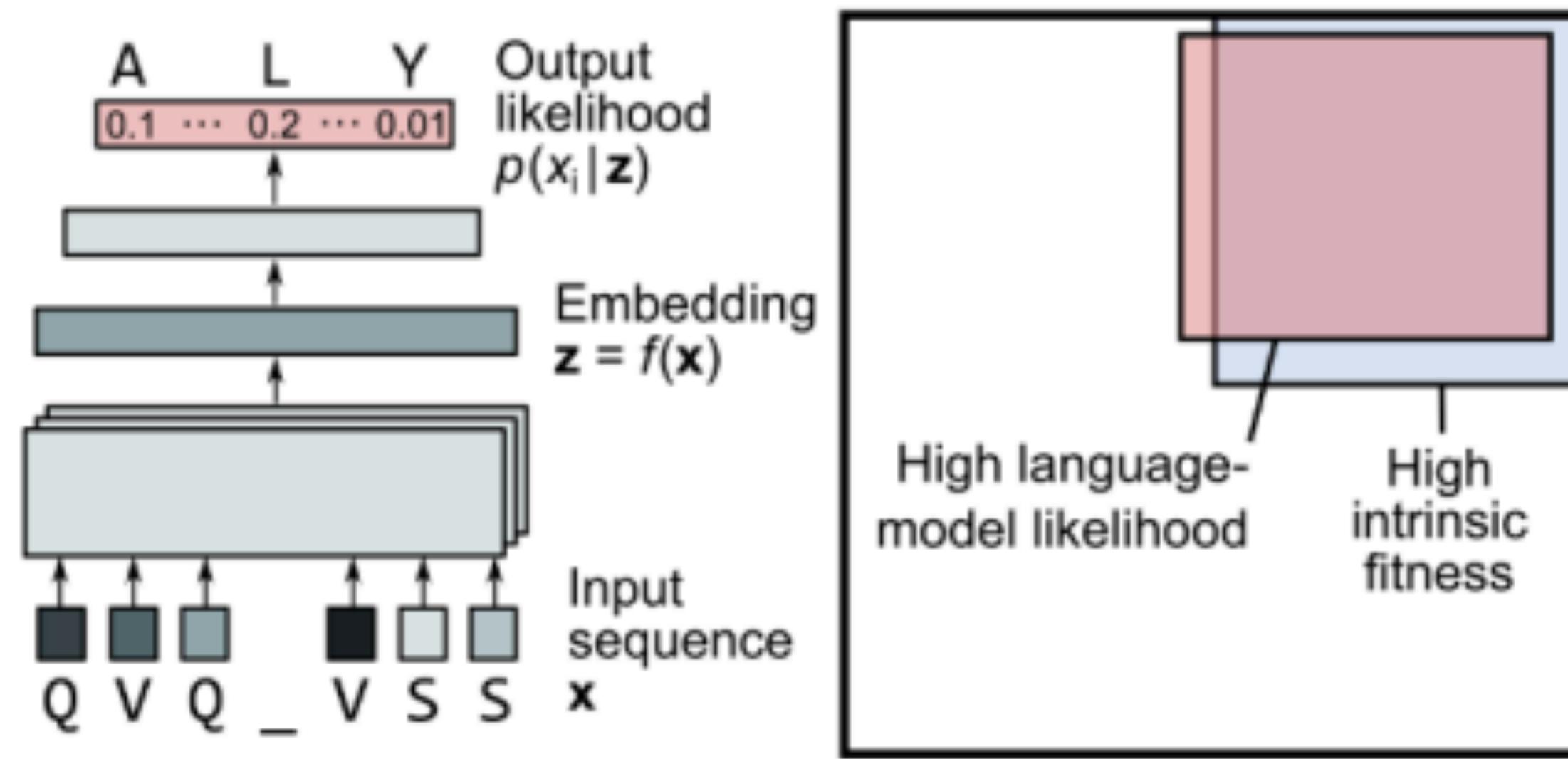
# Pretraining guides search away from loss-of-function sequences



# Pretrained models propose high-fitness sequences

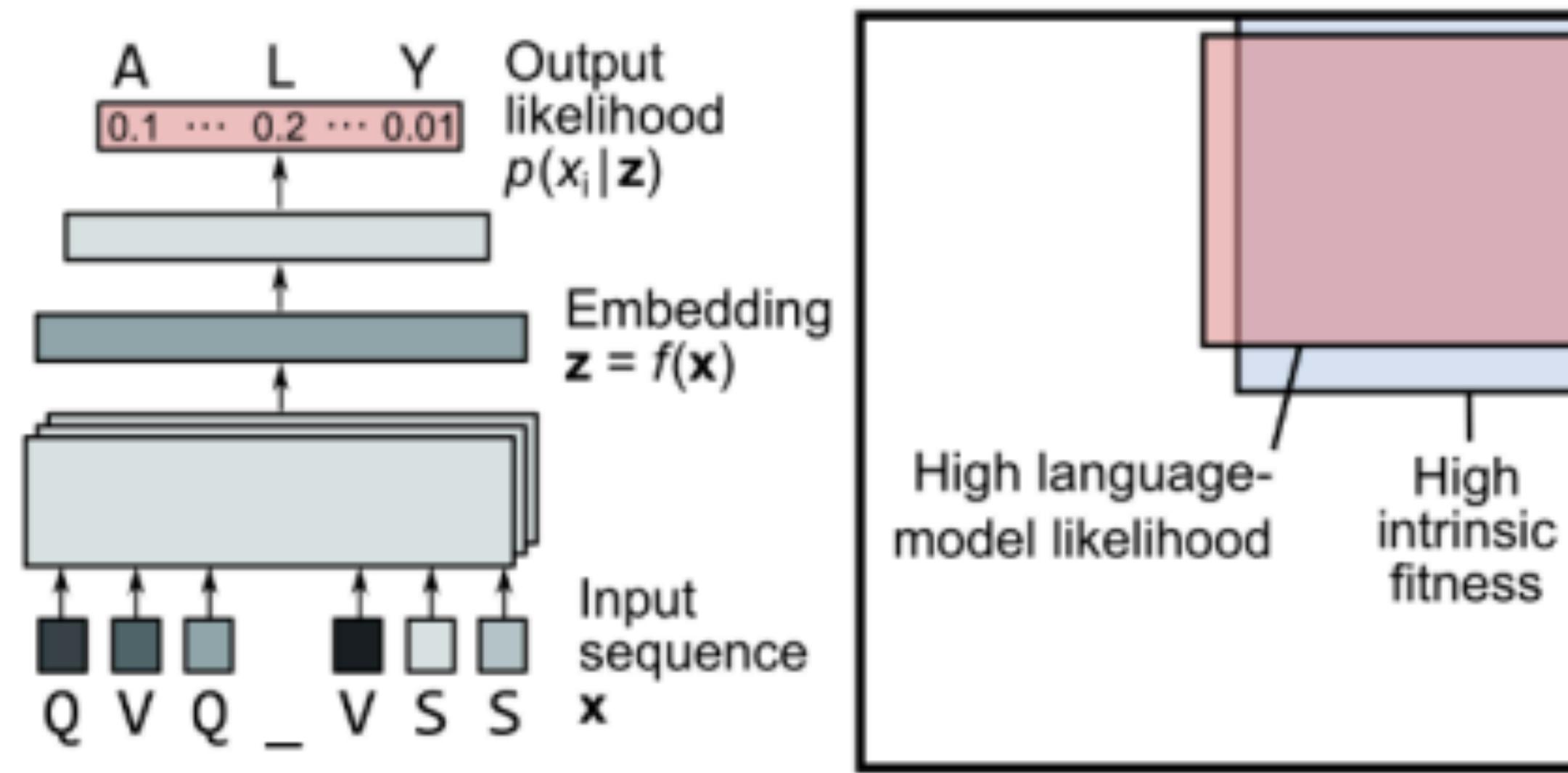
Hie et al. 2022

# Pretrained models propose high-fitness sequences

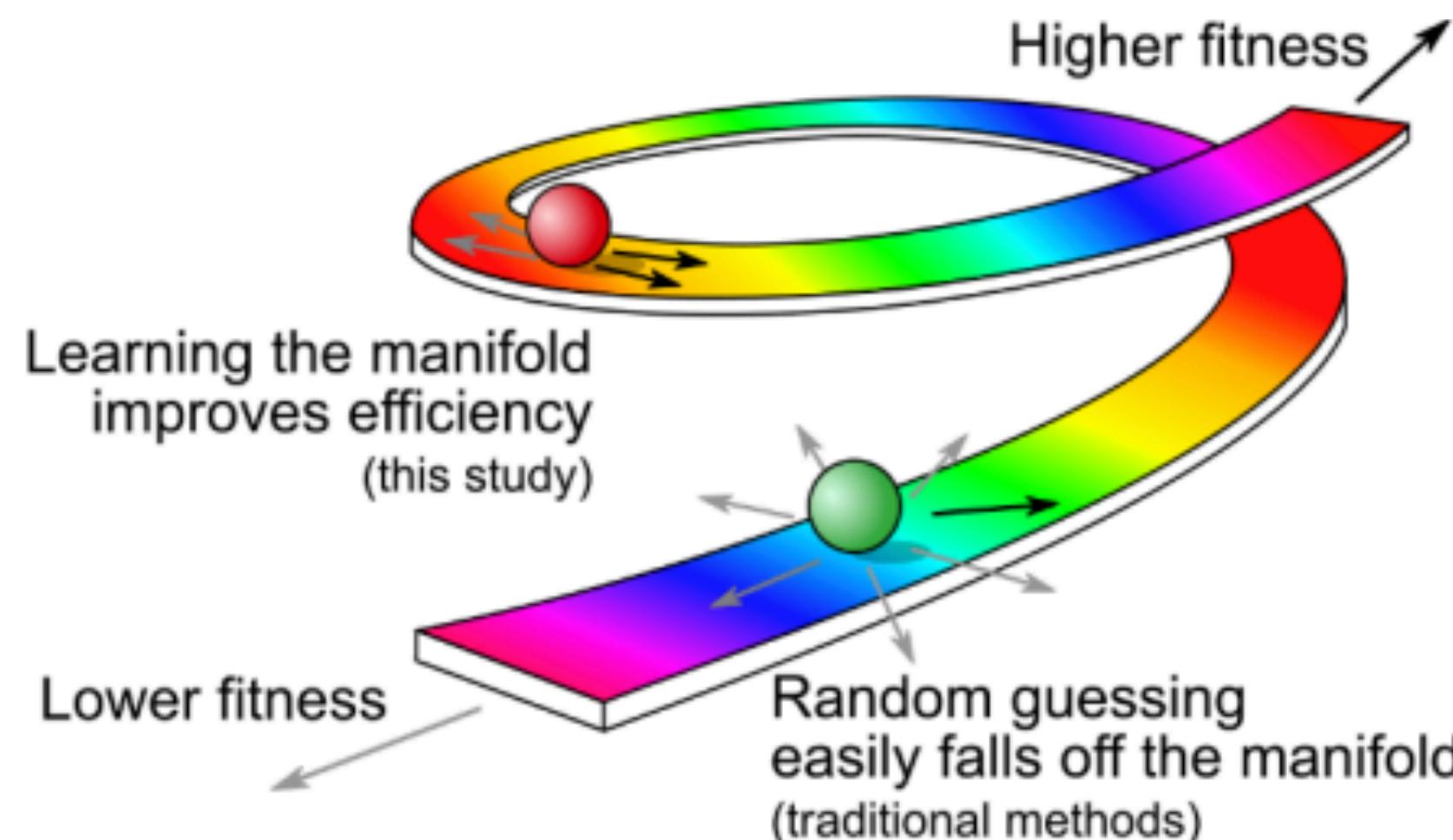


Hie et al. 2022

# Pretrained models propose high-fitness sequences

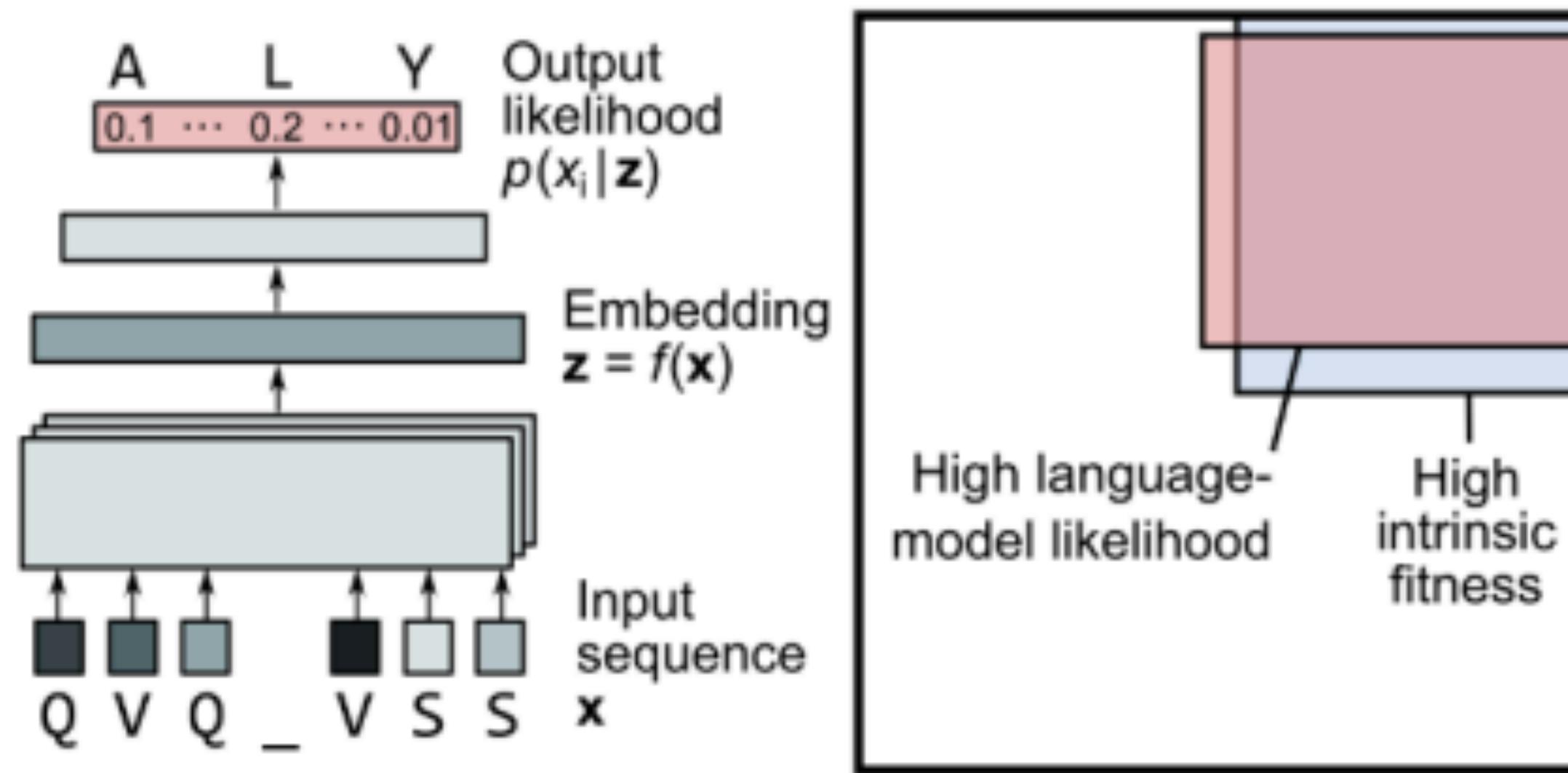


**Experiment:** Approximate intrinsic fitness with language-model likelihood, measure extrinsic fitness

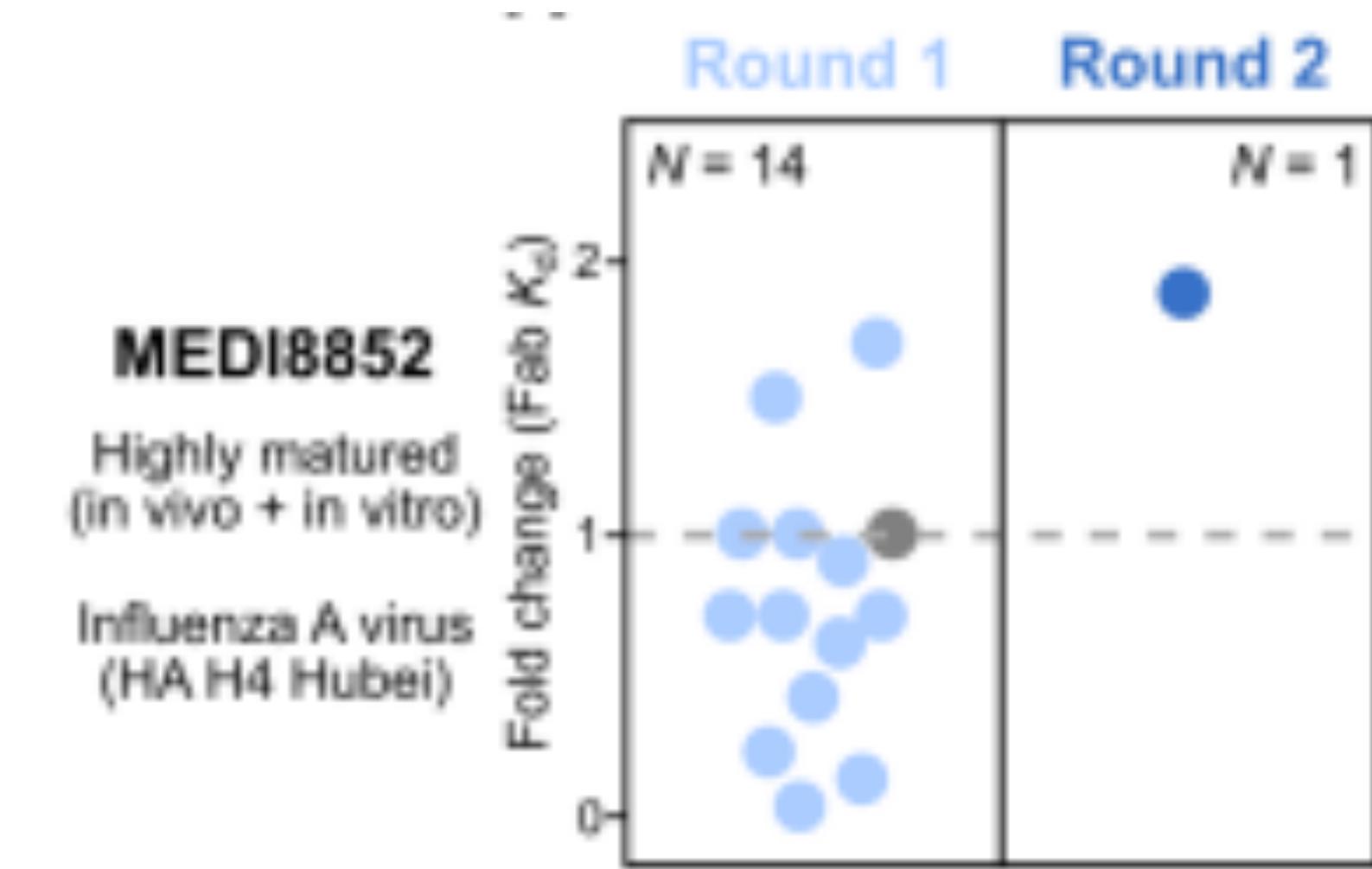
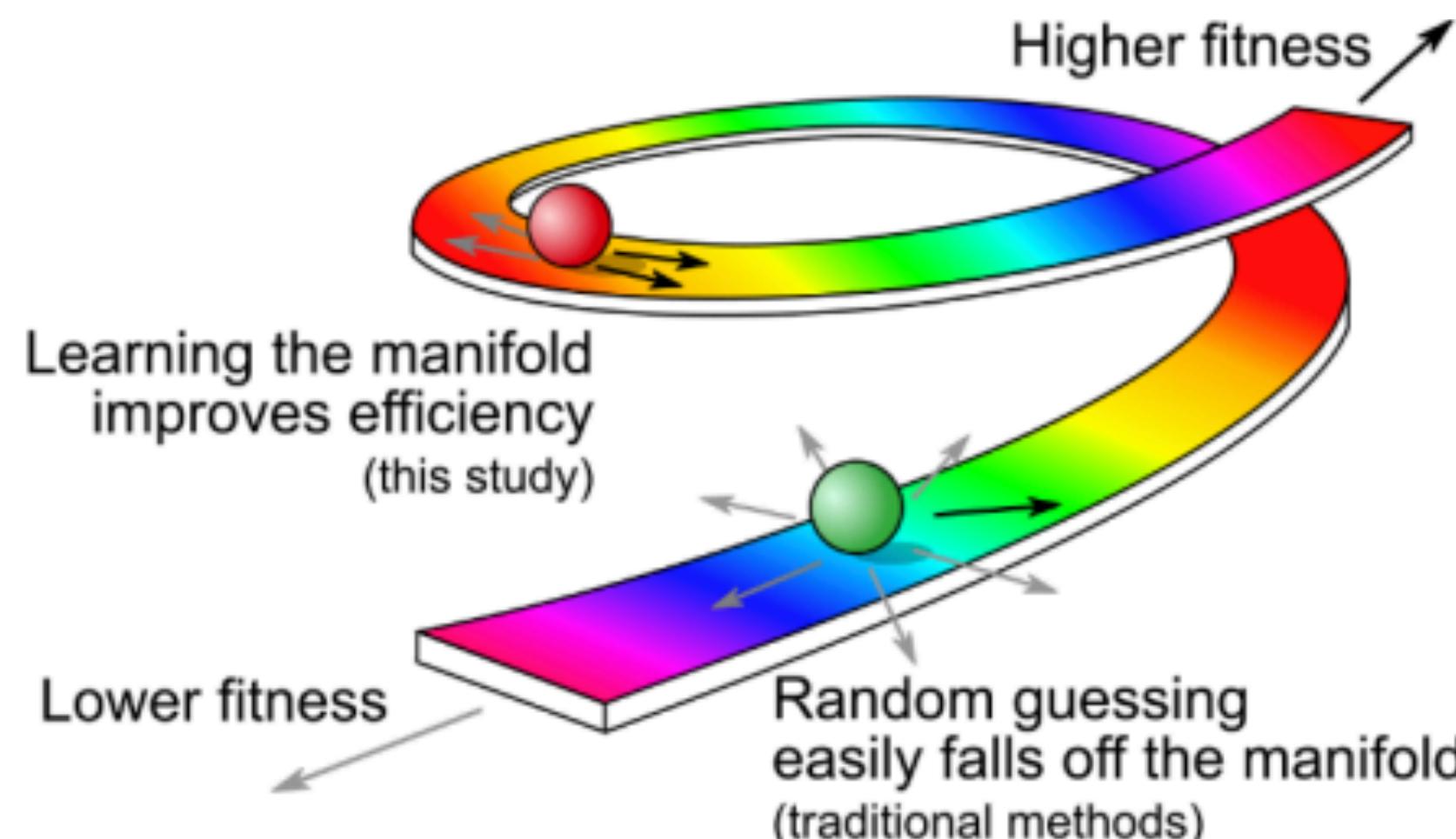


Hie et al. 2022

# Pretrained models propose high-fitness sequences

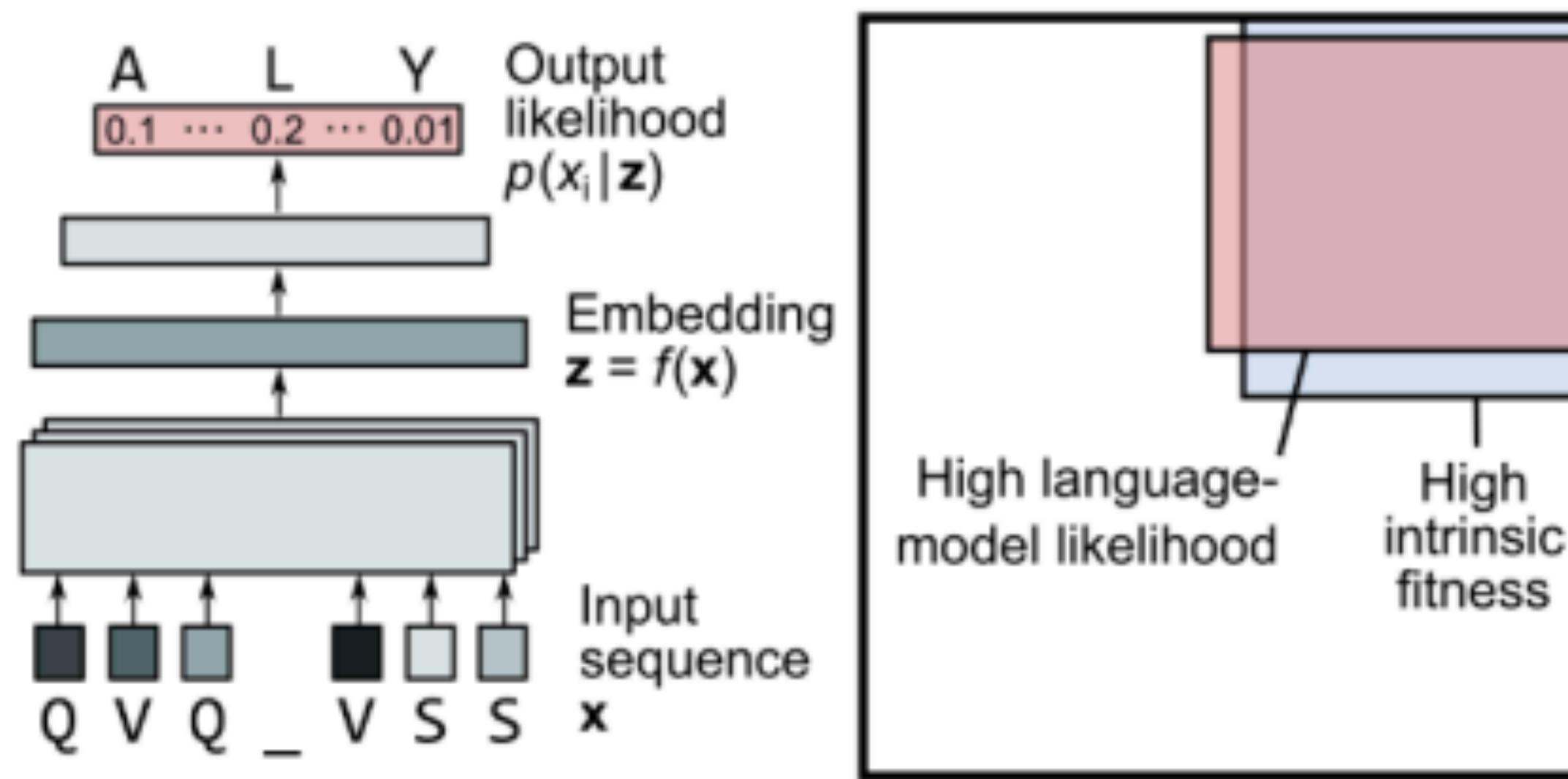


**Experiment:** Approximate intrinsic fitness with language-model likelihood, measure extrinsic fitness

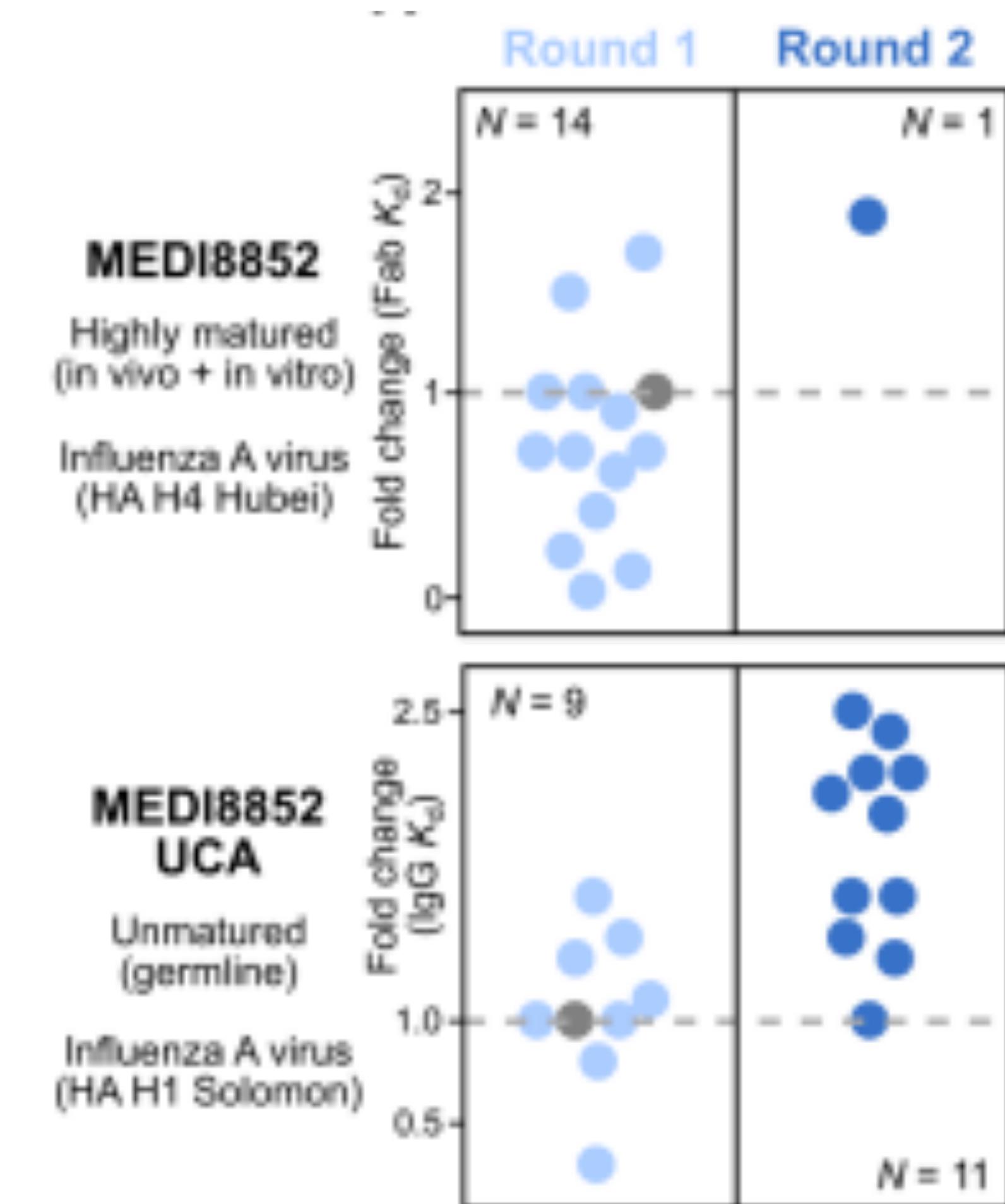
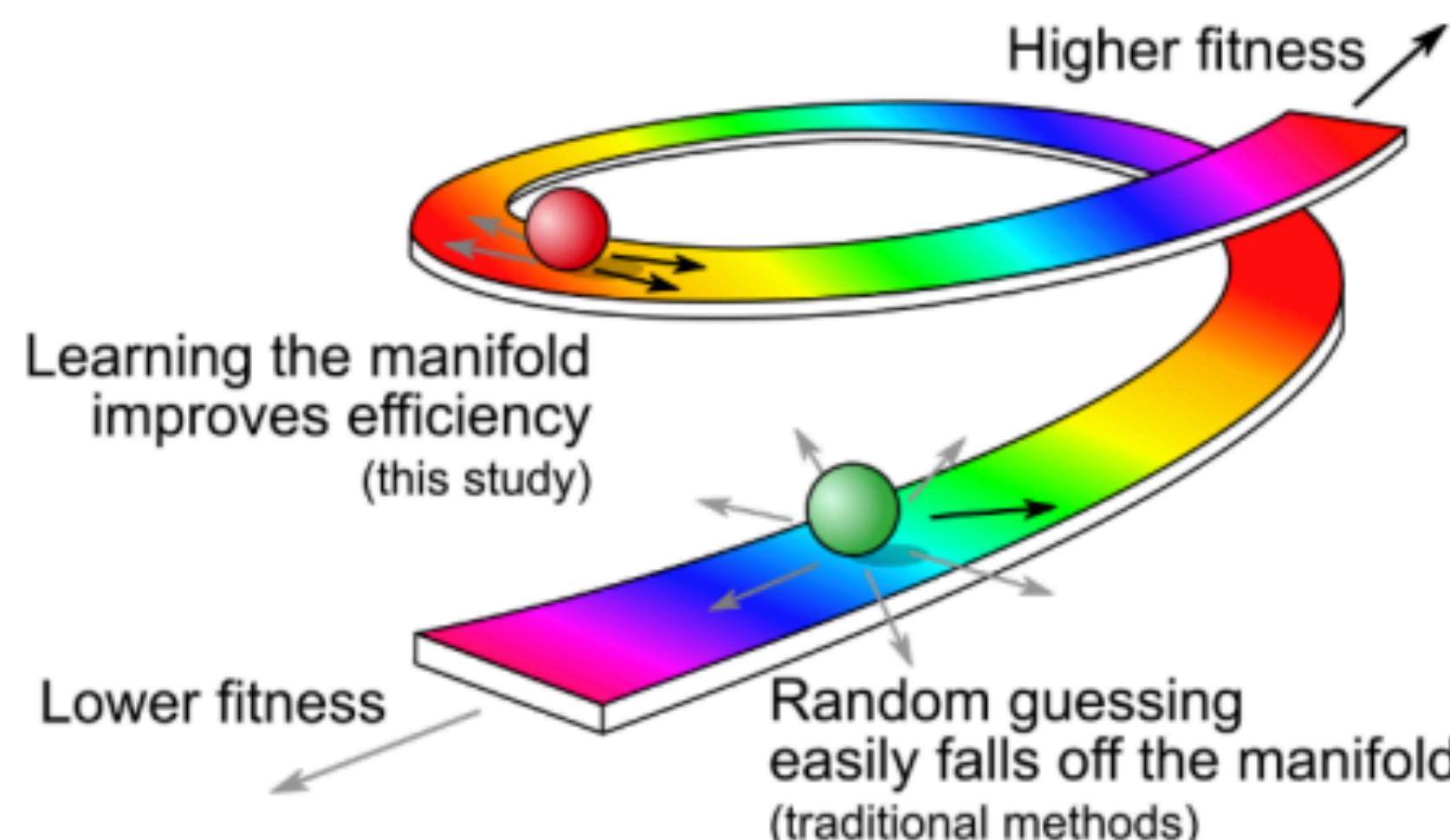


Hie et al. 2022

# Pretrained models propose high-fitness sequences



**Experiment:** Approximate intrinsic fitness with language-model likelihood, measure extrinsic fitness



Hie et al. 2022

# Need benchmarks for protein function ML

# Need benchmarks for protein function ML



# Need benchmarks for protein function ML

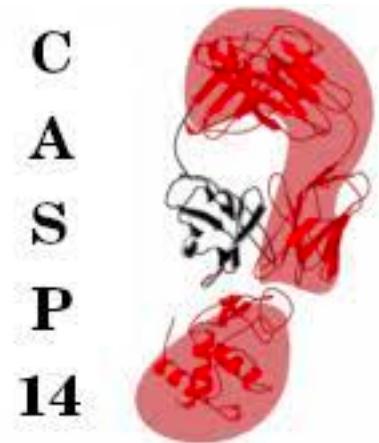
IMAGENET

SuperGLUE

# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



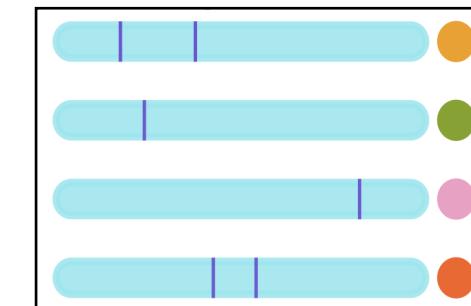
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Focus on sequence -> fitness



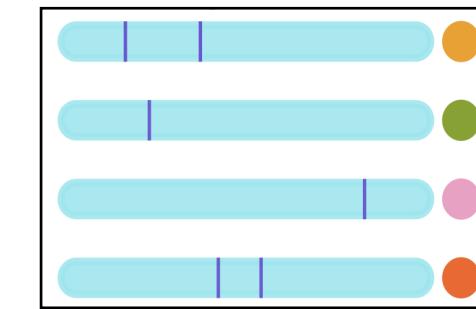
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Focus on sequence -> fitness

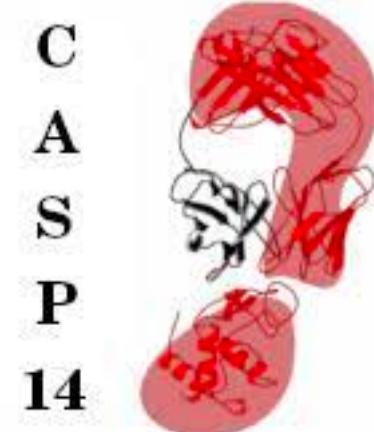


TAPE function datasets do not discriminate well

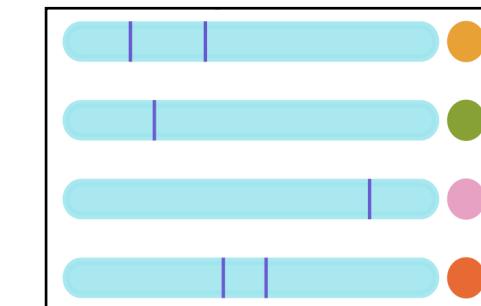
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Focus on sequence -> fitness



TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
-------	-------------	--------------	-----------

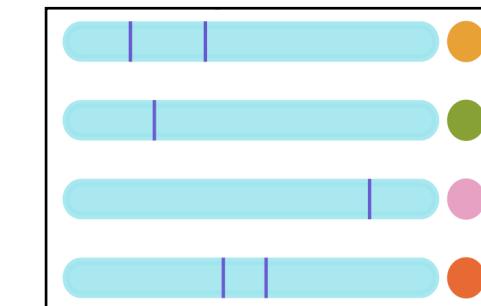
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Focus on sequence -> fitness



TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71

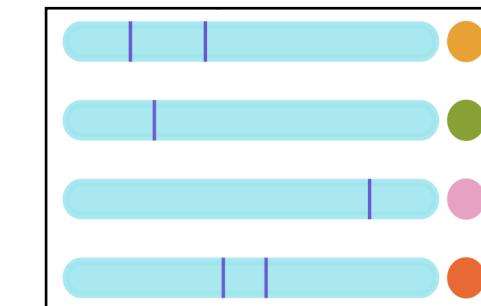
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Focus on sequence -> fitness



TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73

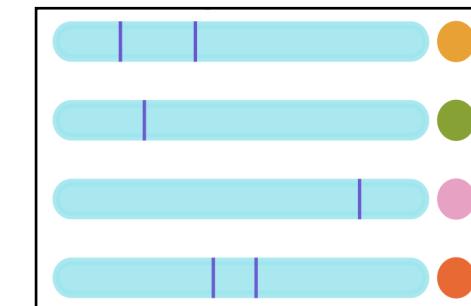
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Focus on sequence -> fitness



TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65

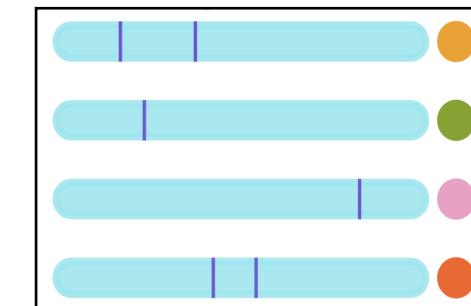
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Focus on sequence -> fitness



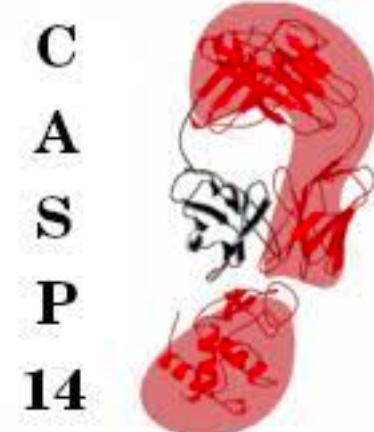
TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65
Linear regression	none	0.68	0.48

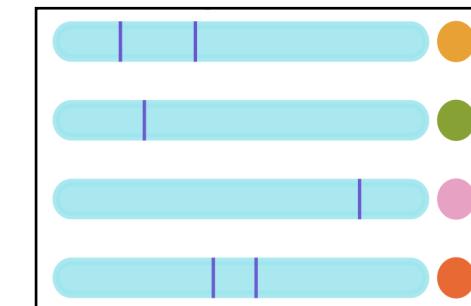
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Focus on sequence -> fitness



TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65
Linear regression	none	0.68	0.48
CNN	none	0.67	0.51

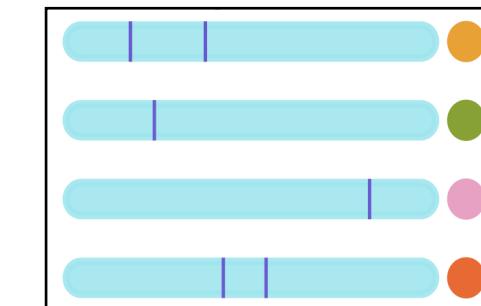
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Focus on sequence -> fitness



TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65
Linear regression	none	0.68	0.48
CNN	none	0.67	0.51
CARP-640M	MLM	0.68	0.72

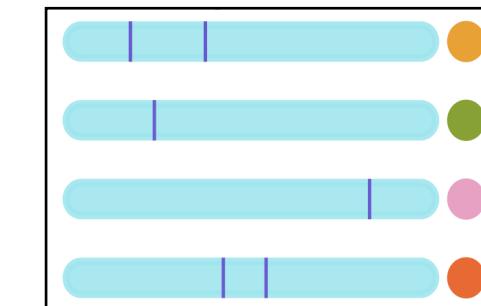
# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



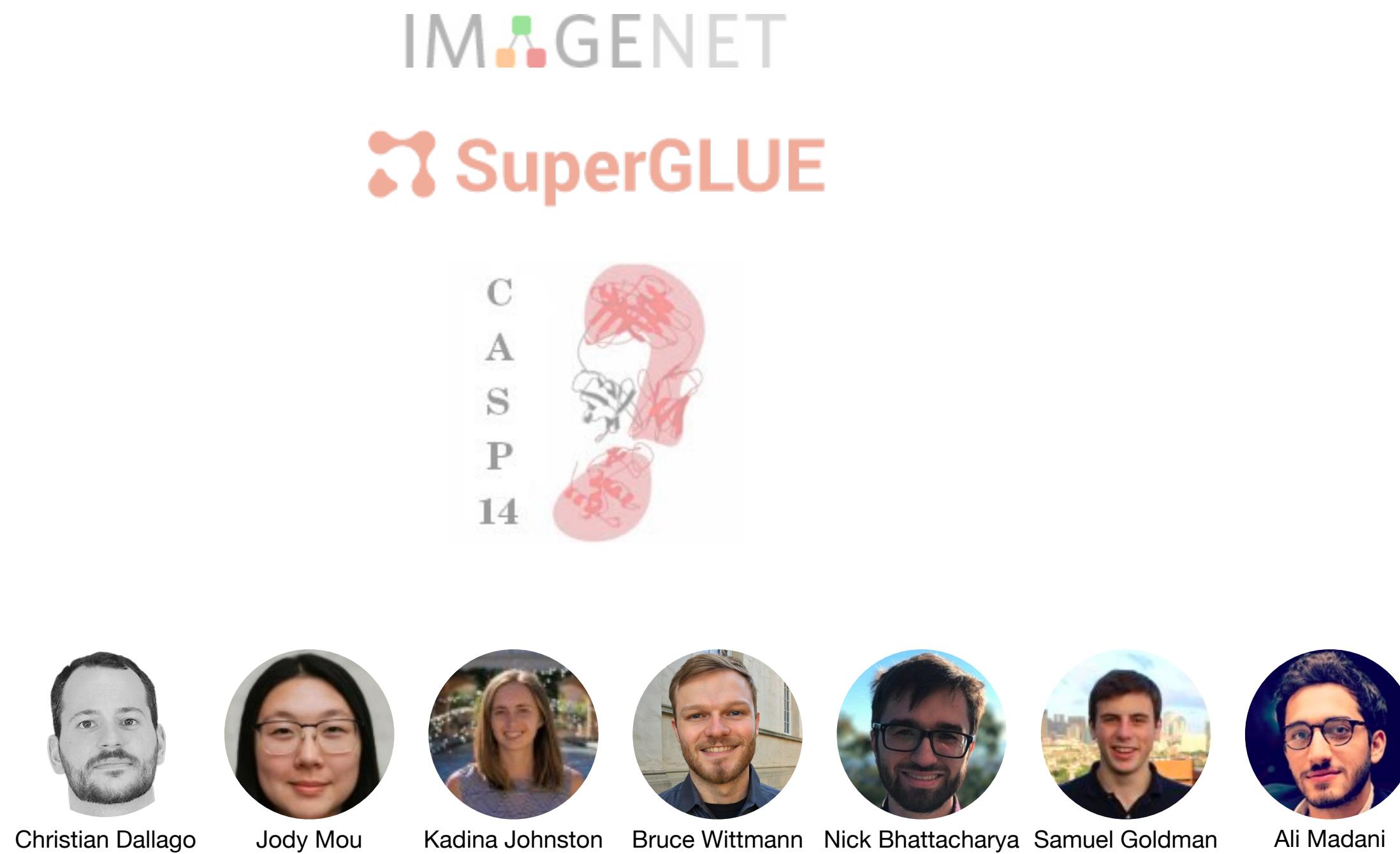
Focus on sequence -> fitness



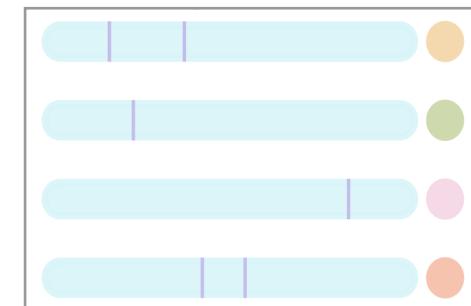
TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65
Linear regression	none	0.68	0.48
CNN	none	0.67	0.51
CARP-640M	MLM	0.68	0.72

# Need benchmarks for protein function ML



Focus on sequence -> fitness



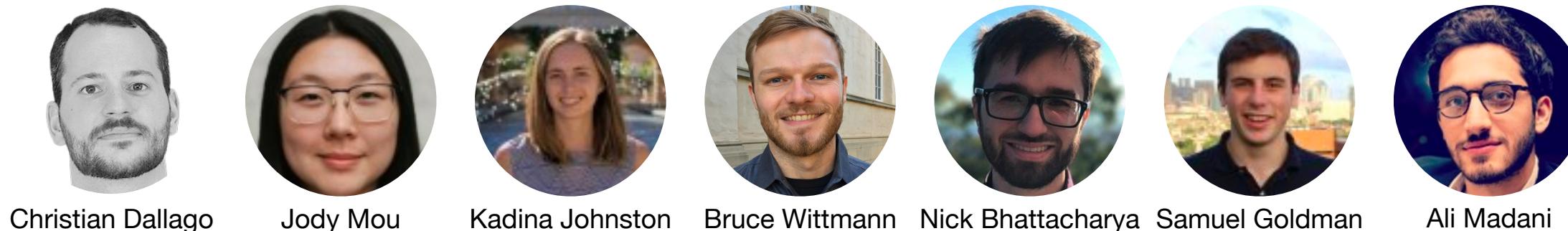
TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65
Linear regression	none	0.68	0.48
CNN	none	0.67	0.51
CARP-640M	MLM	0.68	0.72

# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Christian Dallago

Jody Mou

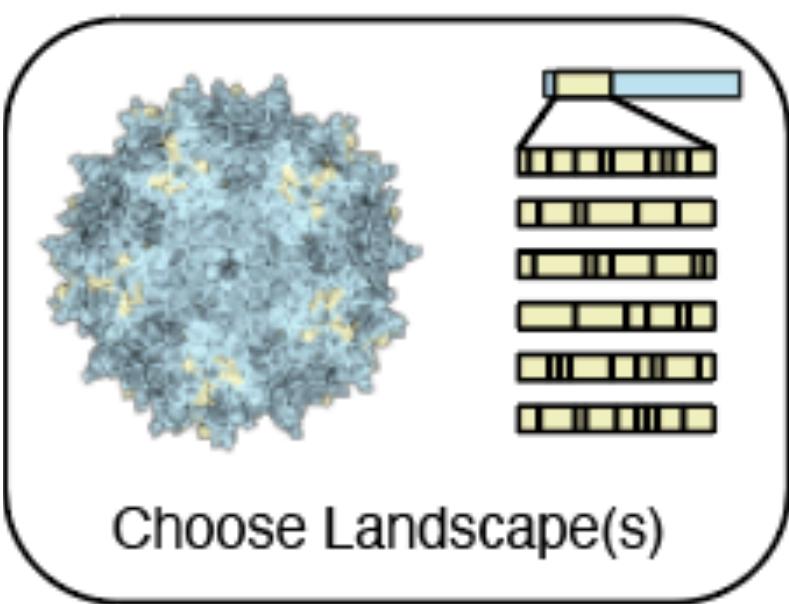
Kadina Johnston

Bruce Wittmann

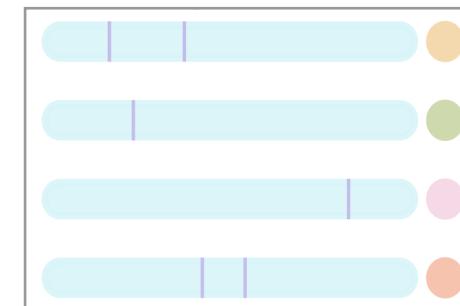
Nick Bhattacharya

Samuel Goldman

Ali Madani



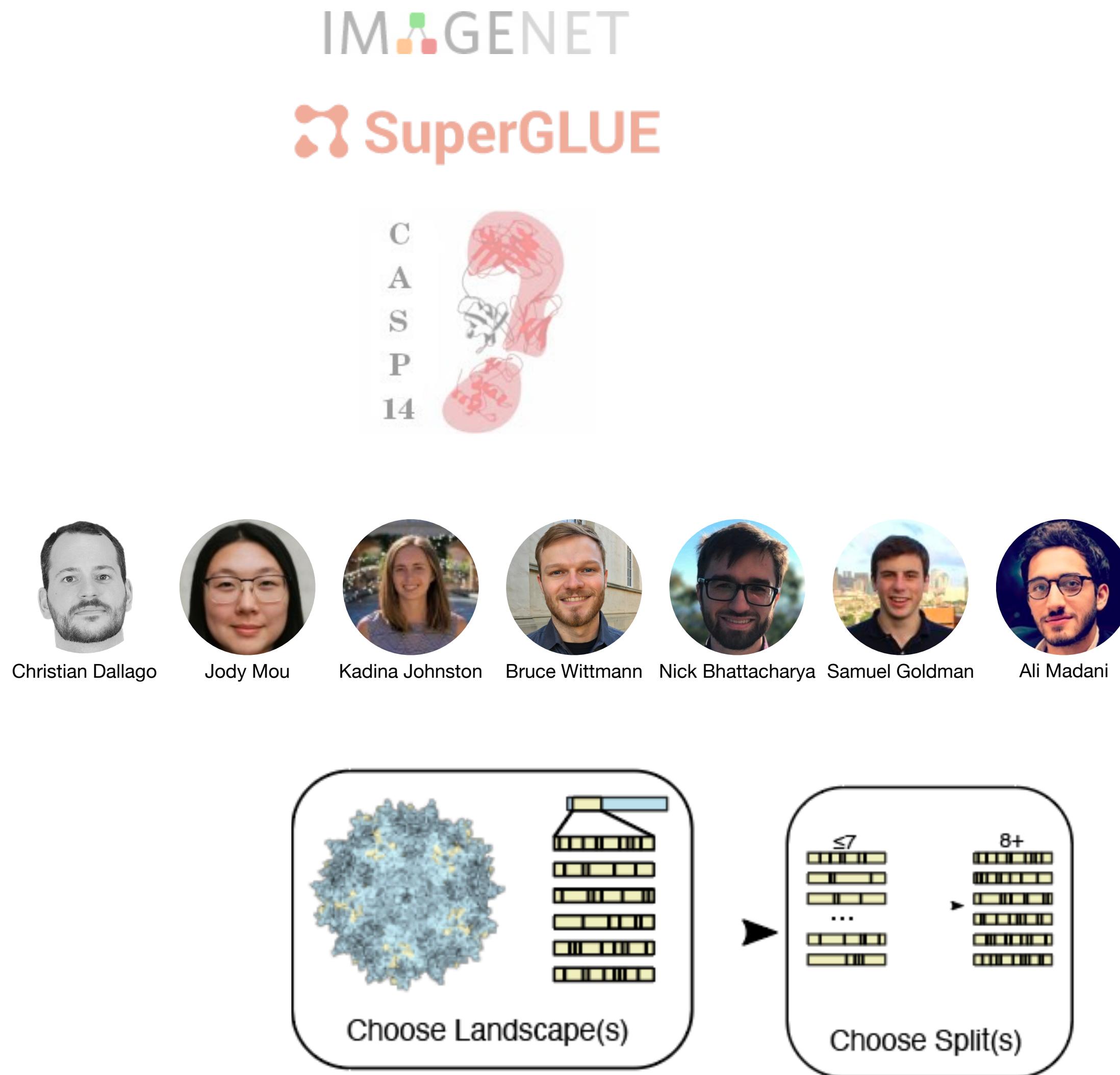
Focus on sequence -> fitness



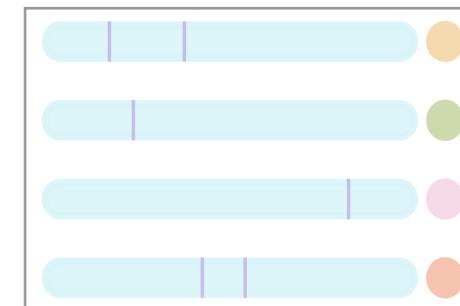
TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65
Linear regression	none	0.68	0.48
CNN	none	0.67	0.51
CARP-640M	MLM	0.68	0.72

# Need benchmarks for protein function ML

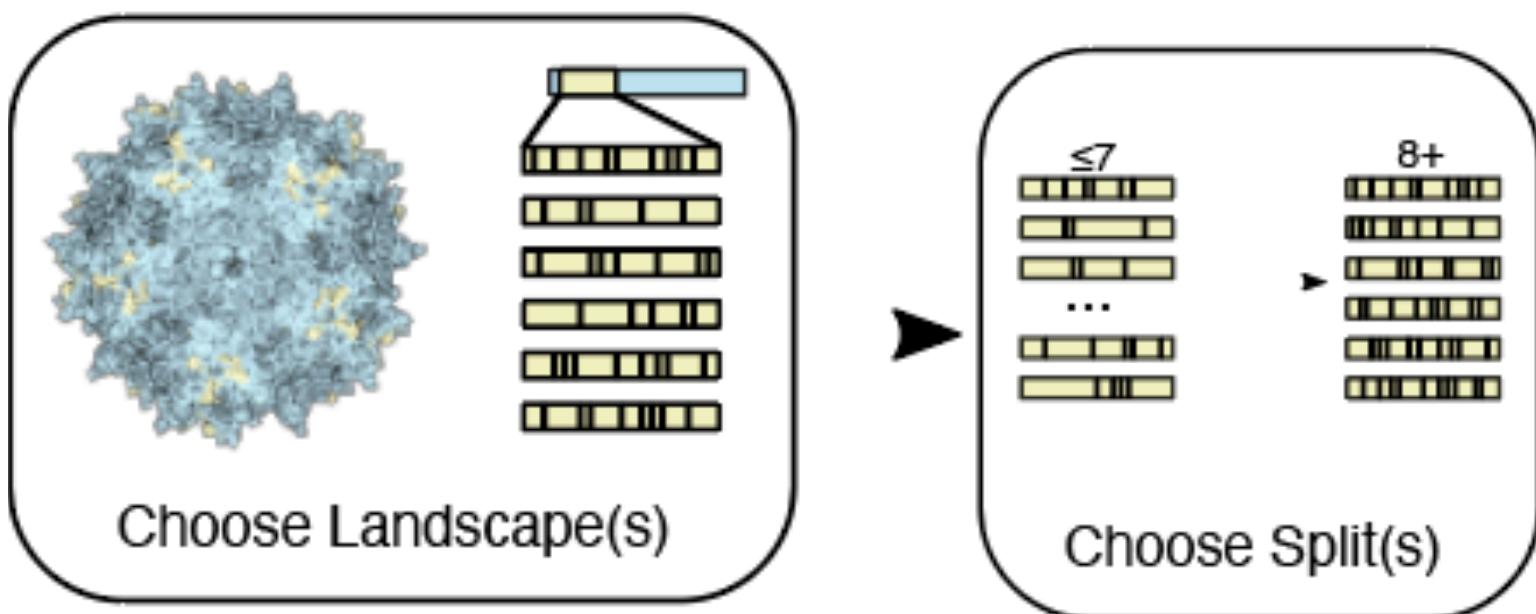


Focus on sequence -> fitness



TAPE function datasets do not discriminate well

Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65
Linear regression	none	0.68	0.48
CNN	none	0.67	0.51
CARP-640M	MLM	0.68	0.72



# Need benchmarks for protein function ML

IMAGENET

SuperGLUE



Christian Dallago

Jody Mou

Kadina Johnston

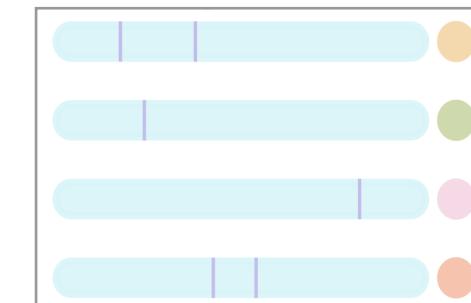
Bruce Wittmann

Nick Bhattacharya

Samuel Goldman

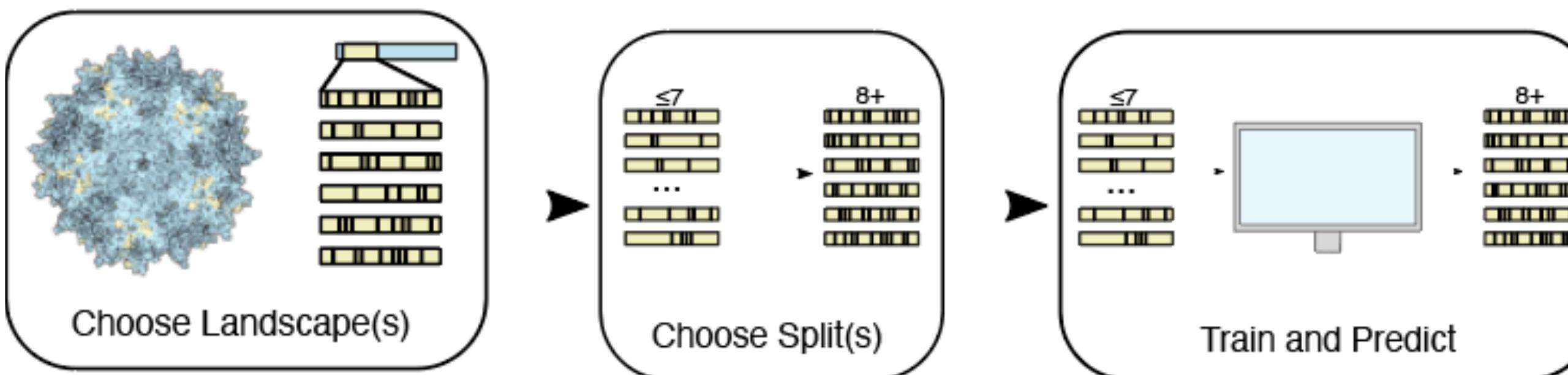
Ali Madani

Focus on sequence -> fitness

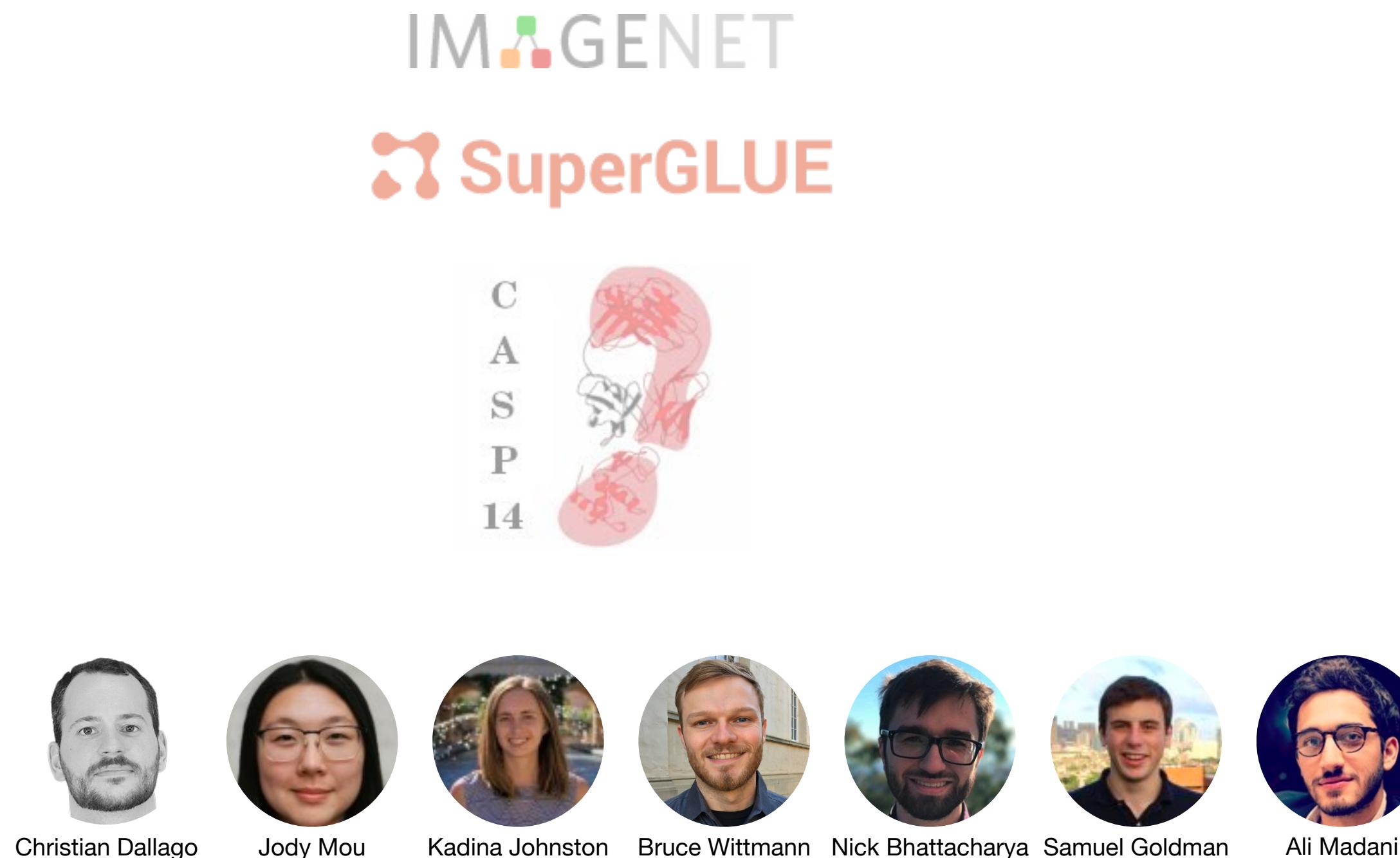


TAPE function datasets do not discriminate well

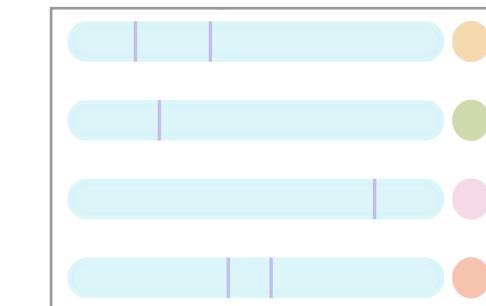
Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65
Linear regression	none	0.68	0.48
CNN	none	0.67	0.51
CARP-640M	MLM	0.68	0.72



# Need benchmarks for protein function ML

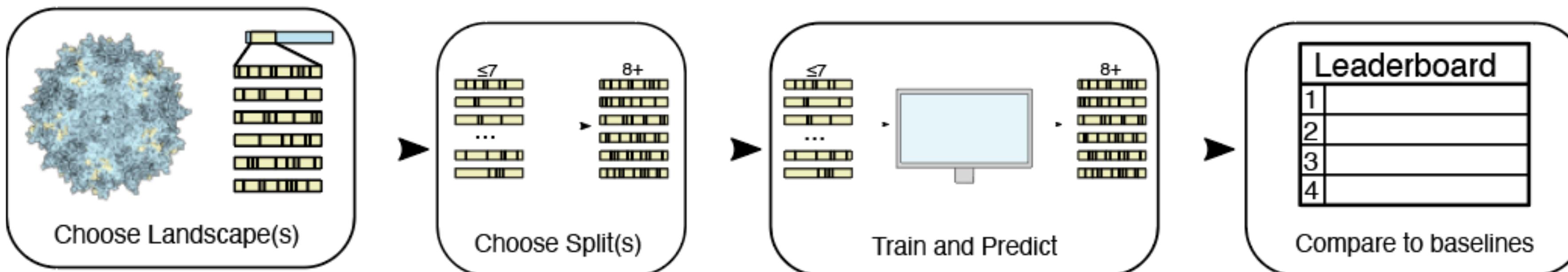


Focus on sequence -> fitness

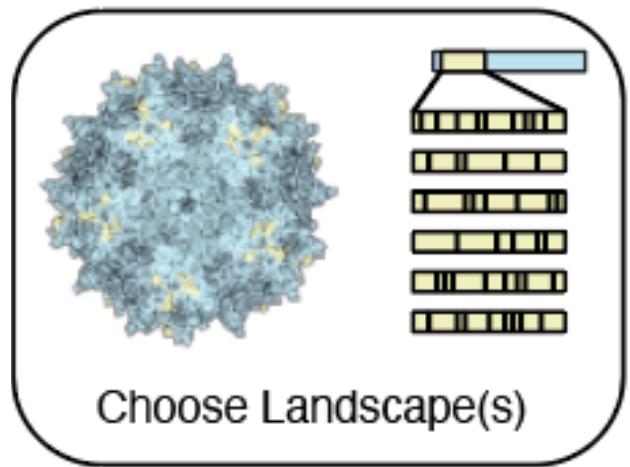


TAPE function datasets do not discriminate well

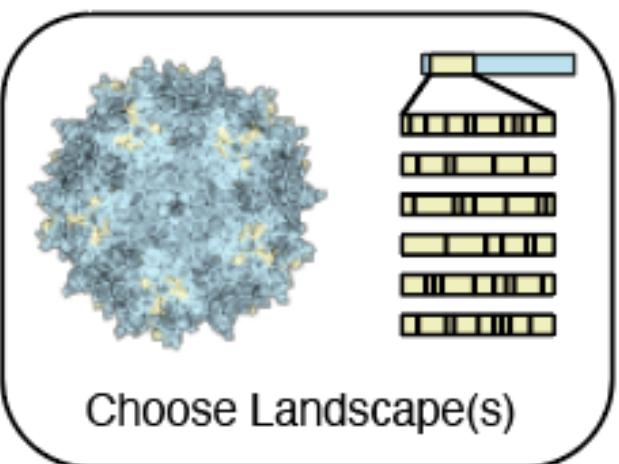
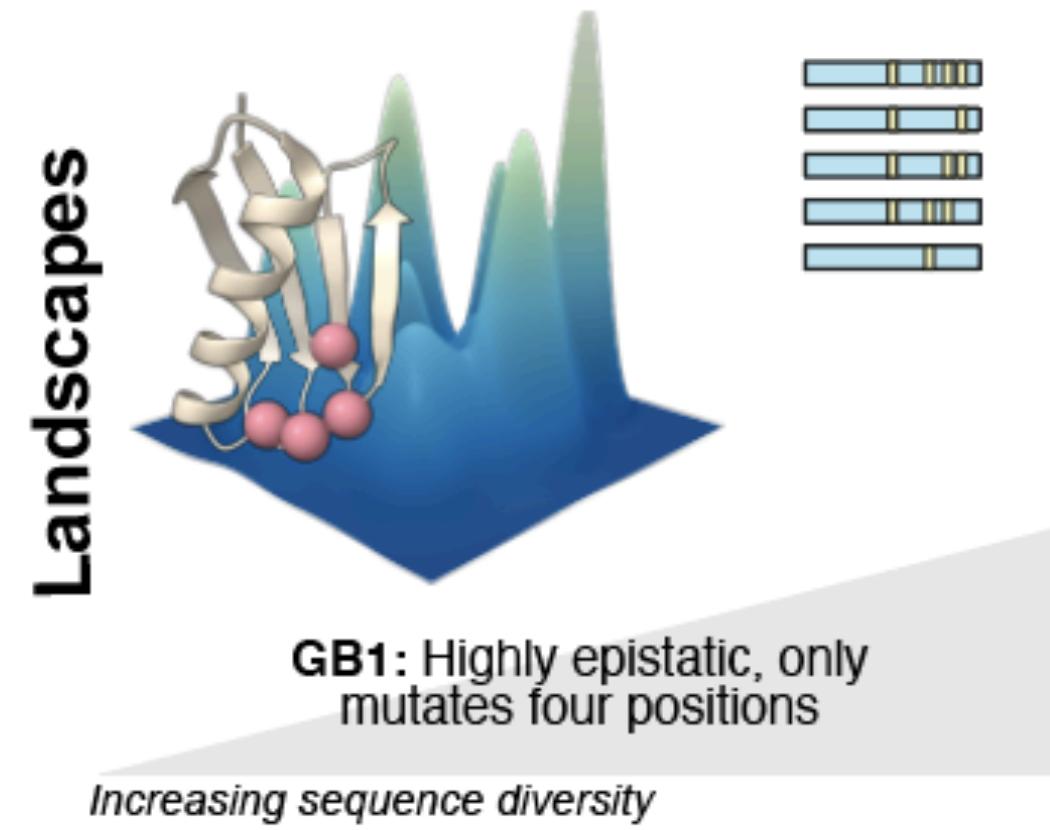
Model	Pretraining	Fluorescence	Stability
ESM (Rives <i>et al.</i> )	MLM	0.68	0.71
TAPE transformer	MLM	0.68	0.73
CPCProt (Lu <i>et al.</i> )	contrastive	0.68	0.65
Linear regression	none	0.68	0.48
CNN	none	0.67	0.51
CARP-640M	MLM	0.68	0.72



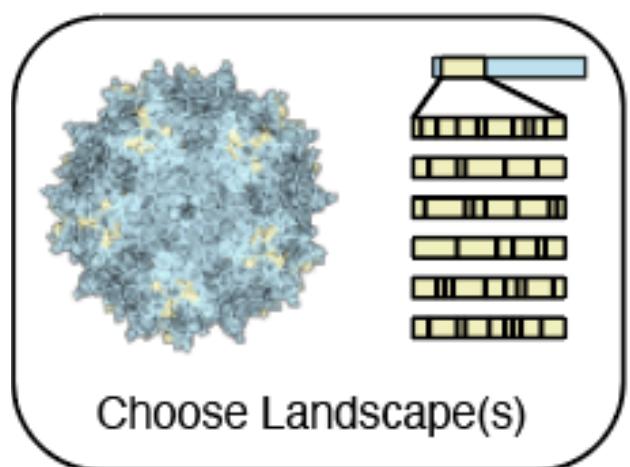
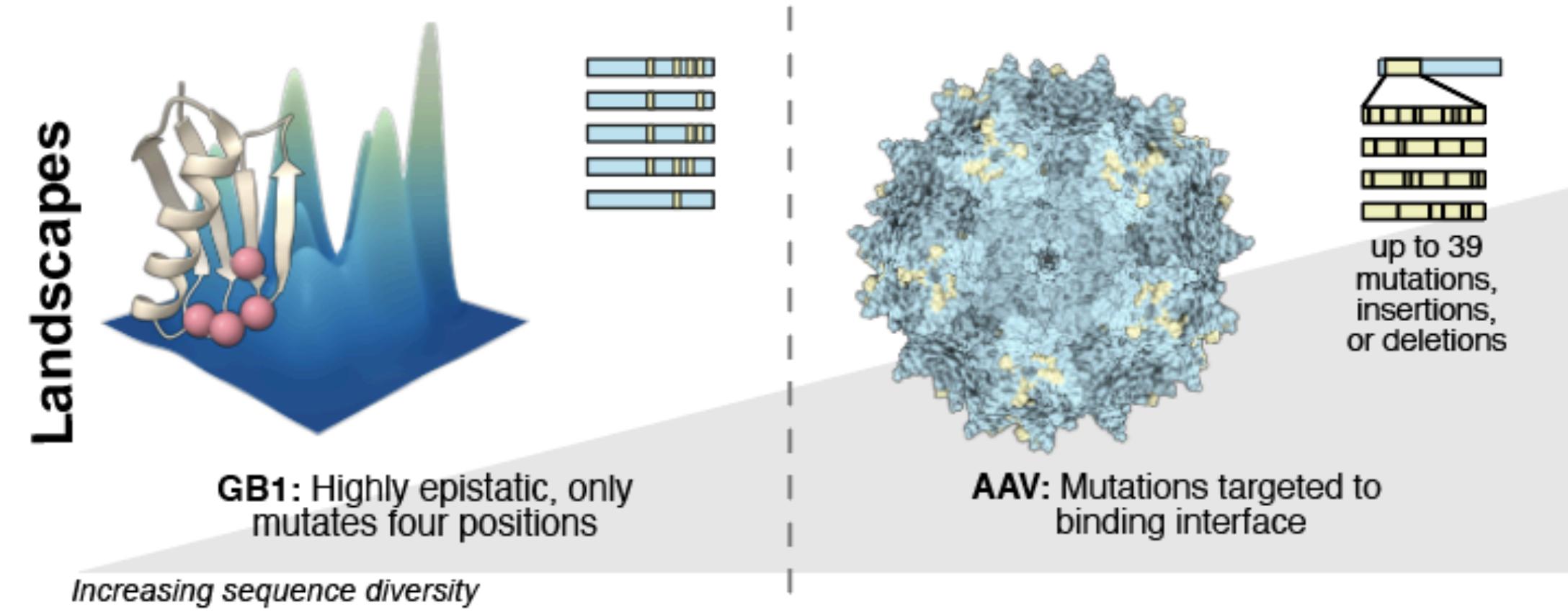
# Choose three landscapes



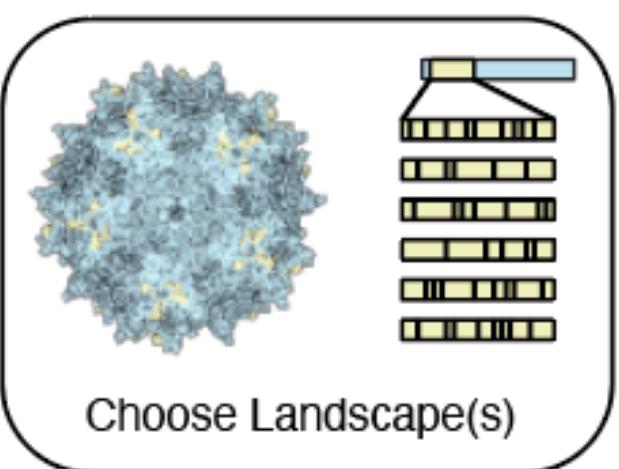
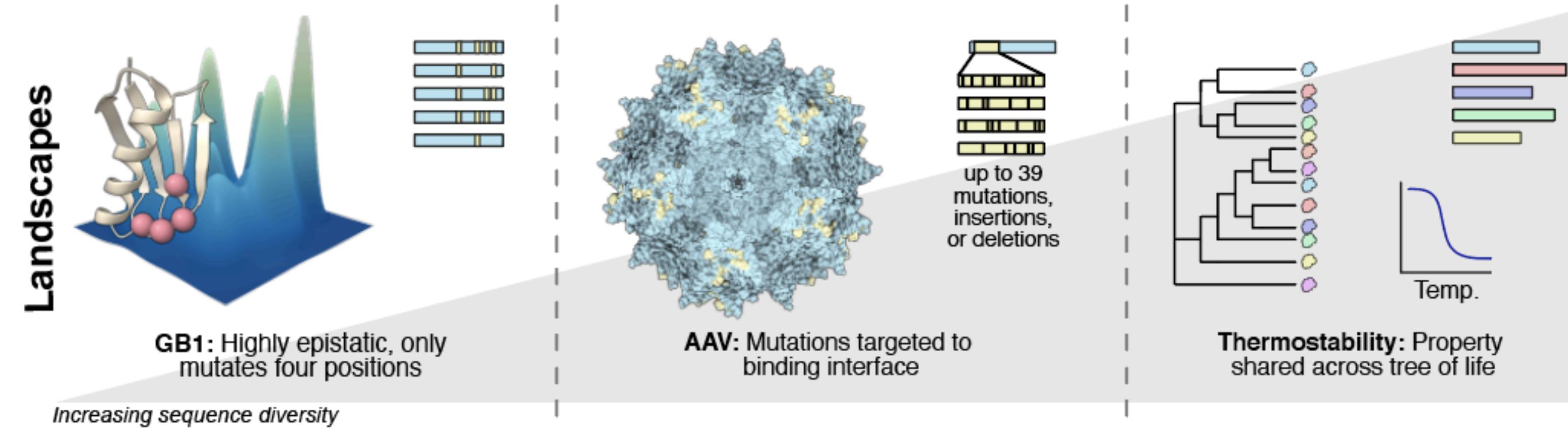
# Choose three landscapes



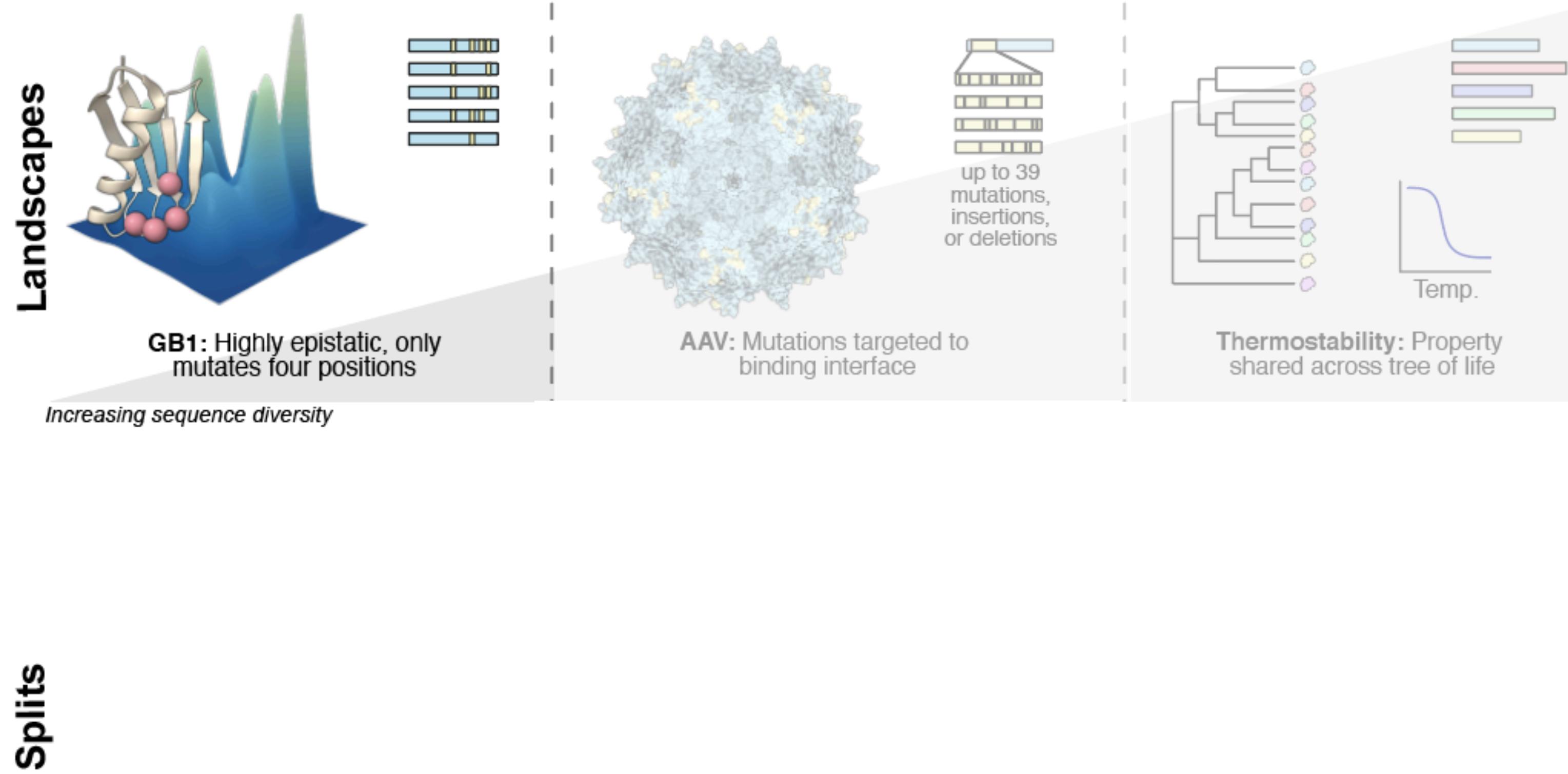
# Choose three landscapes



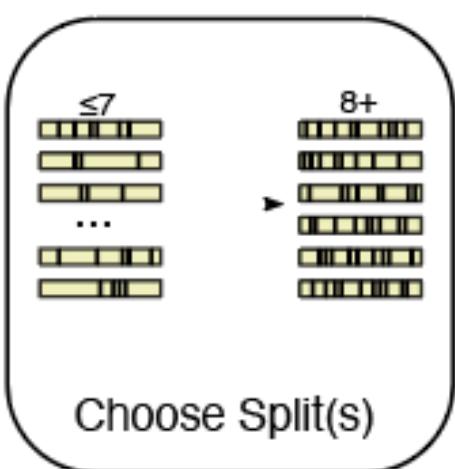
# Choose three landscapes



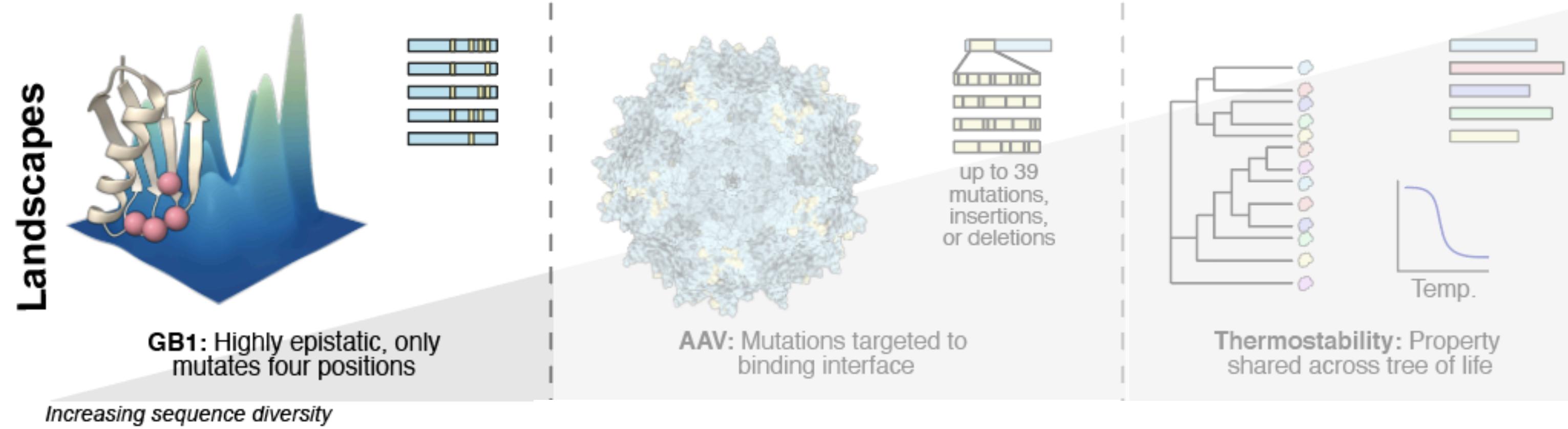
# Make biologically-relevant train/test splits



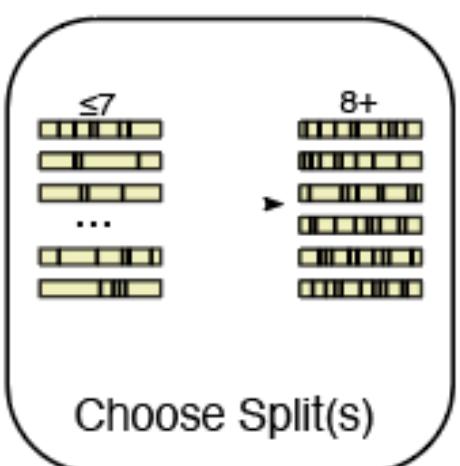
**Splits**



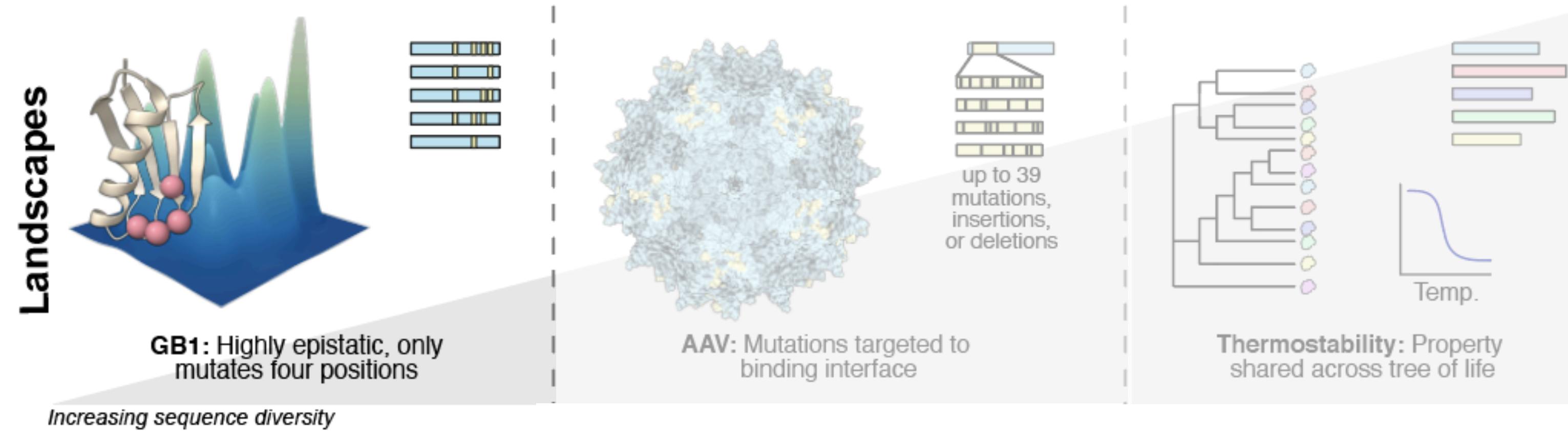
# Make biologically-relevant train/test splits



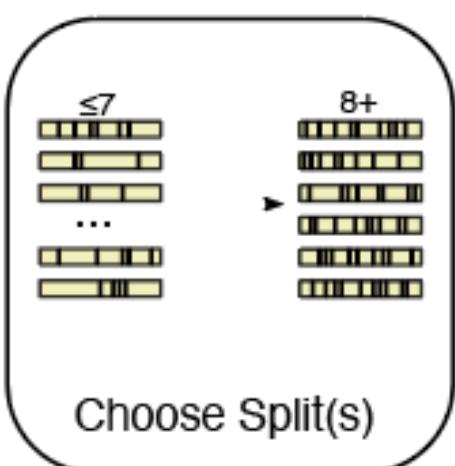
**Splits**



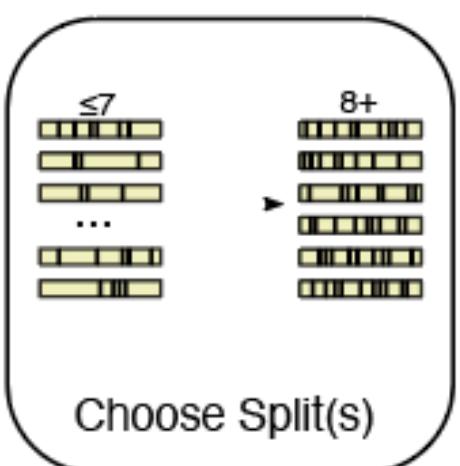
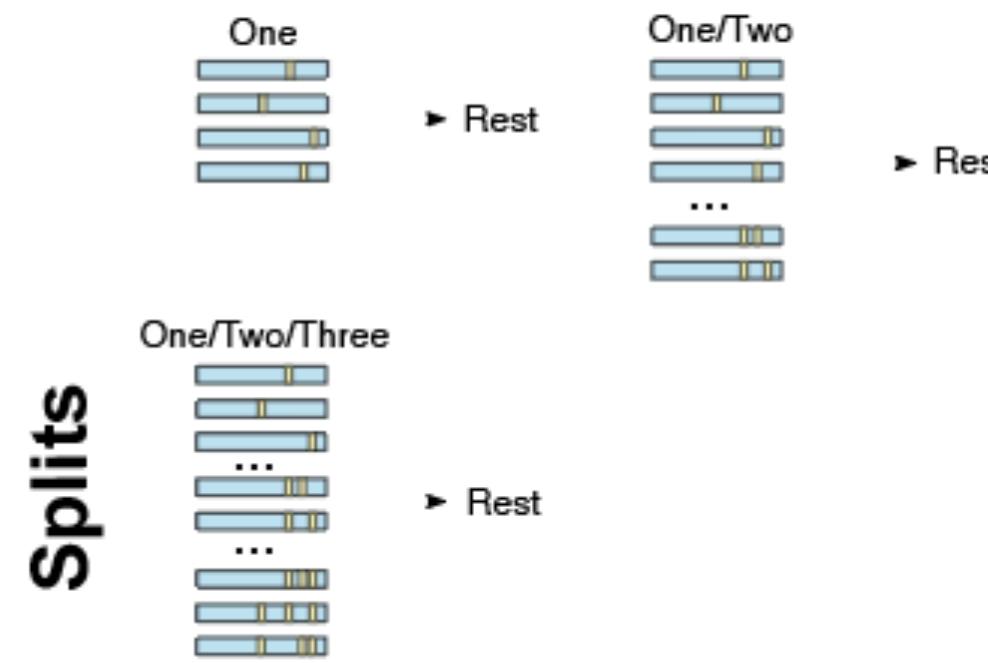
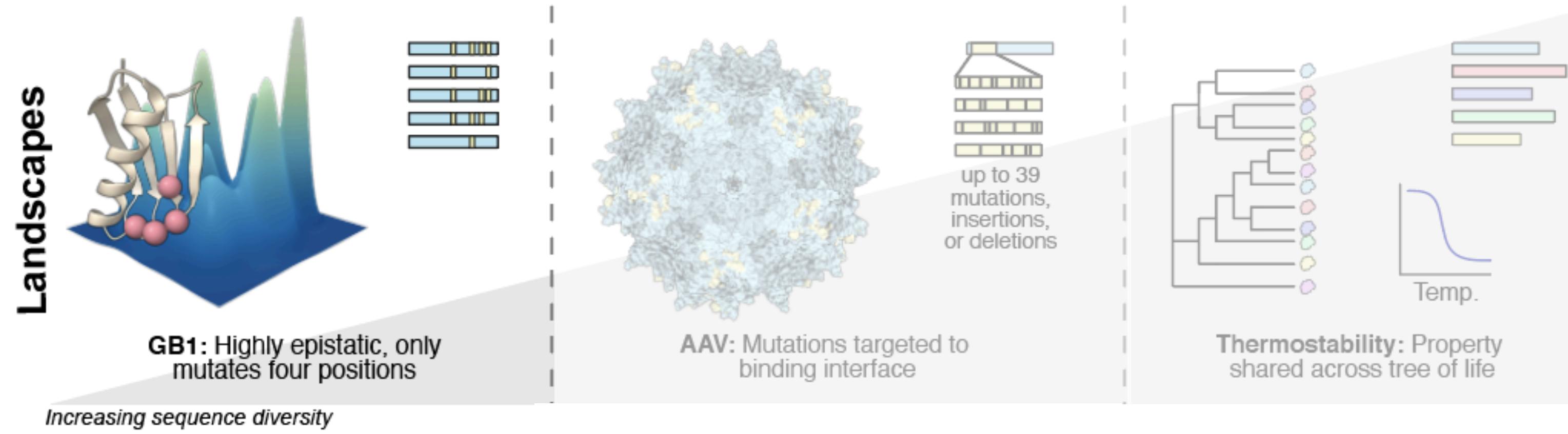
# Make biologically-relevant train/test splits



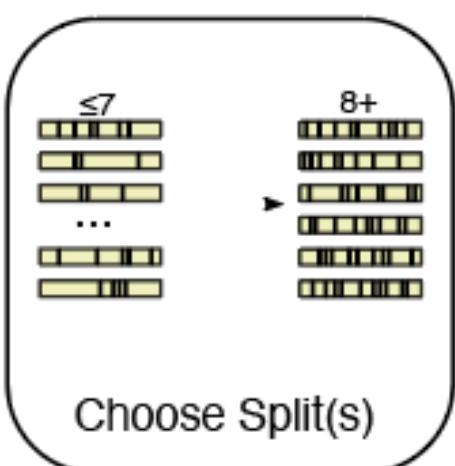
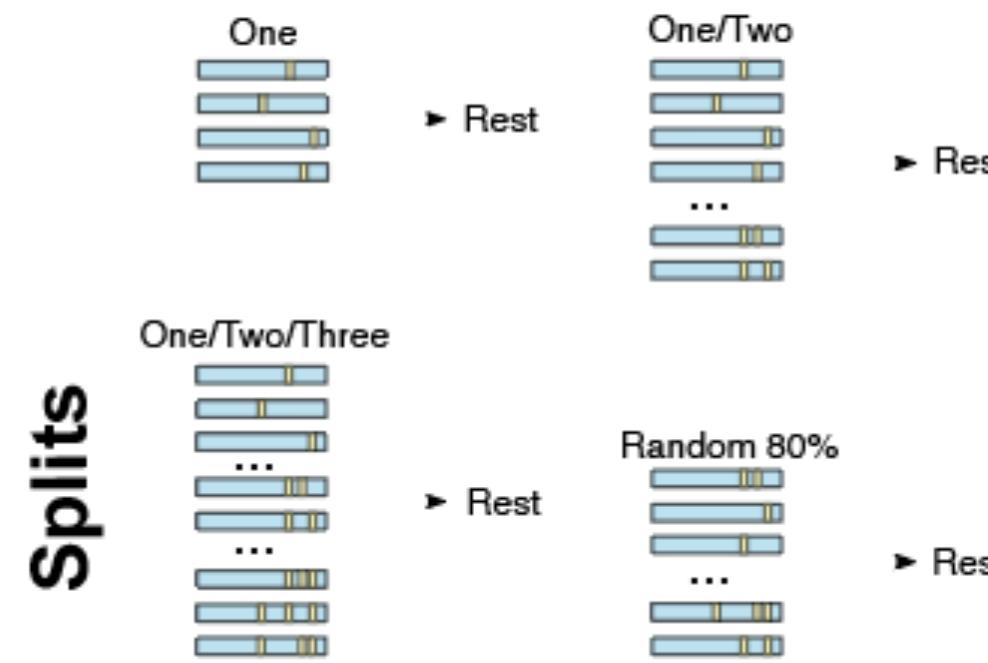
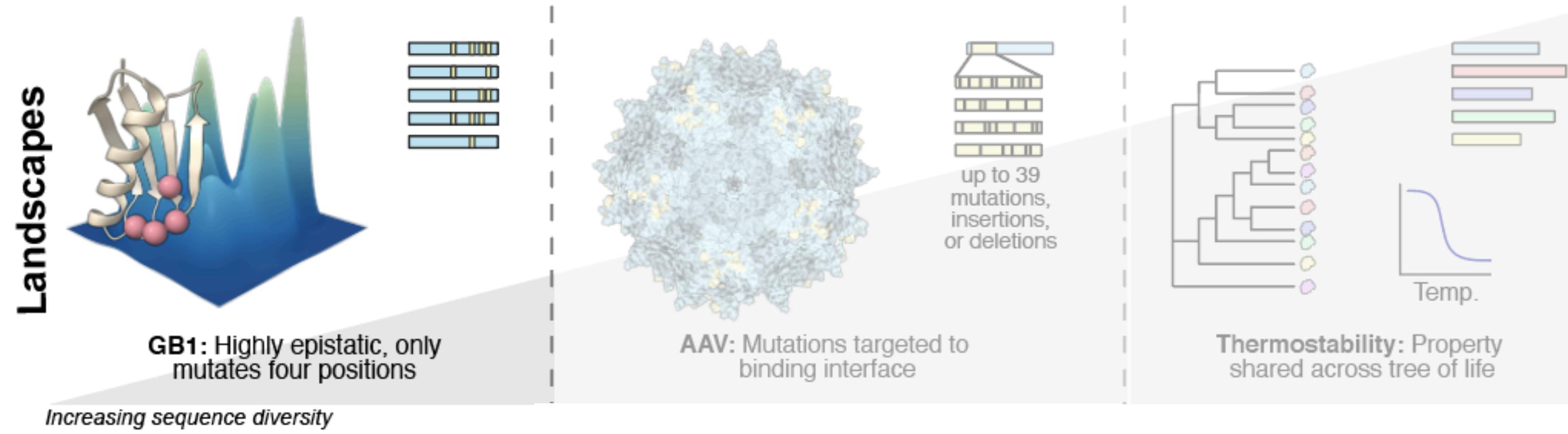
**Splits**



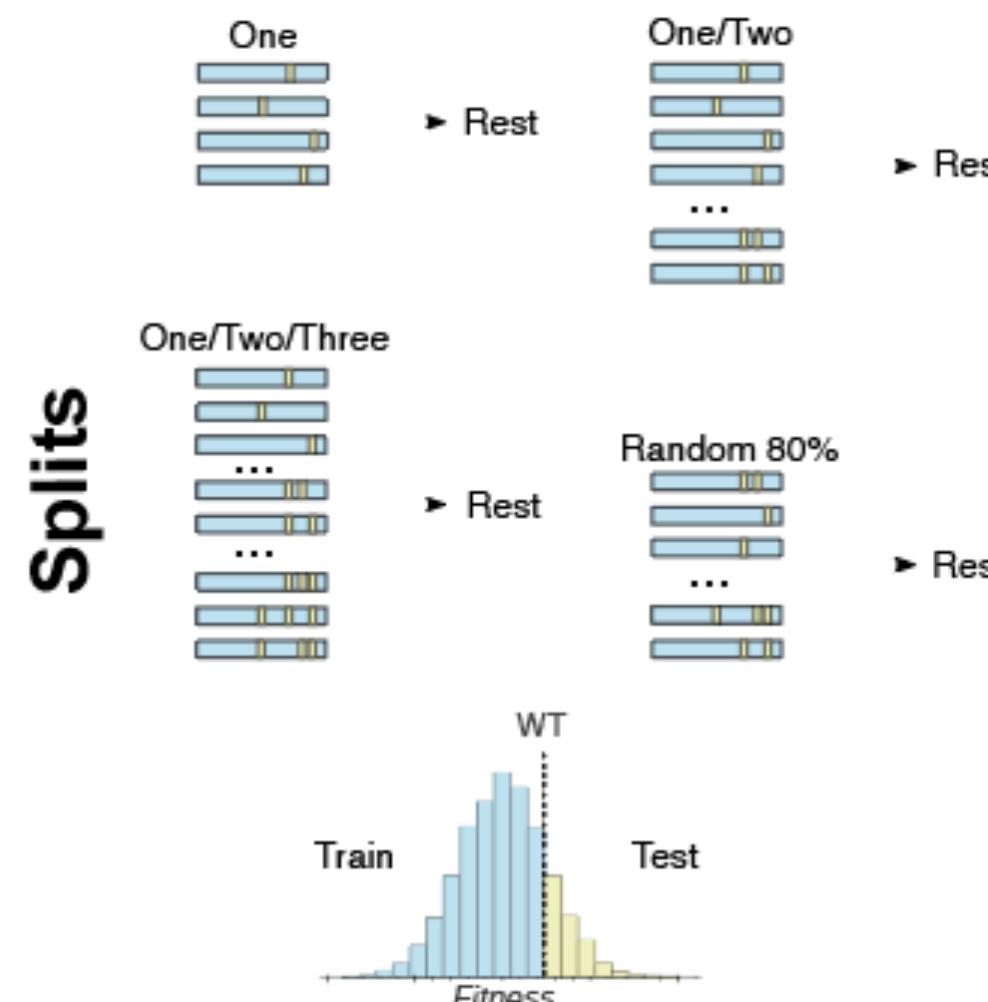
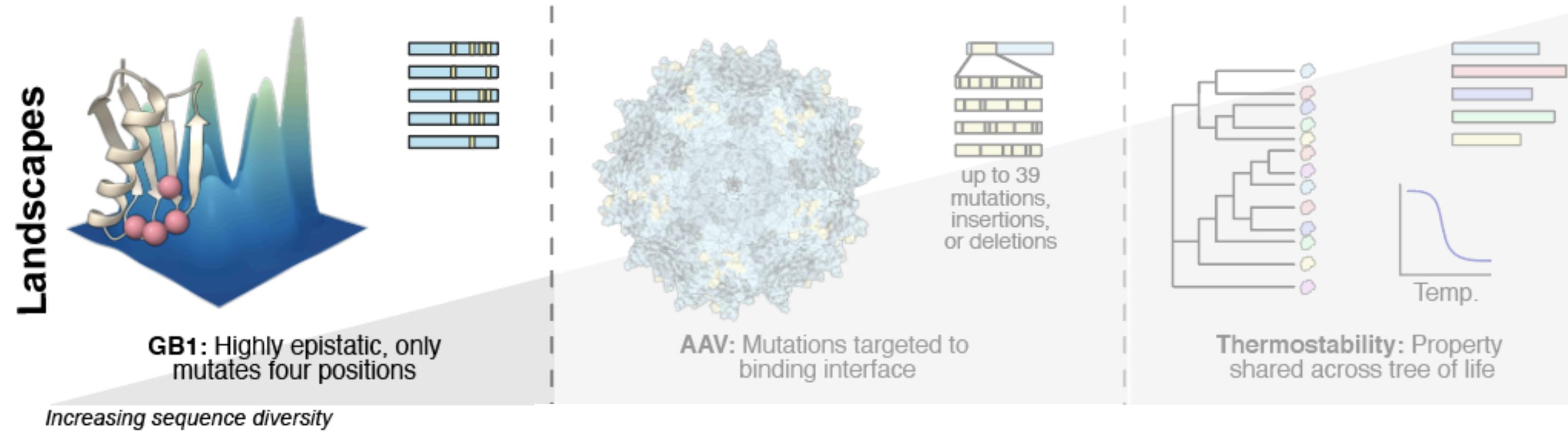
# Make biologically-relevant train/test splits



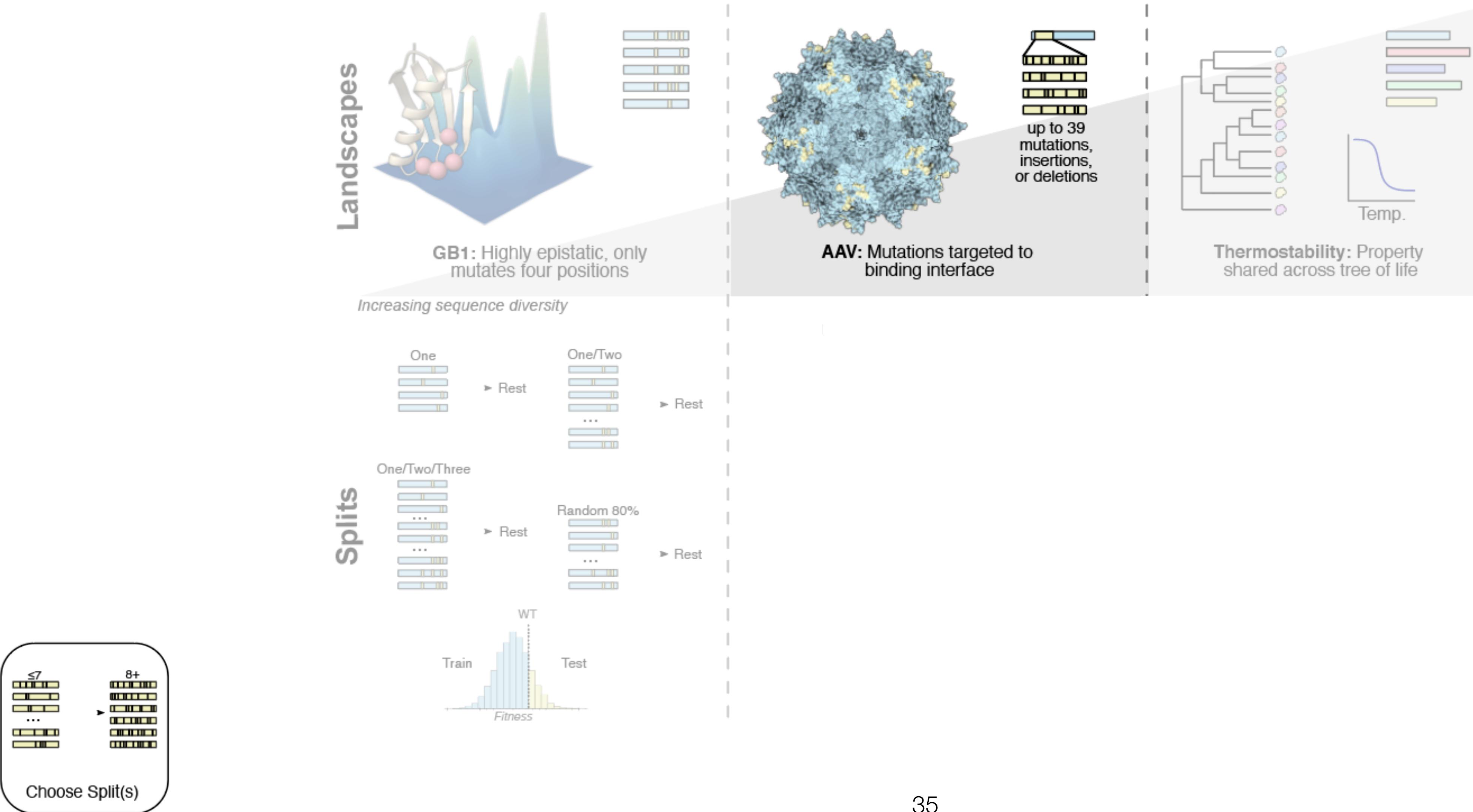
# Make biologically-relevant train/test splits



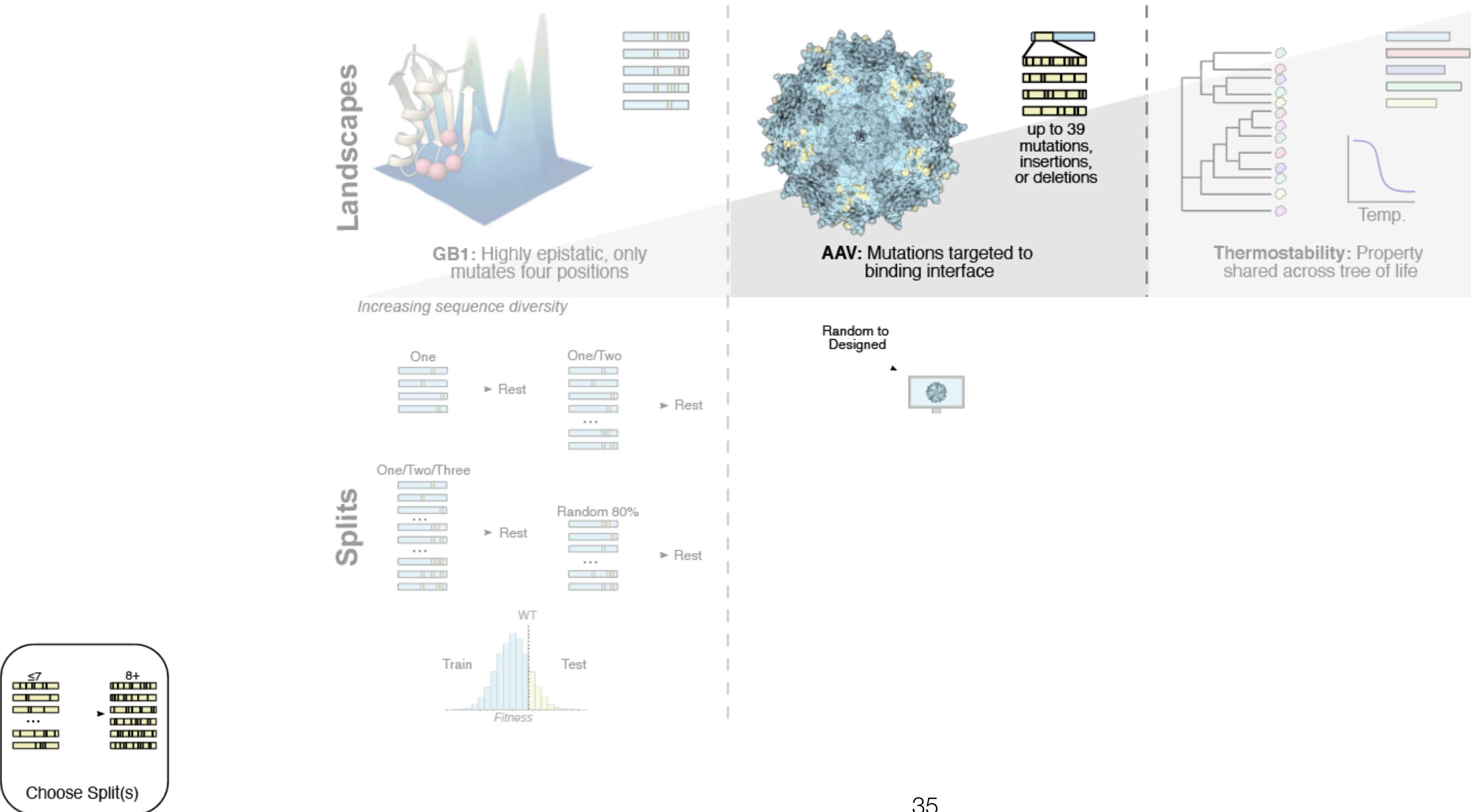
# Make biologically-relevant train/test splits



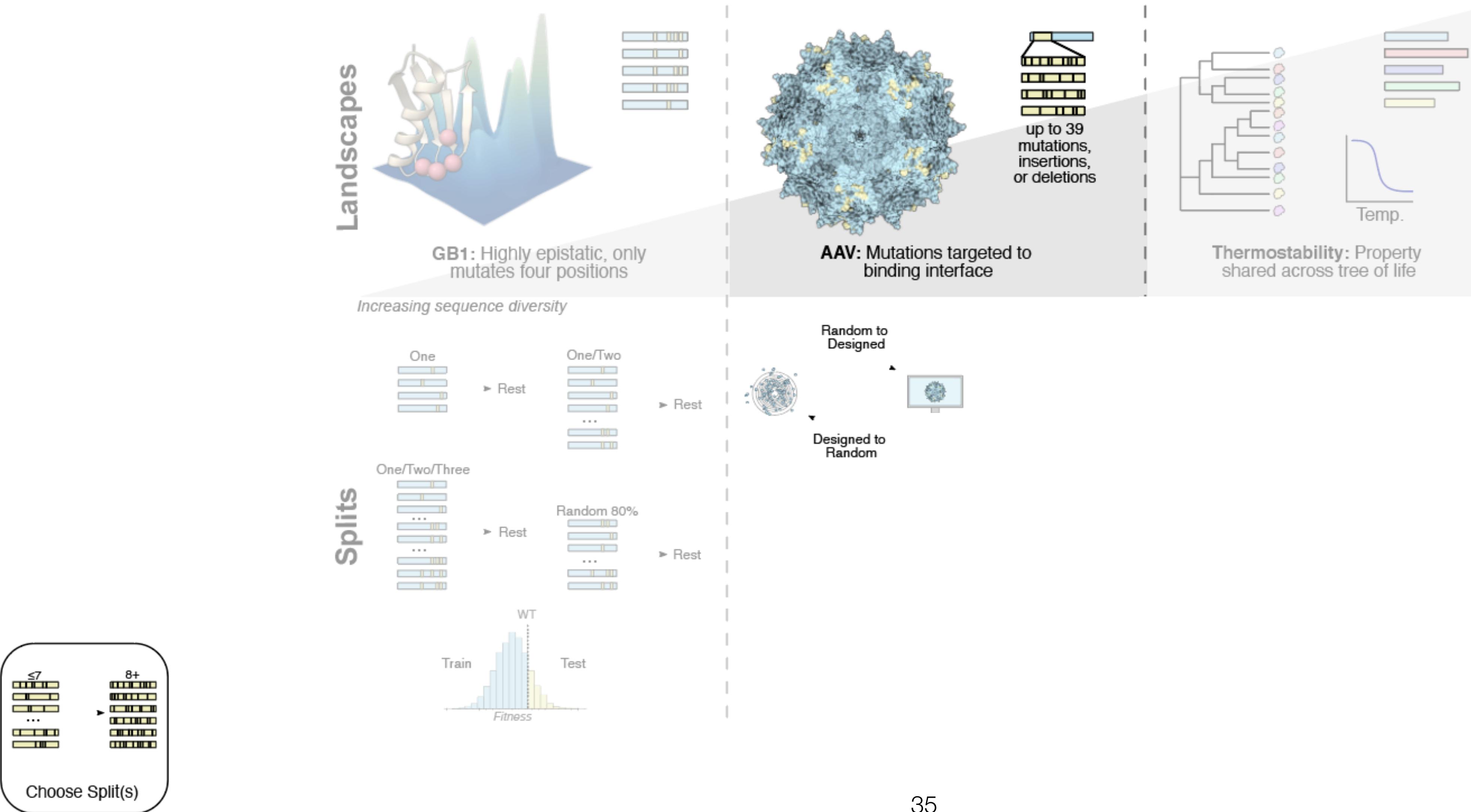
# Make biologically-relevant train/test splits



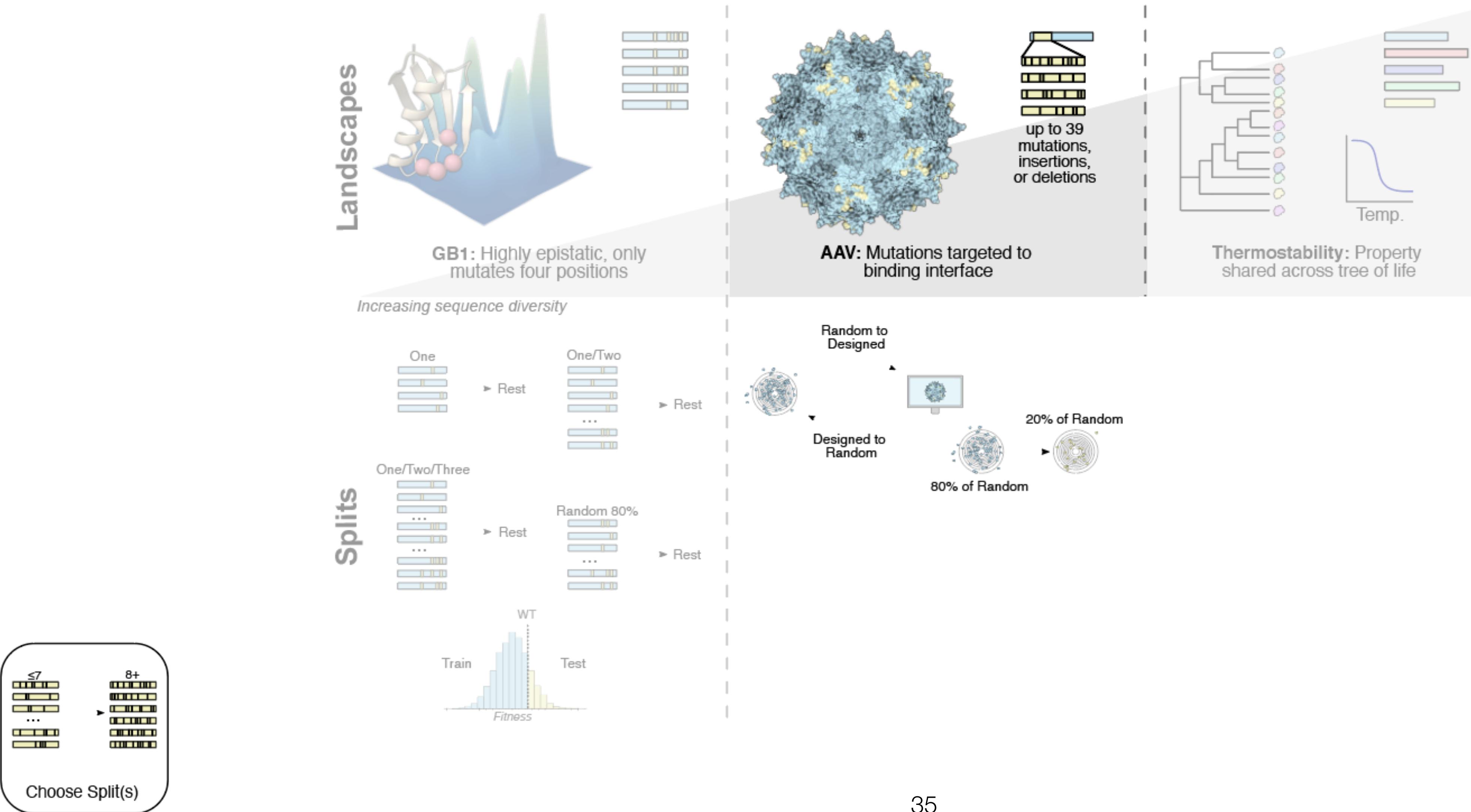
# Make biologically-relevant train/test splits



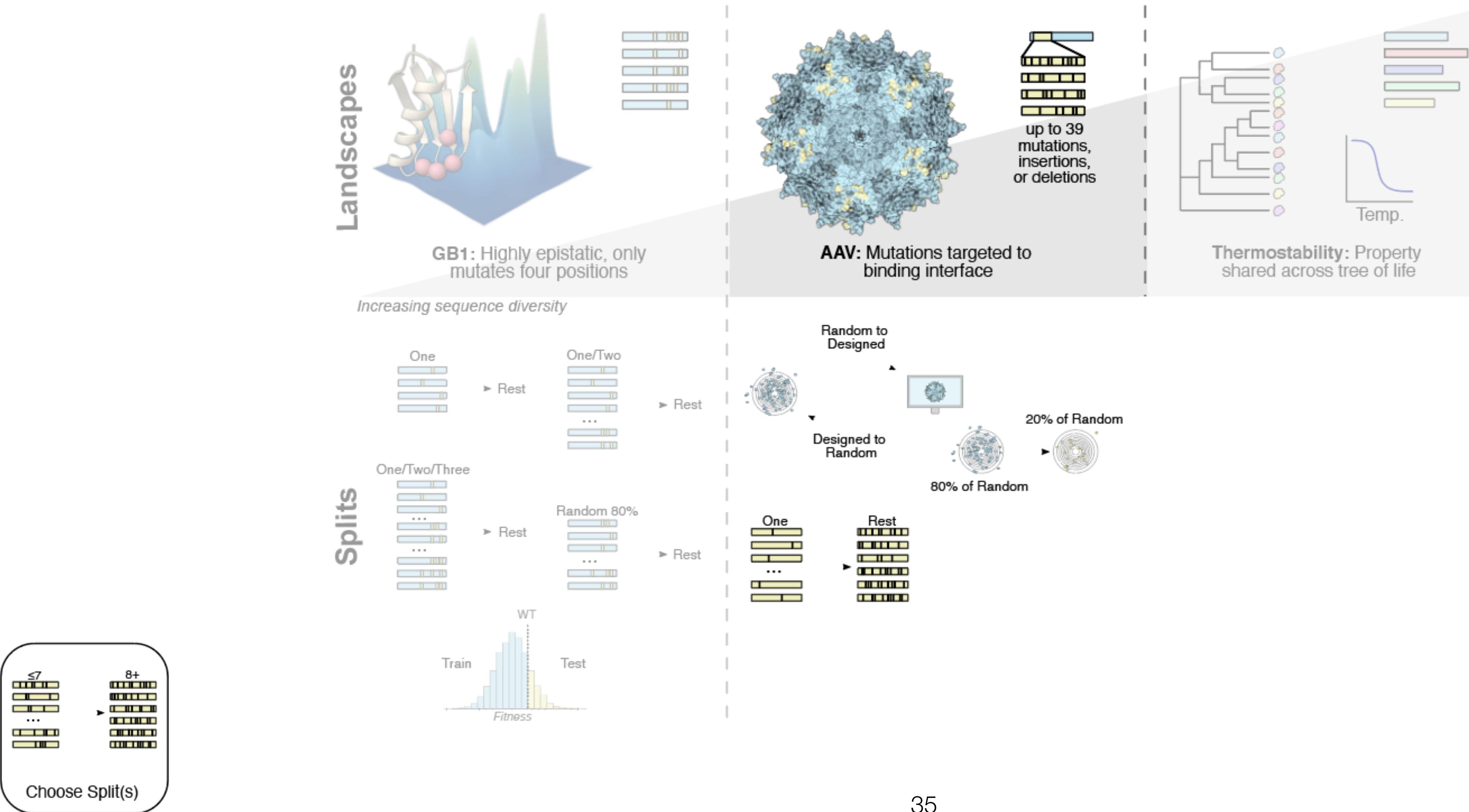
# Make biologically-relevant train/test splits



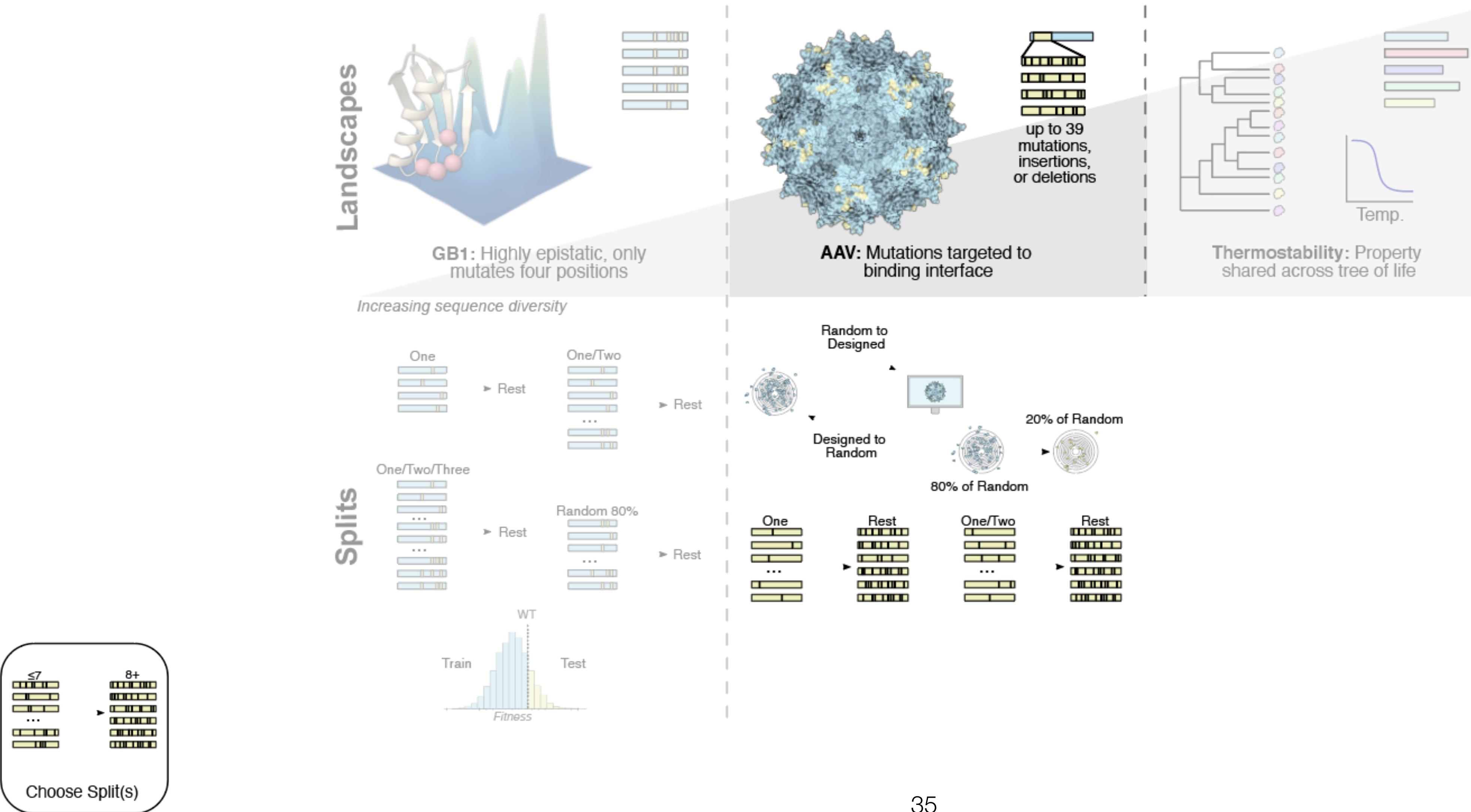
# Make biologically-relevant train/test splits



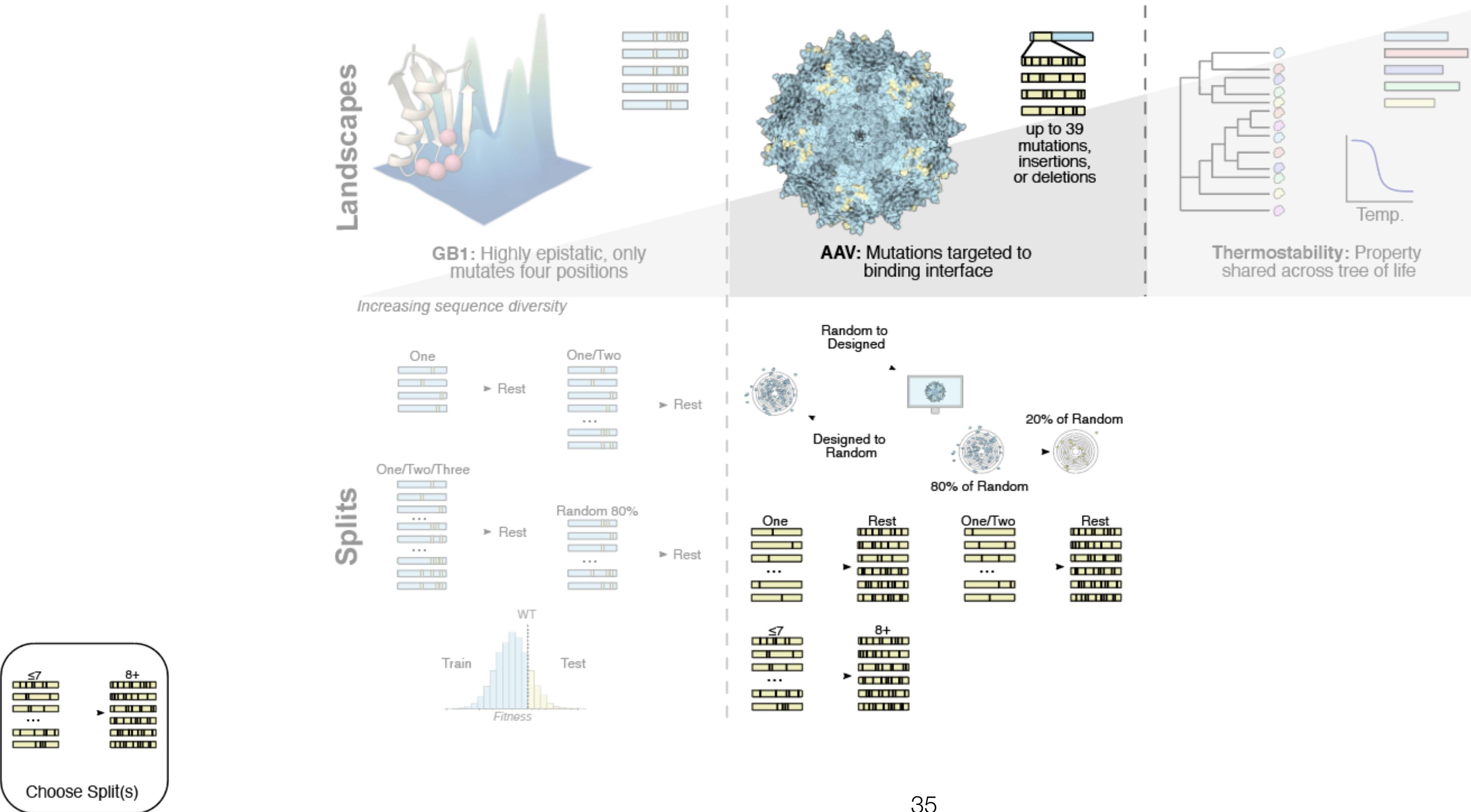
# Make biologically-relevant train/test splits



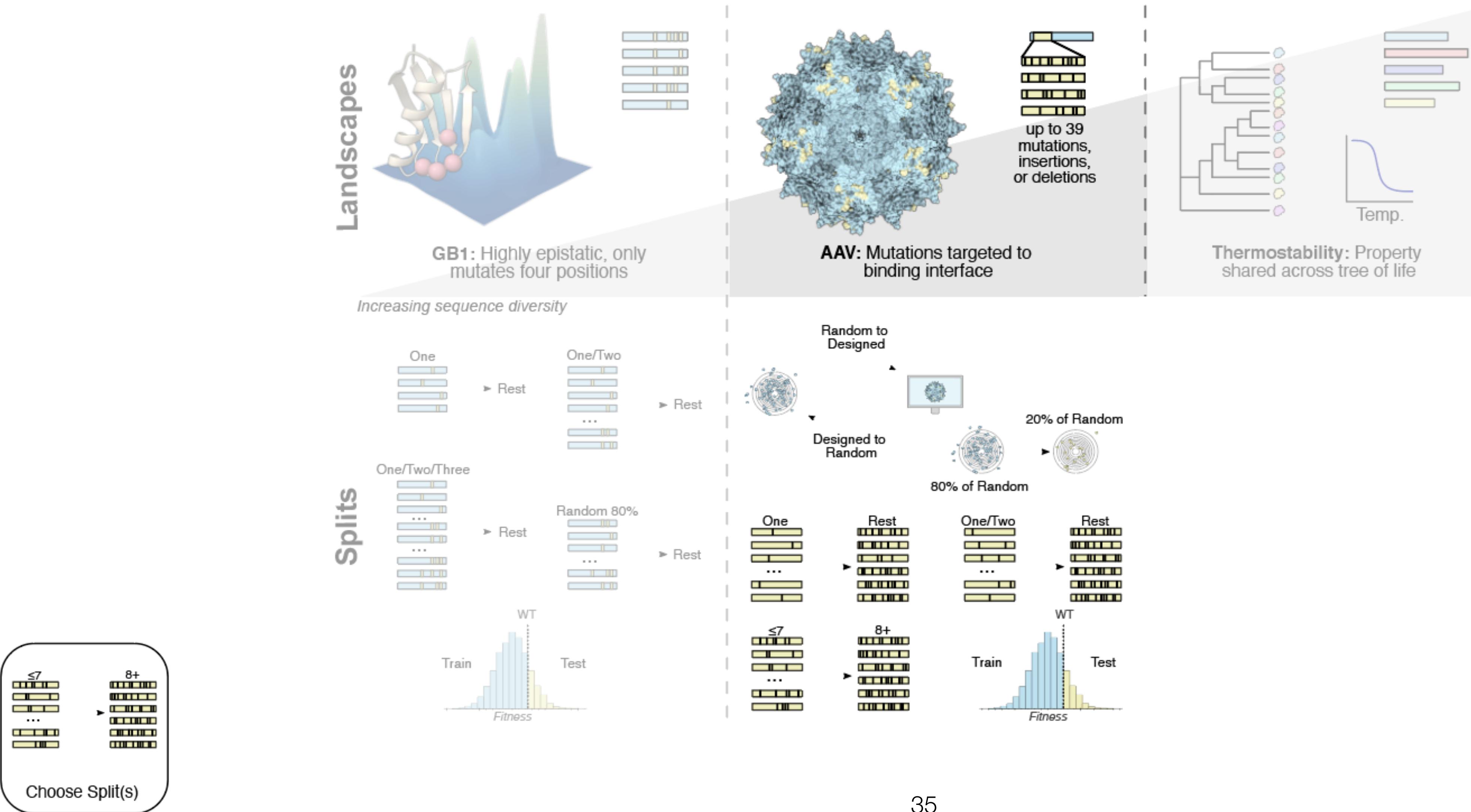
# Make biologically-relevant train/test splits



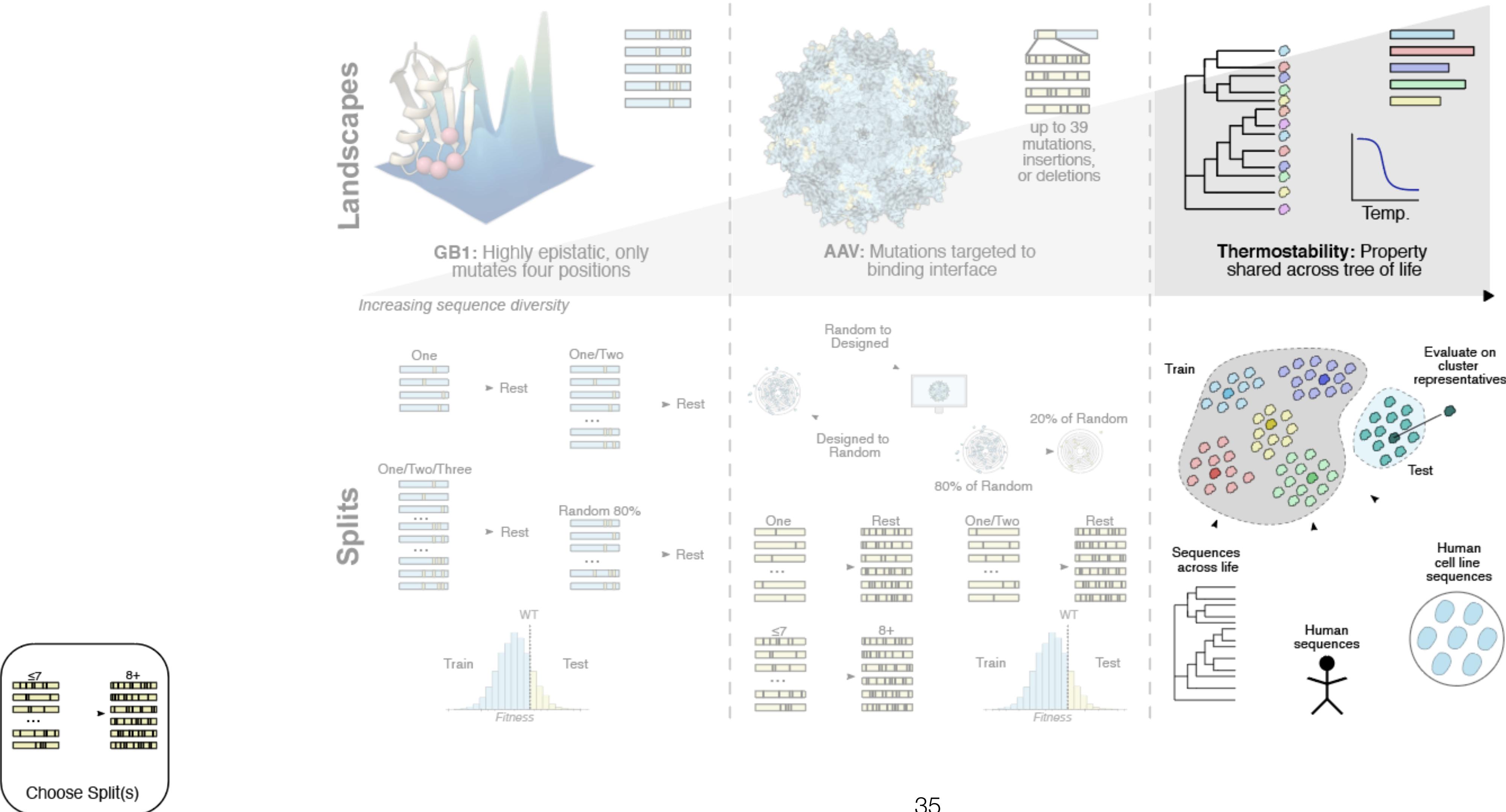
# Make biologically-relevant train/test splits



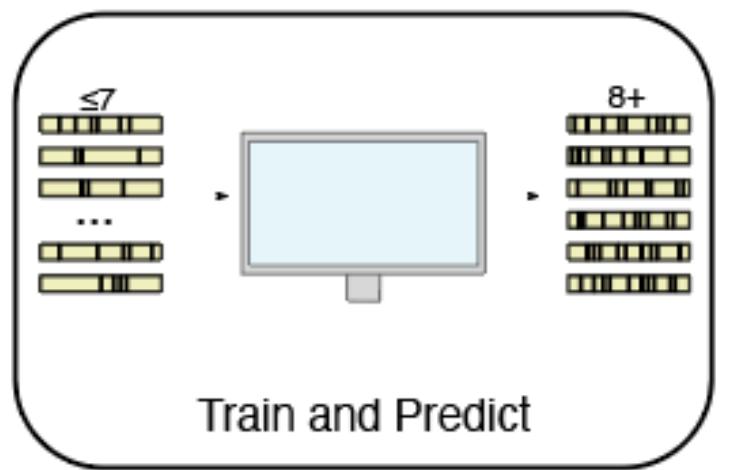
# Make biologically-relevant train/test splits



# Make biologically-relevant train/test splits



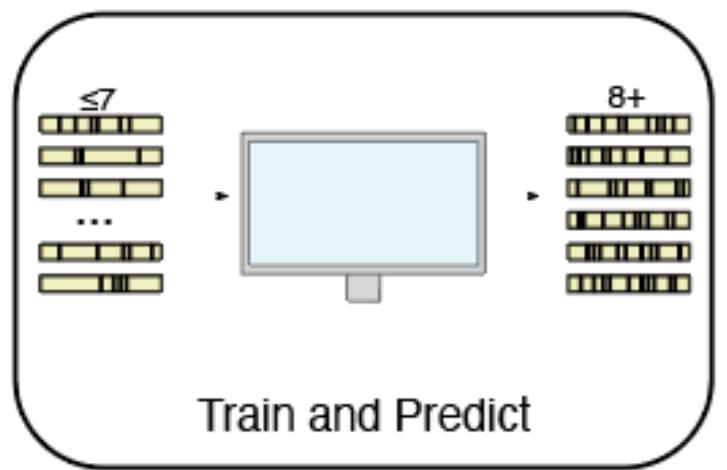
# Evaluate many baseline models



Dallago *et al.* 2021

# Evaluate many baseline models

Sequence distance

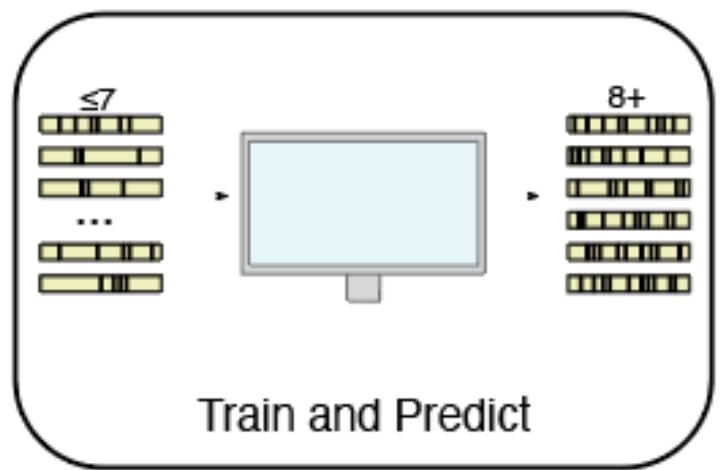


Dallago et al. 2021

# Evaluate many baseline models

Sequence distance

Levenshtein distance



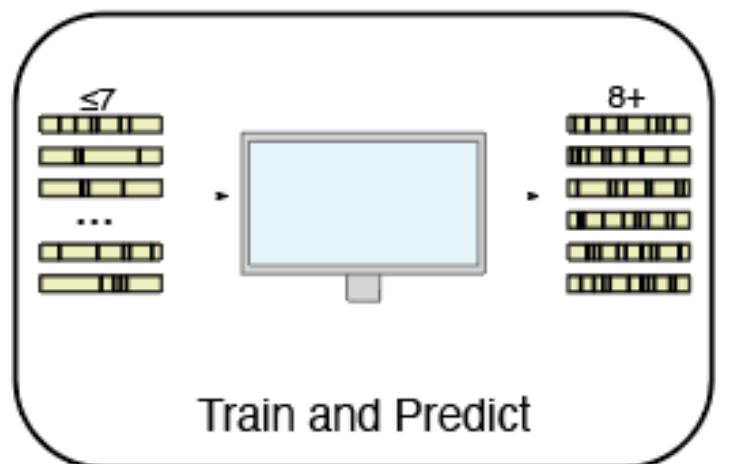
Dallago et al. 2021

# Evaluate many baseline models

Sequence distance

Levenshtein distance

BLOSUM62



Dallago et al. 2021

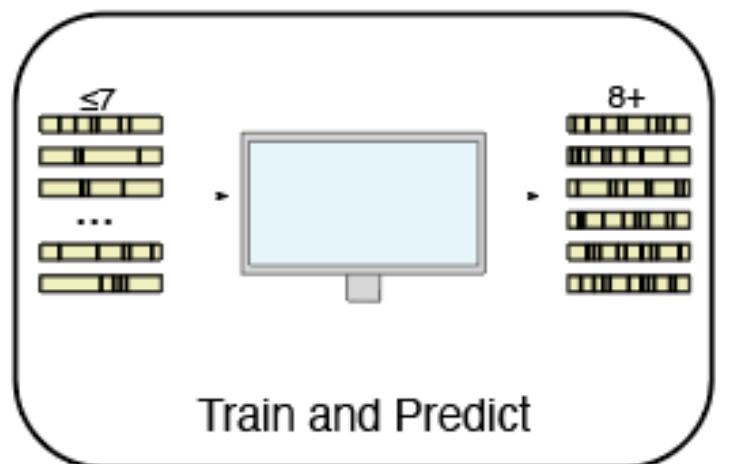
# Evaluate many baseline models

Sequence distance

Levenshtein distance

BLOSUM62

Supervised models



Dallago et al. 2021

# Evaluate many baseline models

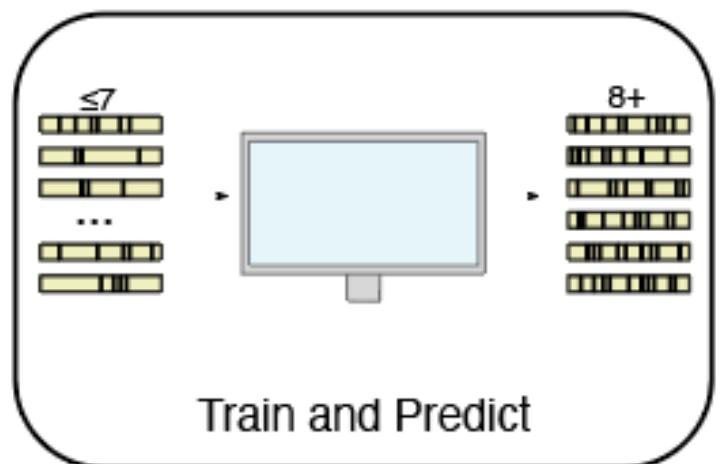
Sequence distance

Levenshtein distance

BLOSUM62

Supervised models

Ridge regression



Dallago et al. 2021

# Evaluate many baseline models

Sequence distance

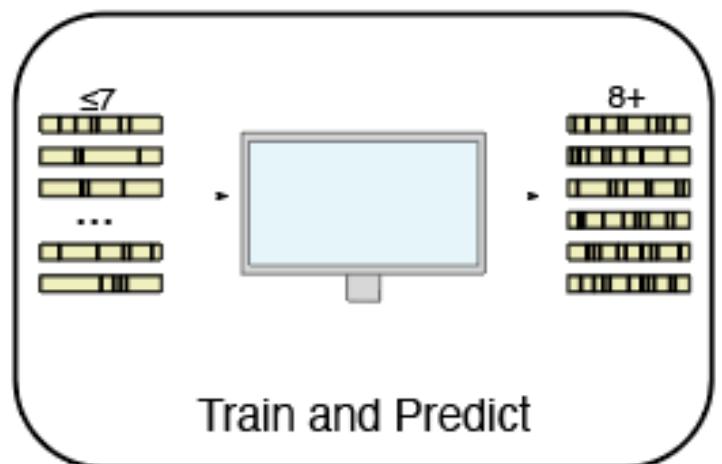
Levenshtein distance

BLOSUM62

Supervised models

Ridge regression

small CNN



Train and Predict

Dallago et al. 2021

# Evaluate many baseline models

Sequence distance

Levenshtein distance

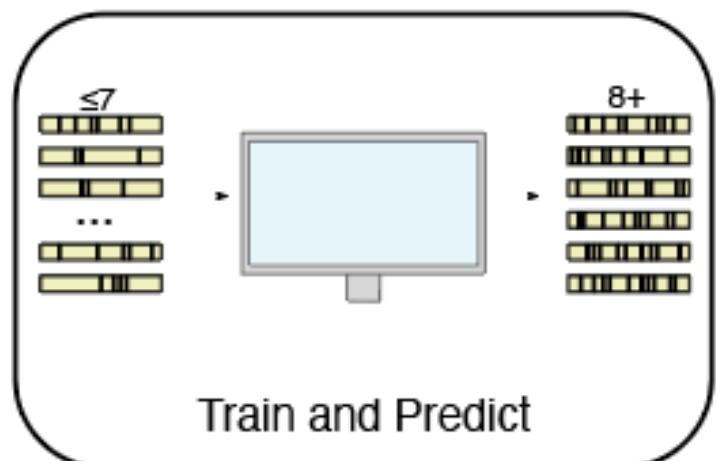
BLOSUM62

Supervised models

Ridge regression

small CNN

ESM-untrained  
(750M par. transformer)



Dallago et al. 2021

# Evaluate many baseline models

Sequence distance

Levenshtein distance

BLOSUM62

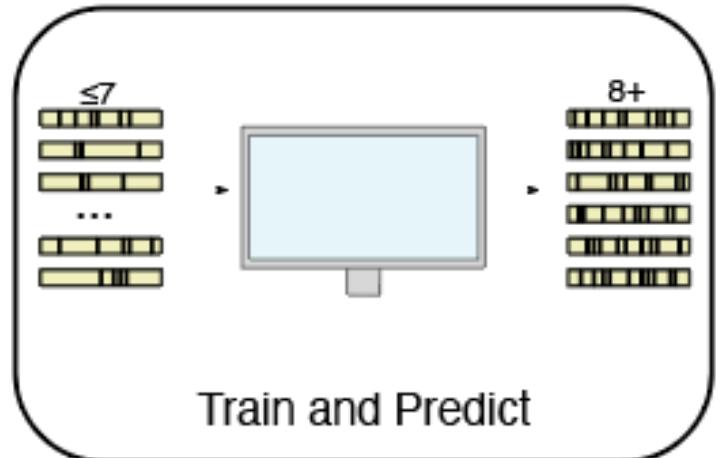
Supervised models

Ridge regression

small CNN

ESM-untrained  
(750M par. transformer)

Pretrained models



Dallago et al. 2021

# Evaluate many baseline models

Sequence distance

Levenshtein distance

BLOSUM62

Supervised models

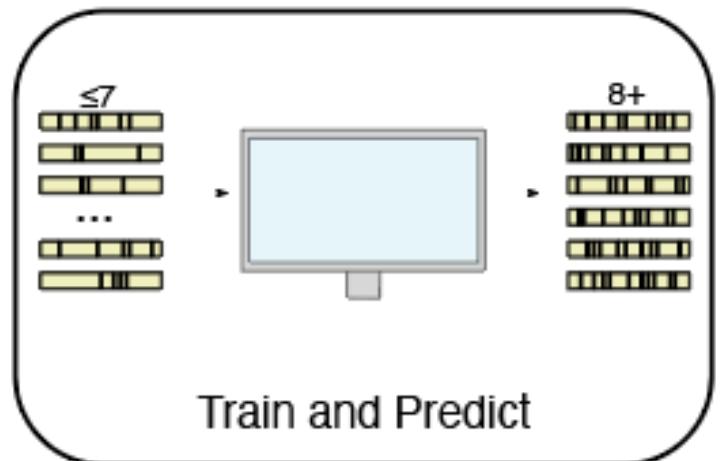
Ridge regression

small CNN

ESM-untrained  
(750M par. transformer)

Pretrained models

ESM-1b  
(pretrained on UniRef50)



Dallago et al. 2021

# Evaluate many baseline models

## Sequence distance

Levenshtein distance

BLOSUM62

## Supervised models

Ridge regression

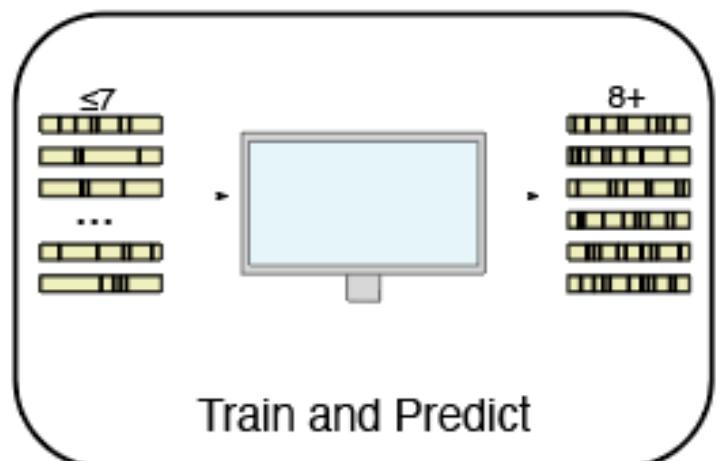
small CNN

ESM-untrained  
(750M par. transformer)

## Pretrained models

ESM-1b  
(pretrained on UniRef50)

ESM-1v  
(pretrained on UniRef90)



Dallago et al. 2021

# Evaluate many baseline models

## Sequence distance

Levenshtein distance

BLOSUM62

## Supervised models

Ridge regression

small CNN

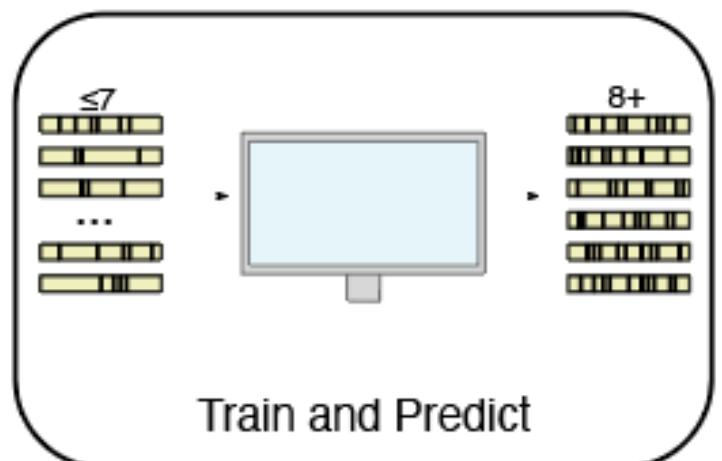
ESM-untrained  
(750M par. transformer)

## Pretrained models

ESM-1b  
(pretrained on UniRef50)

ESM-1v  
(pretrained on UniRef90)

Best on highly-diverse landscapes (thermostability)  
or when data is very limited



Dallago et al. 2021

# Evaluate many baseline models

## Sequence distance

Levenshtein distance

BLOSUM62

## Supervised models

Ridge regression

small CNN

ESM-untrained  
(750M par. transformer)

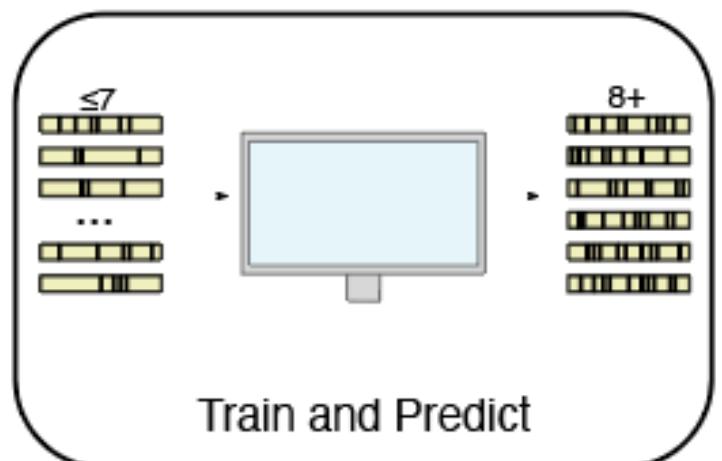
Strong performance on most landscapes

## Pretrained models

ESM-1b  
(pretrained on UniRef50)

ESM-1v  
(pretrained on UniRef90)

Best on highly-diverse landscapes (thermostability)  
or when data is very limited

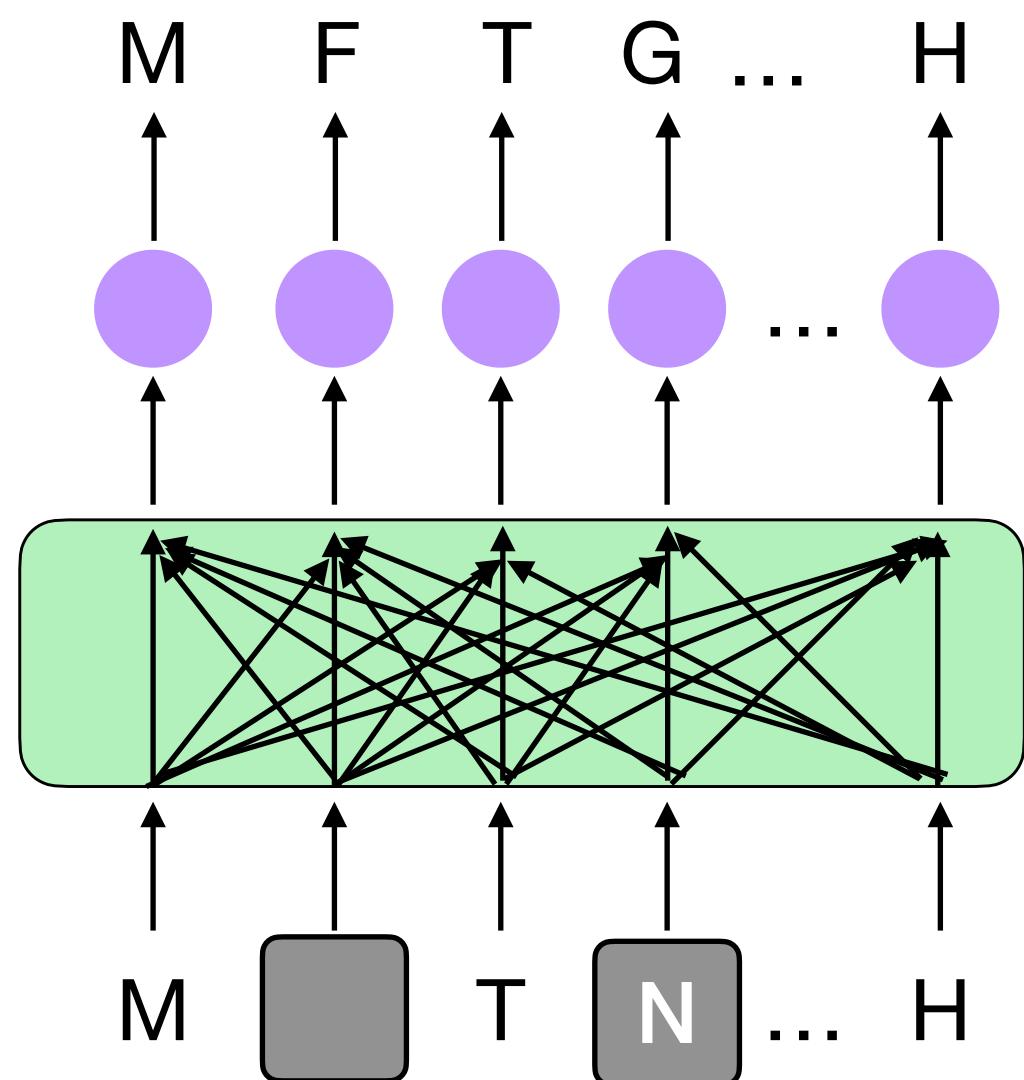


Dallago et al. 2021

# Try your models!

**<https://benchmark.protein.properties/home>**

# Local-to-global aggregation affects model performance



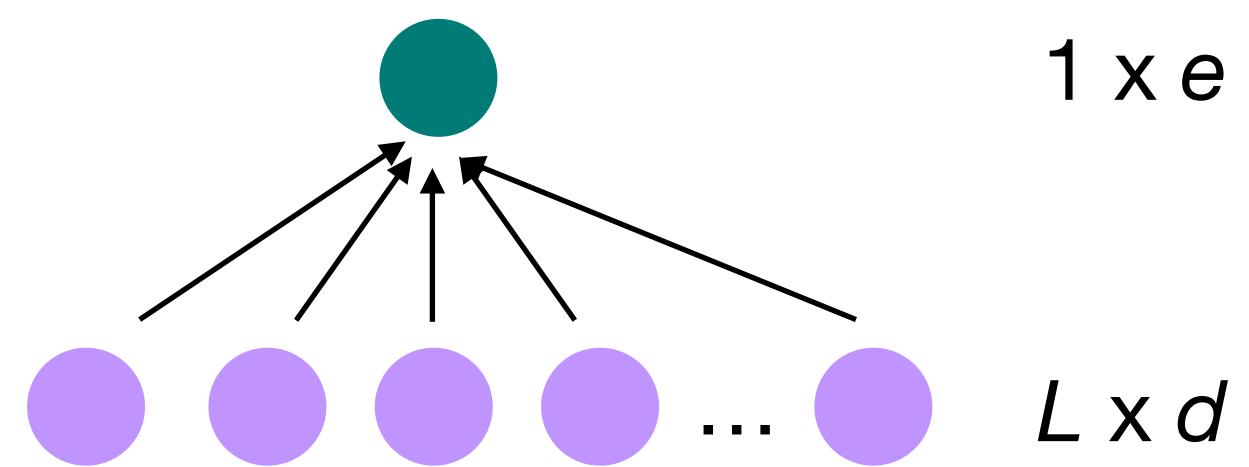
# Local-to-global aggregation affects model performance



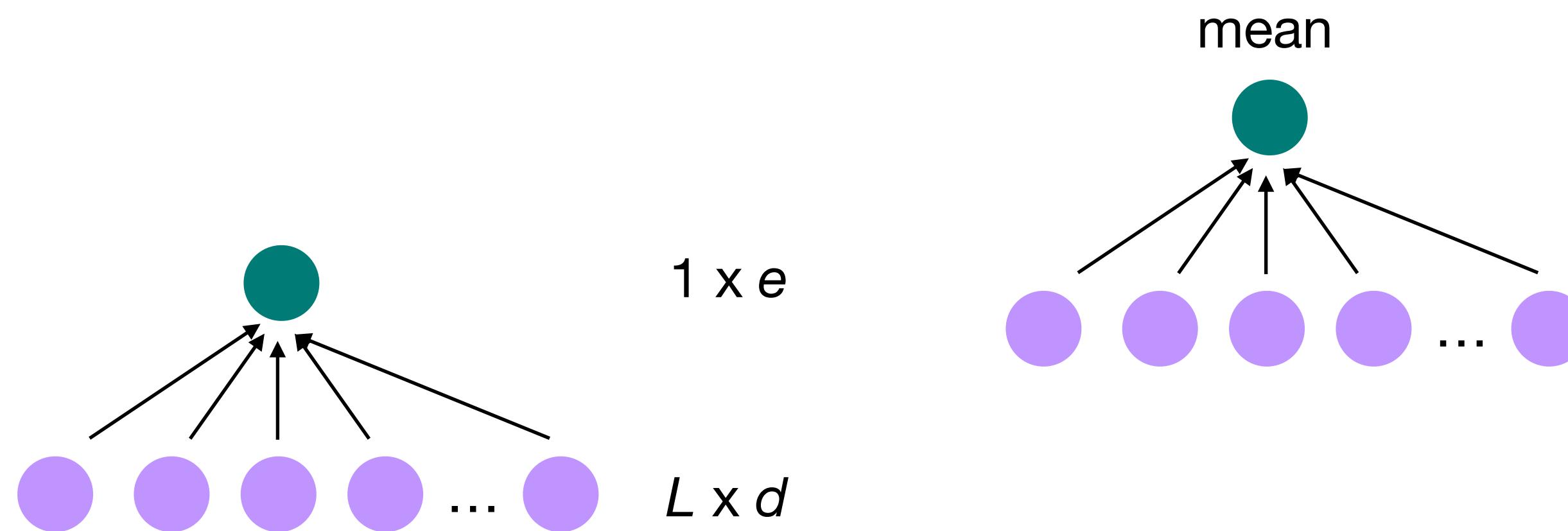
# Local-to-global aggregation affects model performance

  $L \times d$

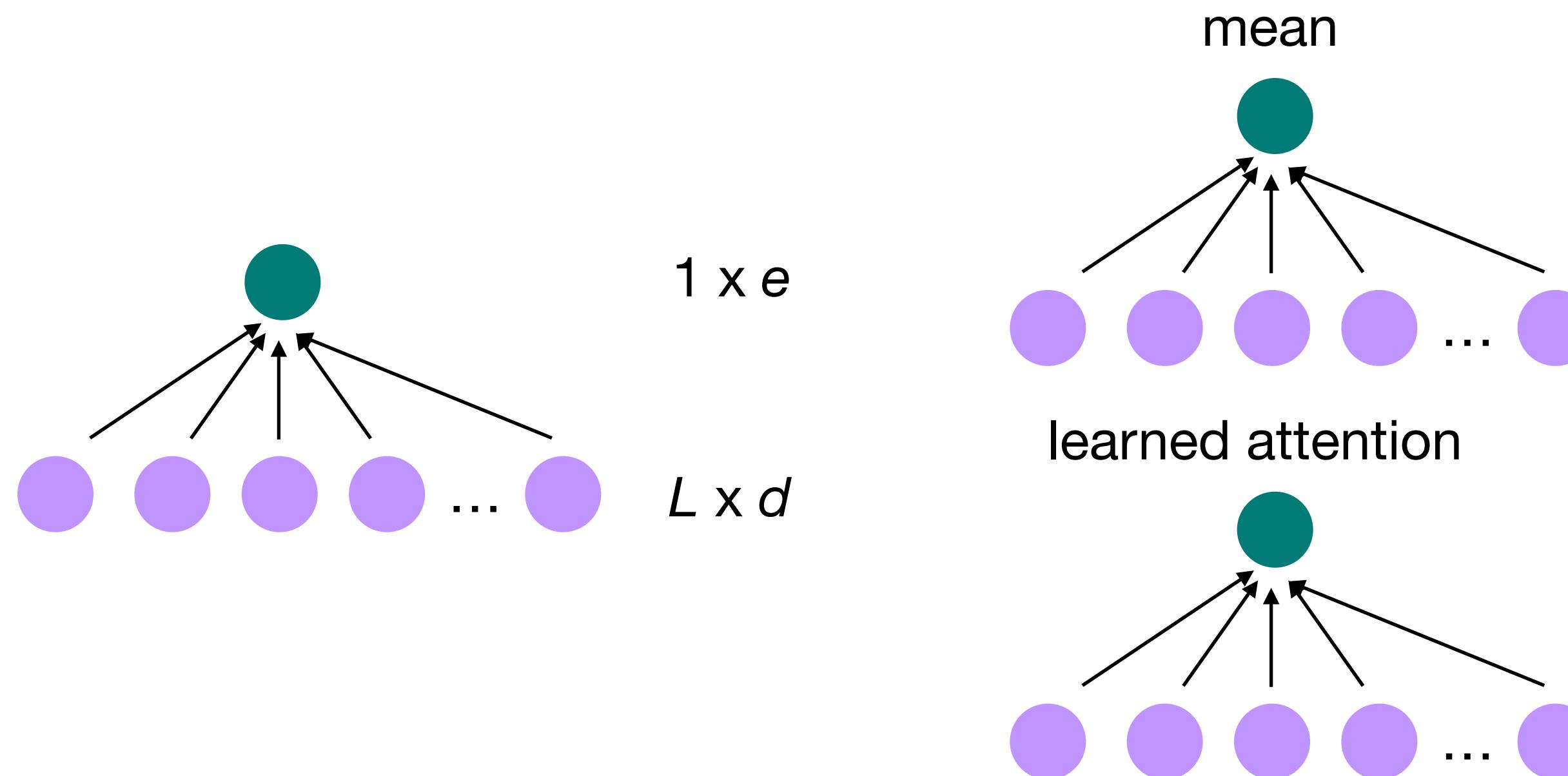
# Local-to-global aggregation affects model performance



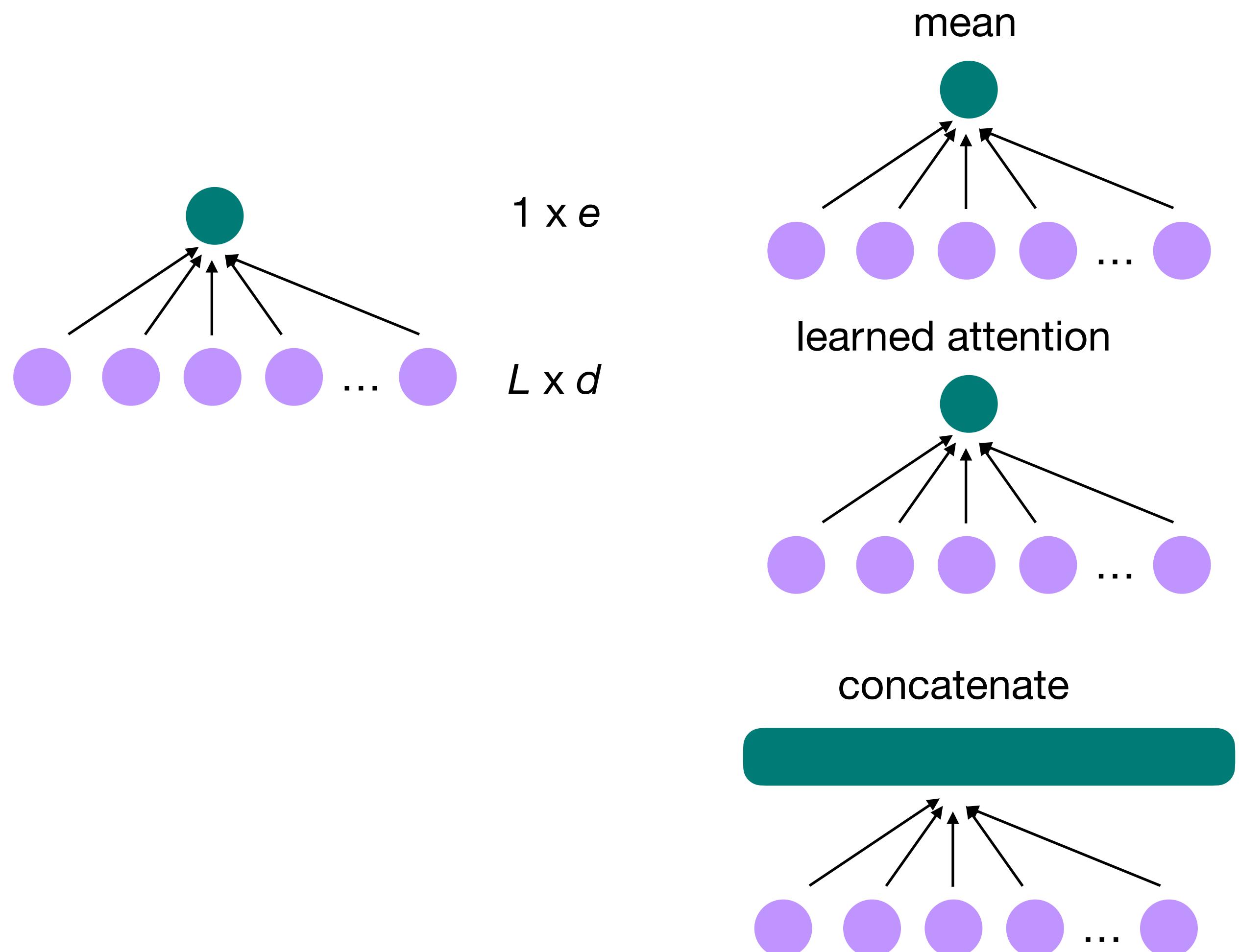
# Local-to-global aggregation affects model performance



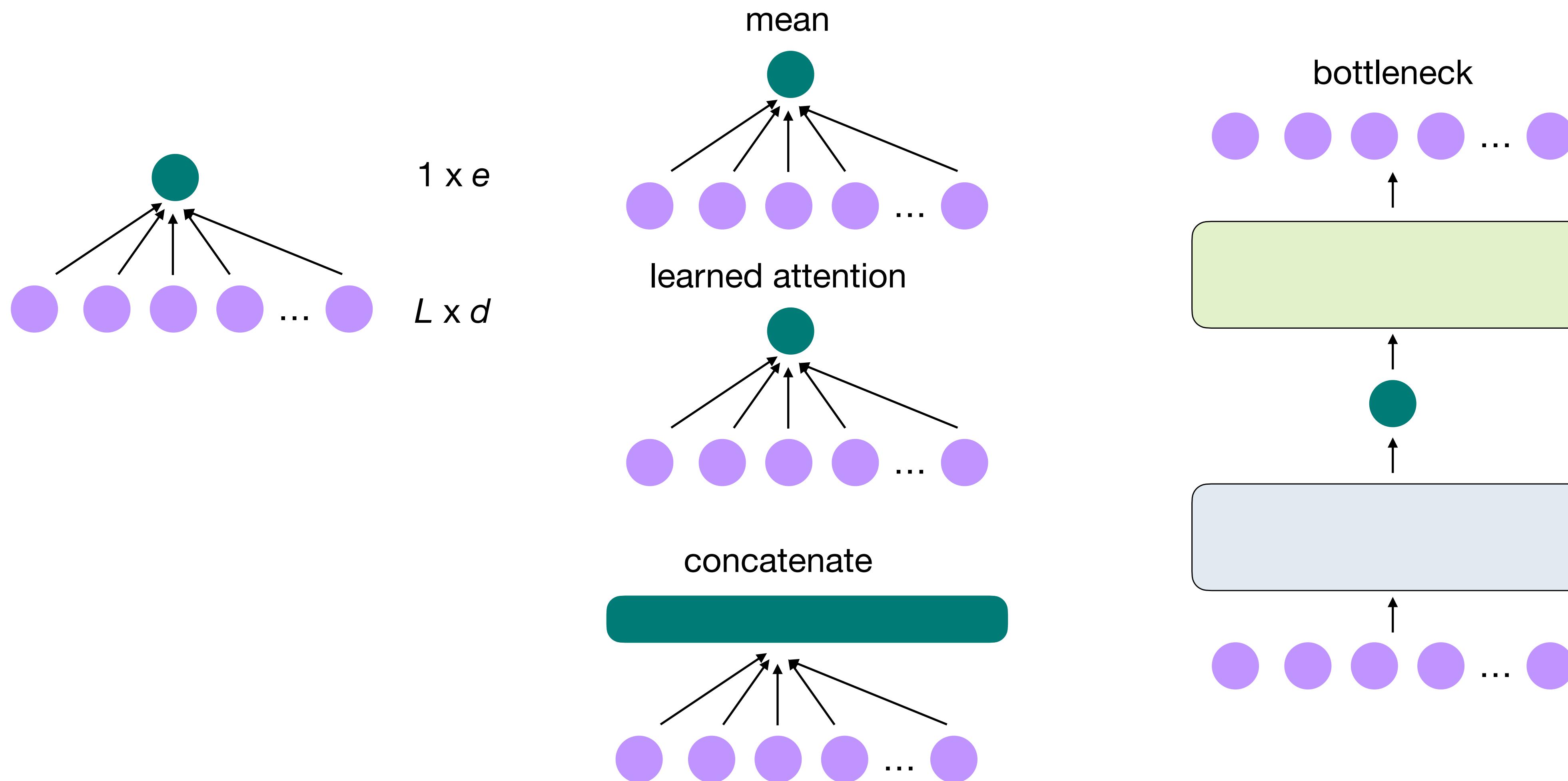
# Local-to-global aggregation affects model performance



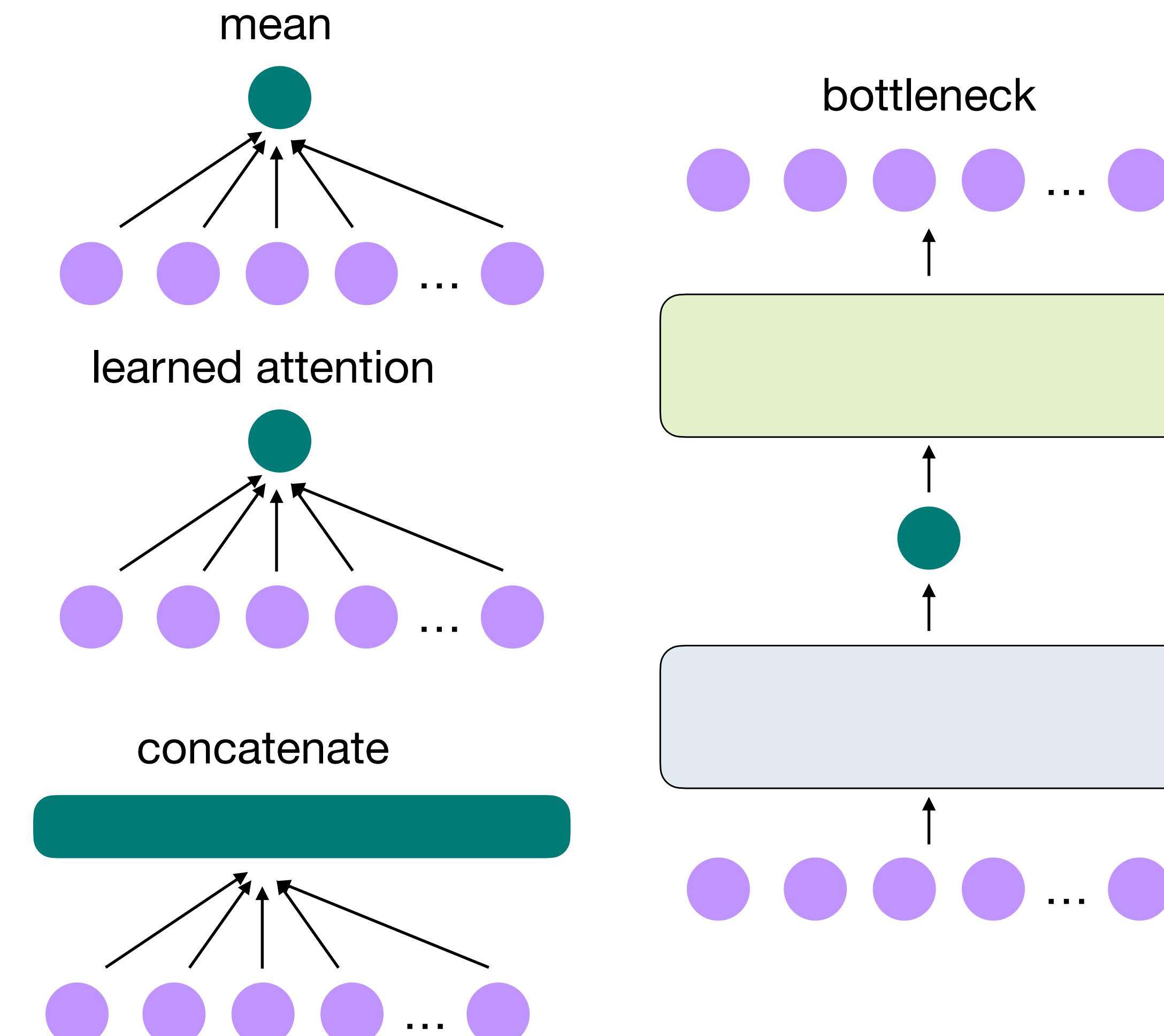
# Local-to-global aggregation affects model performance



# Local-to-global aggregation affects model performance



# Local-to-global aggregation affects model performance

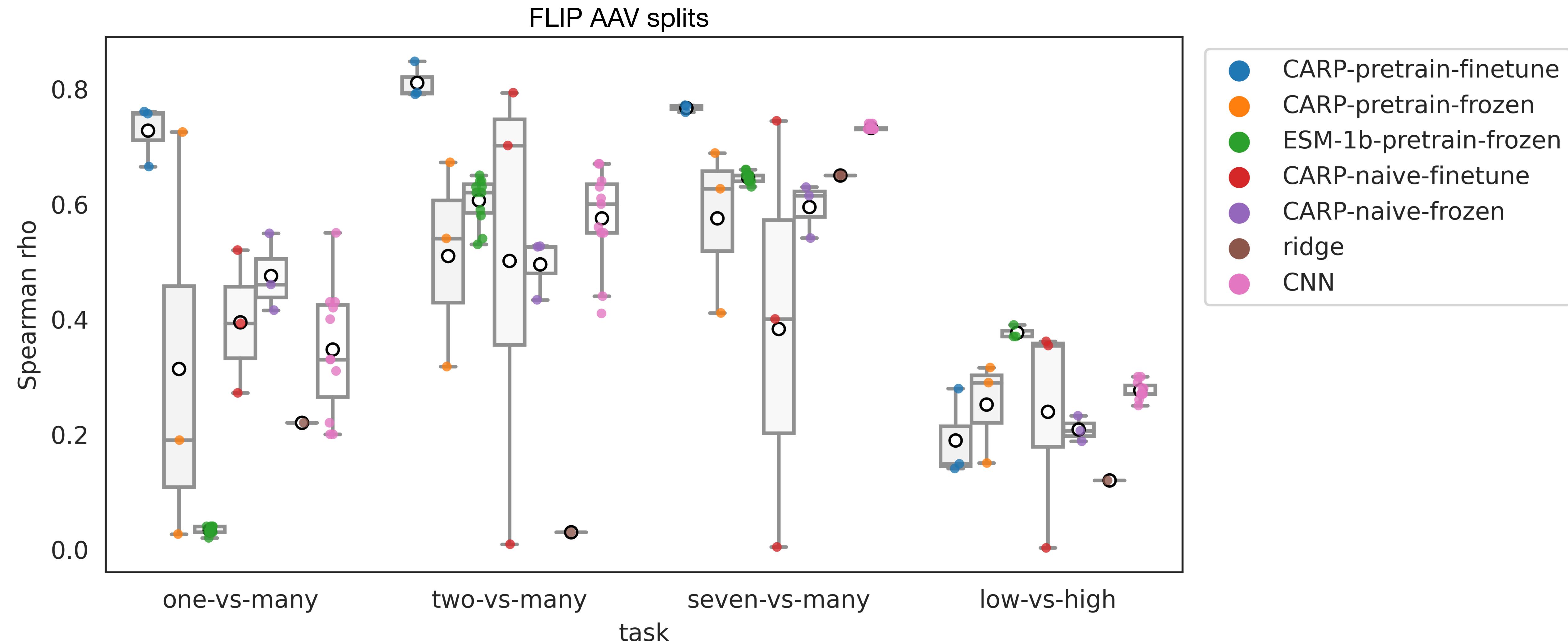


	Stability (Corr.)	Fluorescence (Corr.)	Homology (Acc.)
Mean	0.42	0.19	0.27
Attention	0.65	0.23	0.27
Light Att.	0.66	0.23	0.27
Maximum	0.02	0.02	0.28
MeanMax	0.37	0.15	0.26
KMax	0.10	0.11	0.27
Concat	0.74	0.69	0.34
Bottleneck	<b>0.79</b>	<b>0.78</b>	<b>0.41</b>

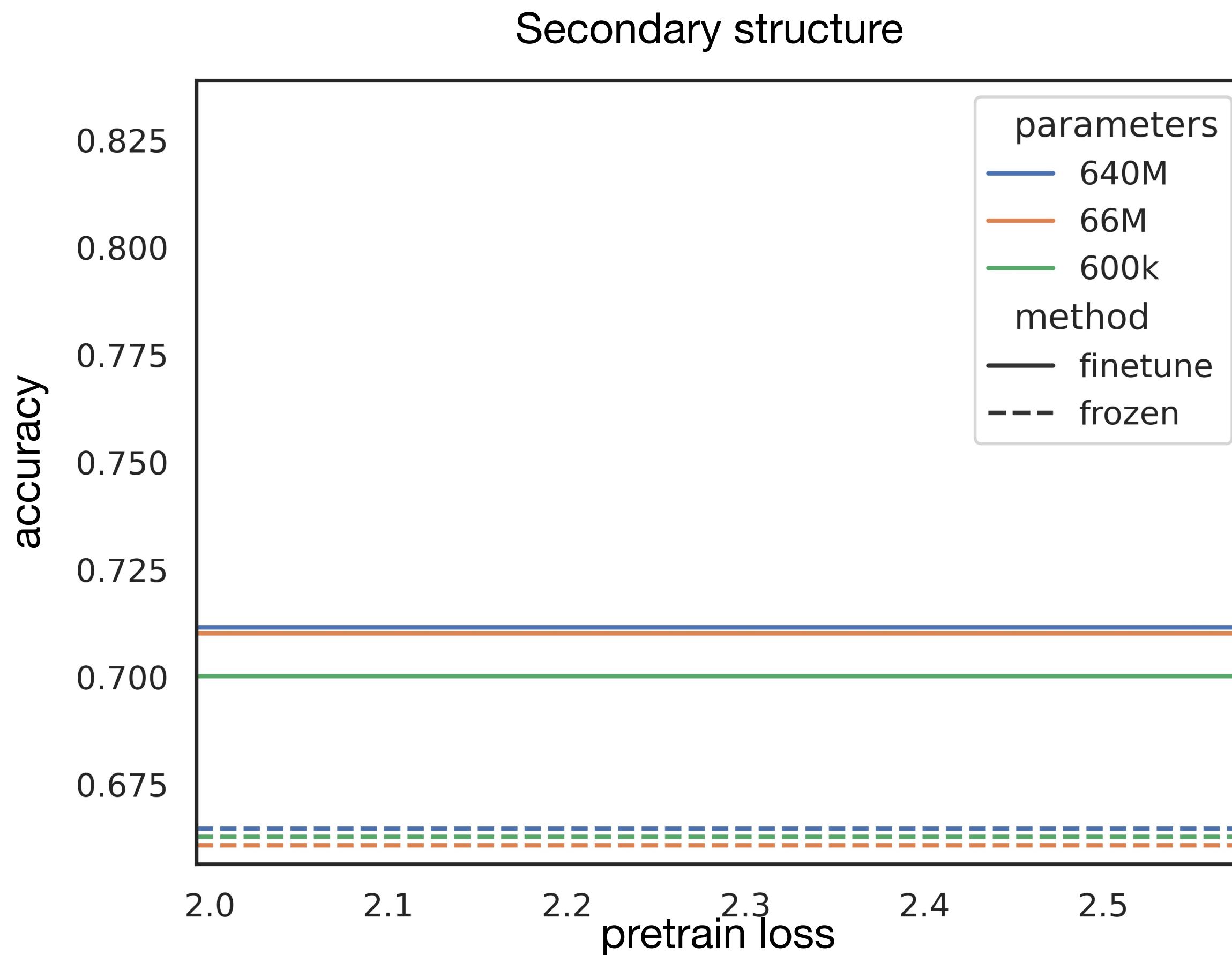
# Finetuning/fixing pretrained weights also matters

	Remote Homology			Fluorescence			Stability		
	Resnet	LSTM	Trans	Resnet	LSTM	Trans	Resnet	LSTM	Trans
PRE+Fix	0.27	<b>0.37</b>	0.27	0.23	<b>0.74</b>	0.48	0.65	0.70	0.62
PRE+FIN	0.17	0.26	0.21	0.21	0.67	0.68	<b>0.73</b>	0.69	<b>0.73</b>
RNG+Fix	0.03	0.10	0.04	0.25	0.63	0.14	0.21	0.61	-
RNG+FIN	0.10	0.12	0.09	-0.28	0.21	0.22	0.61	0.28	-0.06
Baseline	0.09 (Accuracy)			0.14 (Correlation)			0.19 (Correlation)		

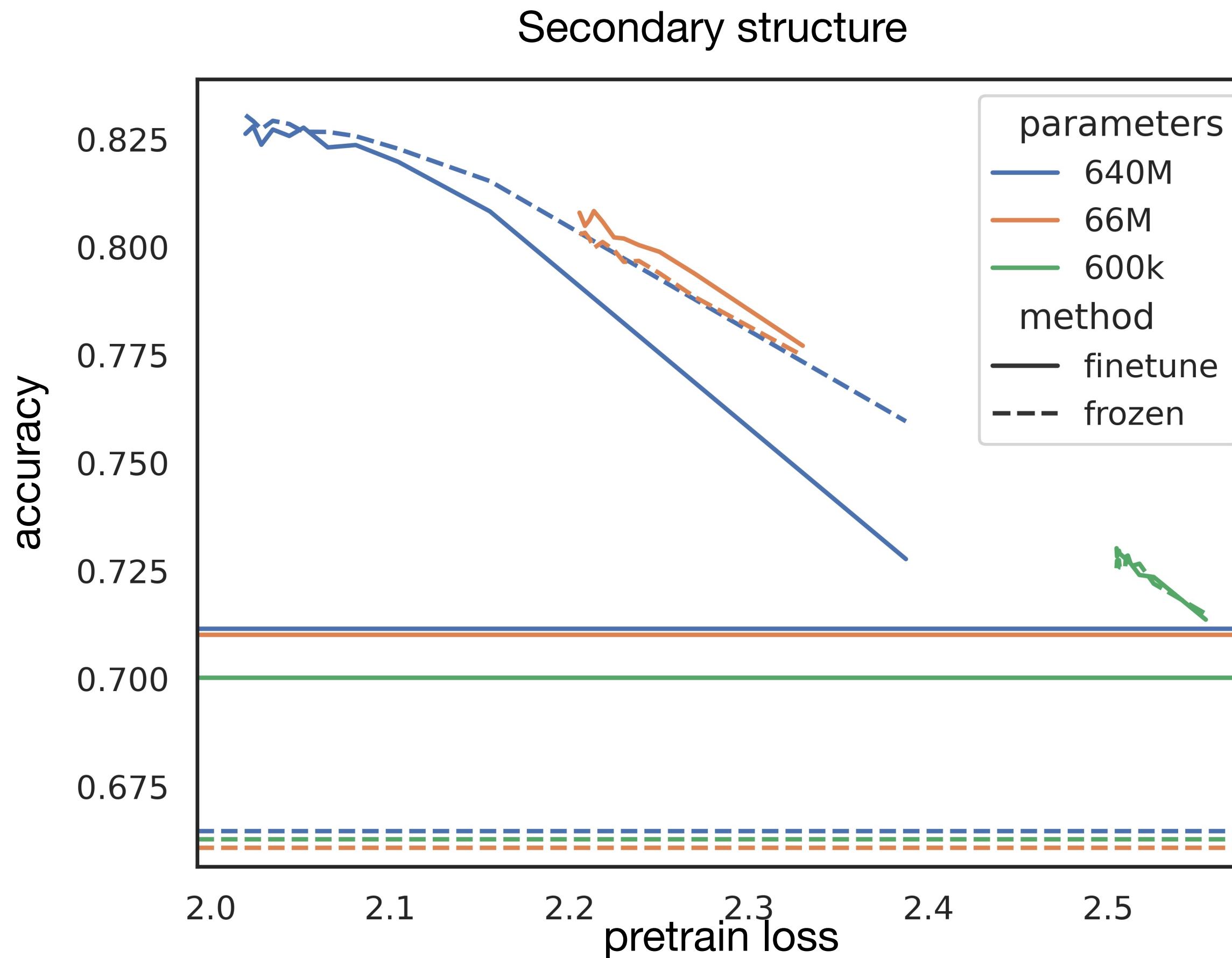
# Finetuning/fixing pretrained weights also matters



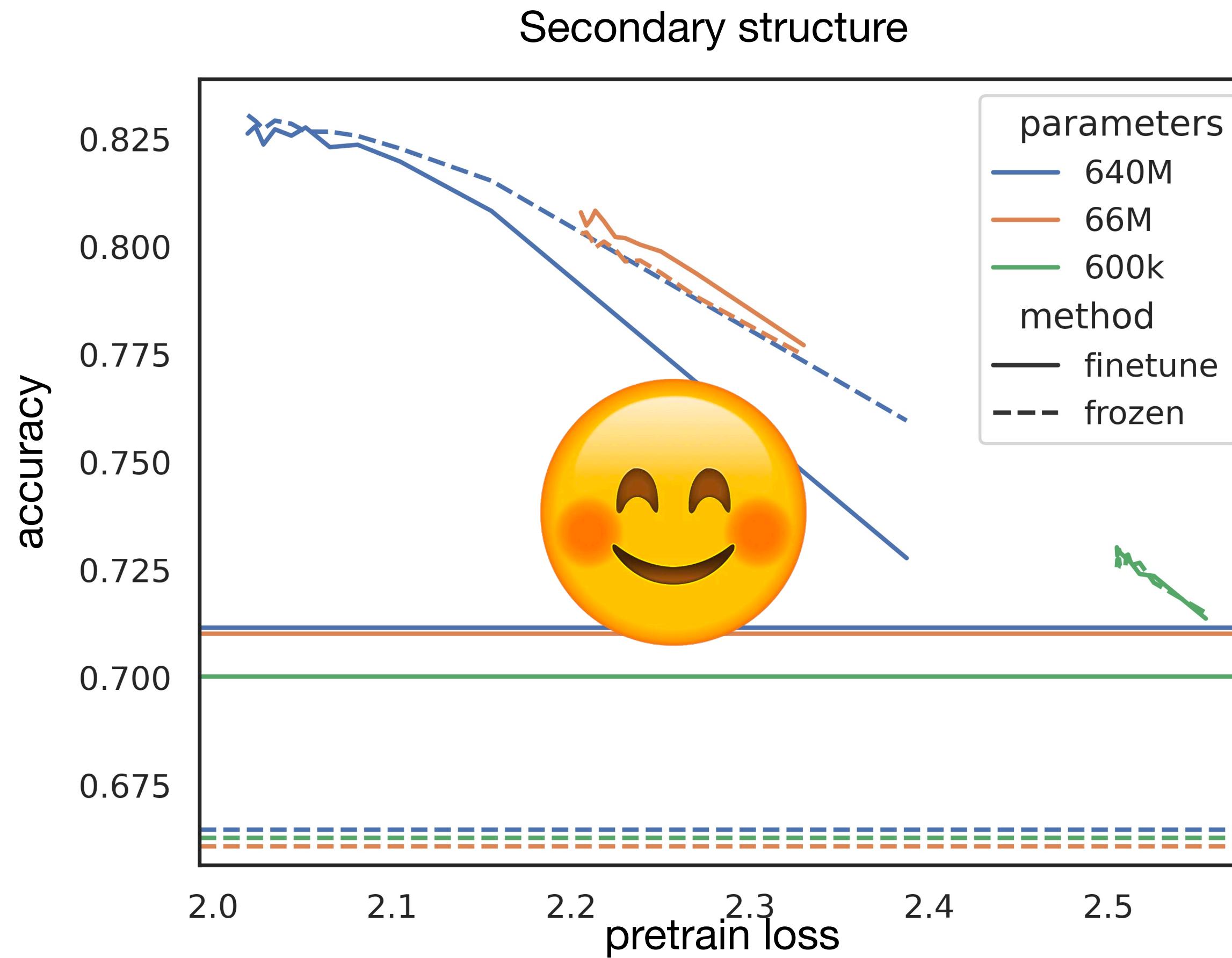
# When does pretraining help?



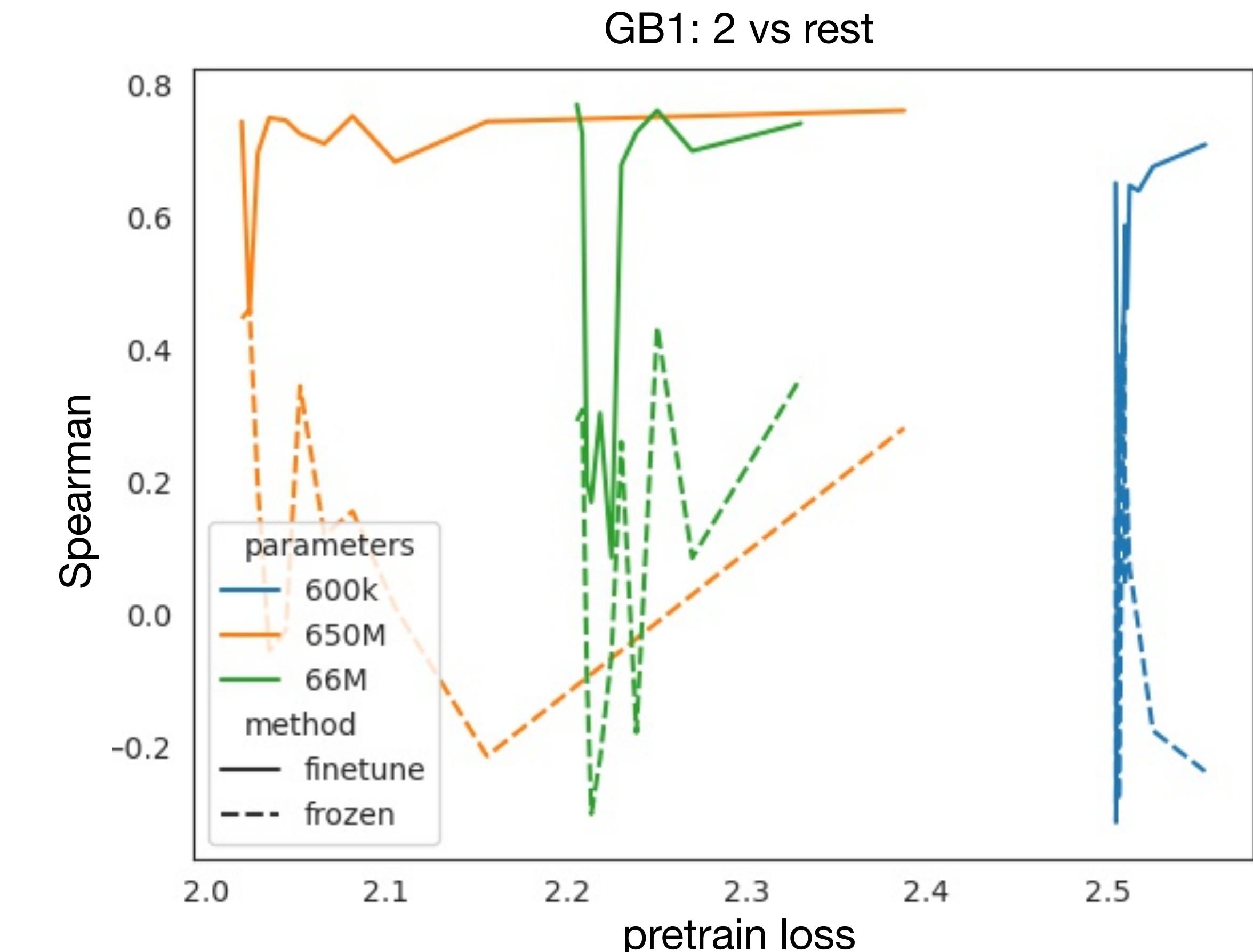
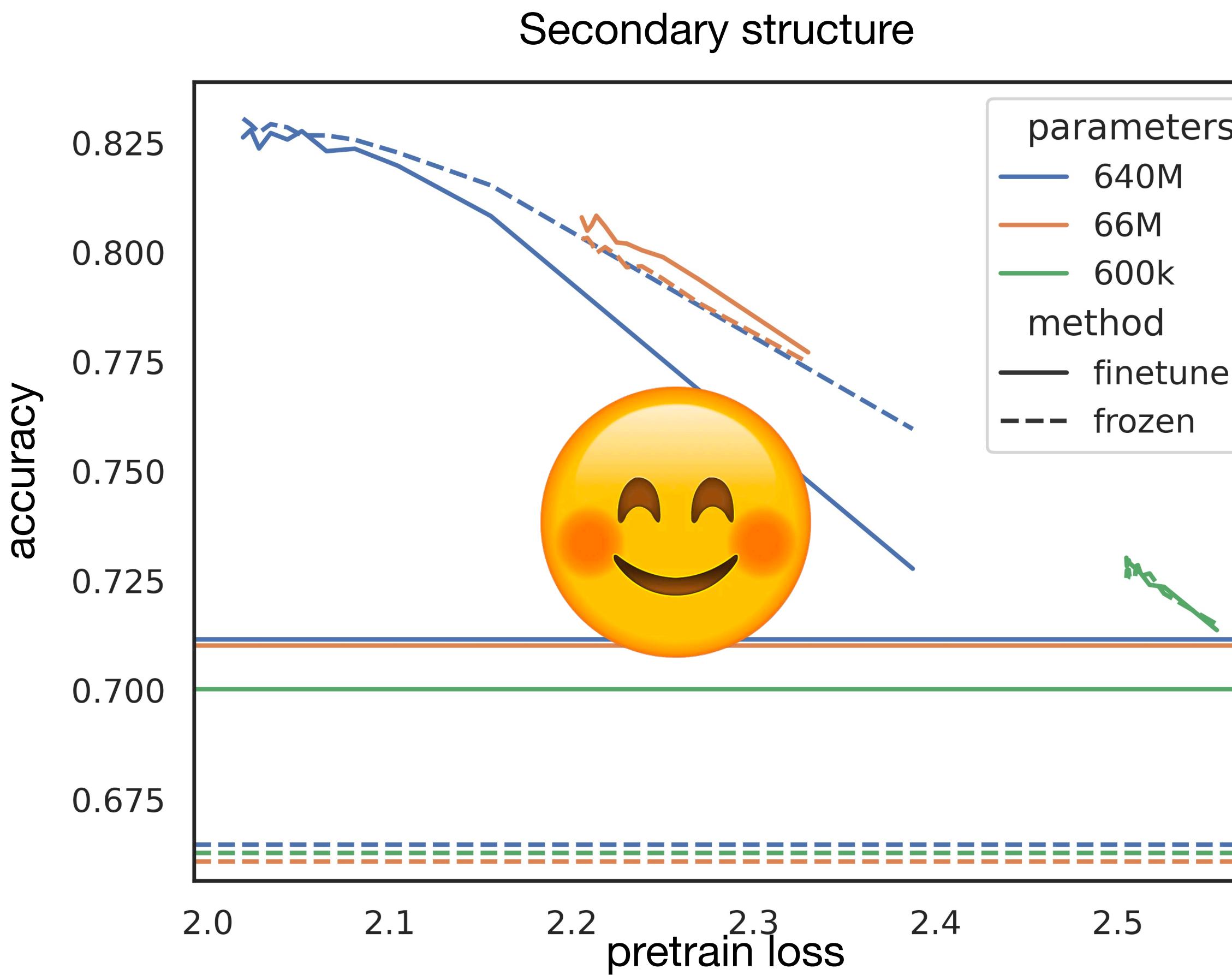
# When does pretraining help?



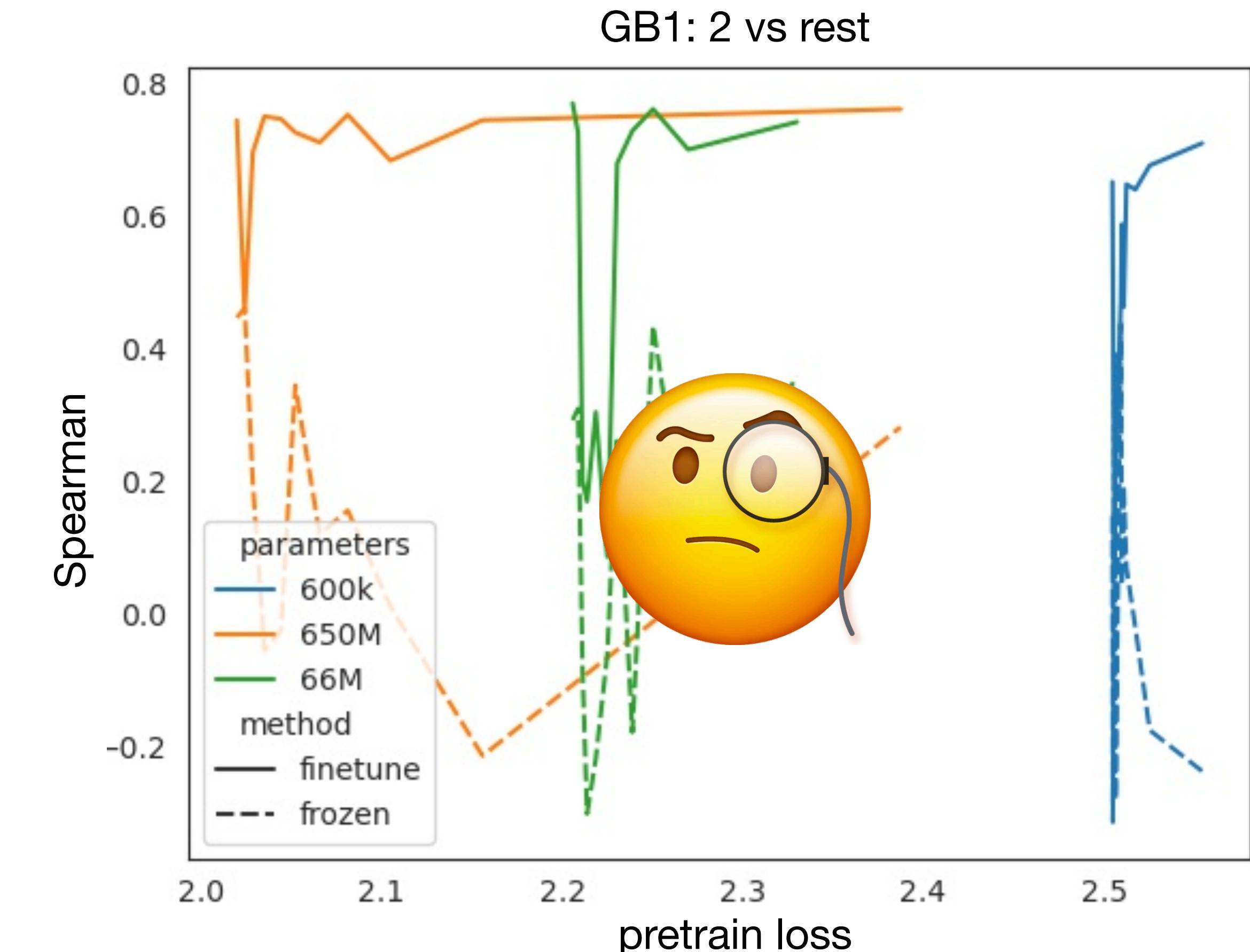
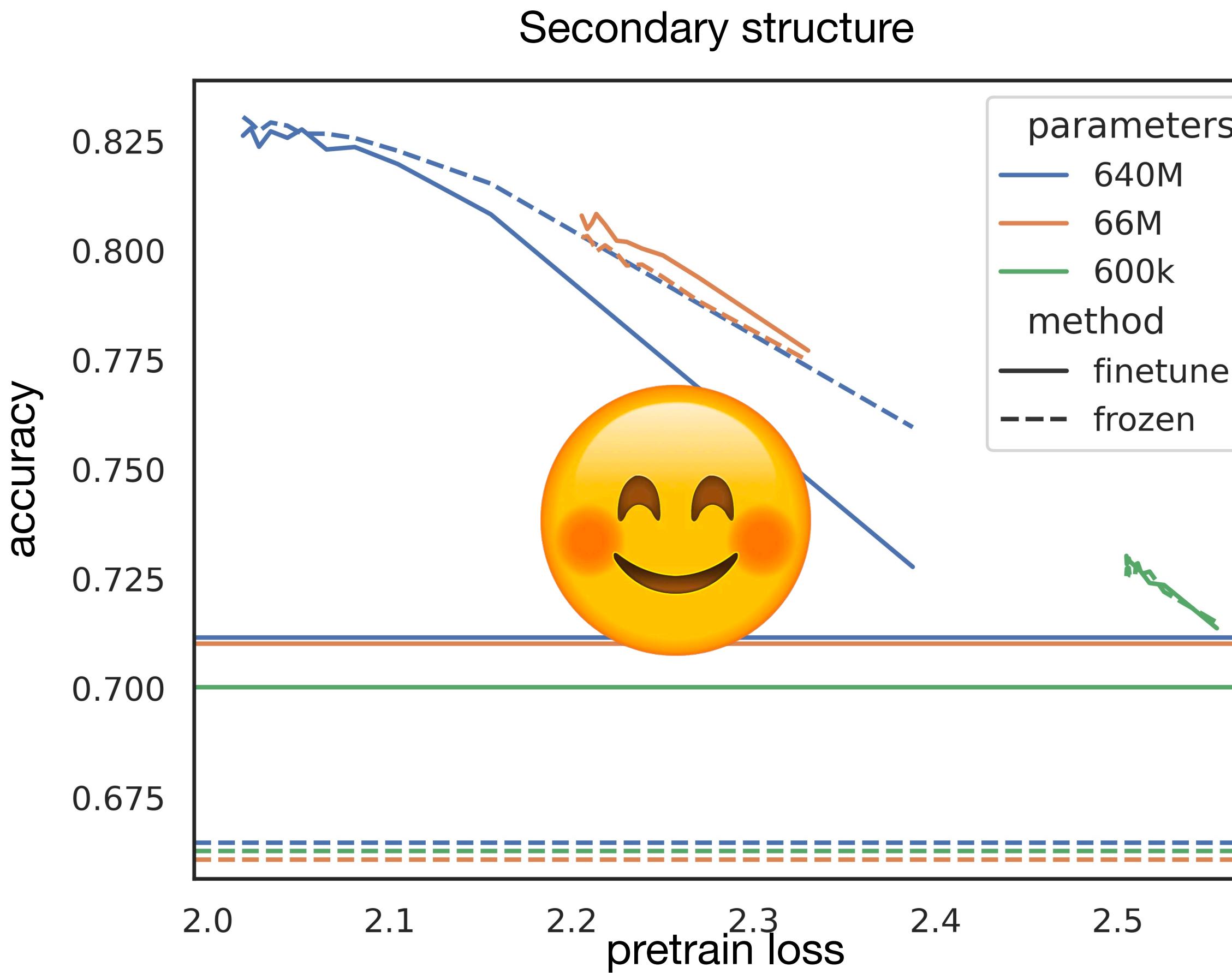
# When does pretraining help?



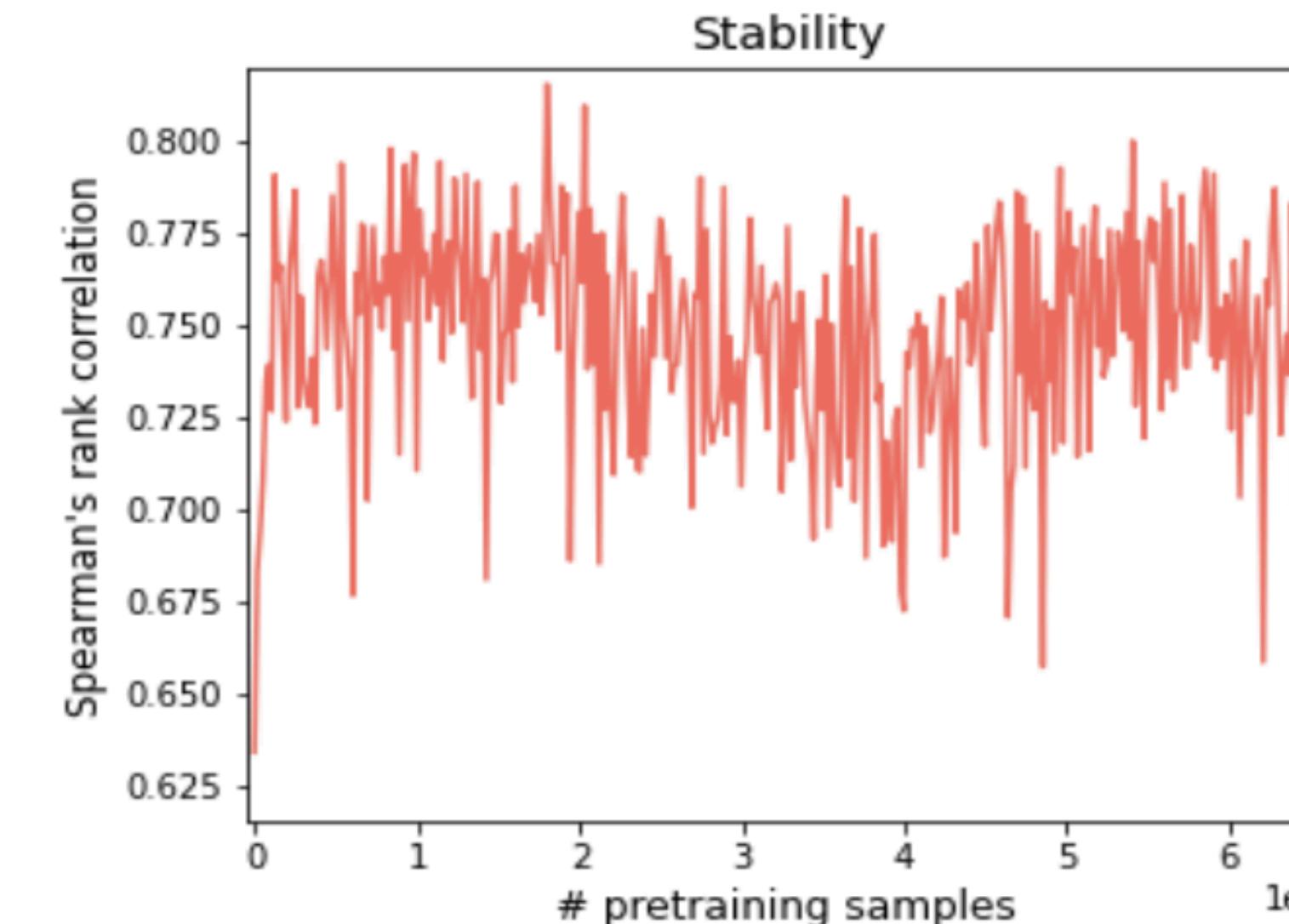
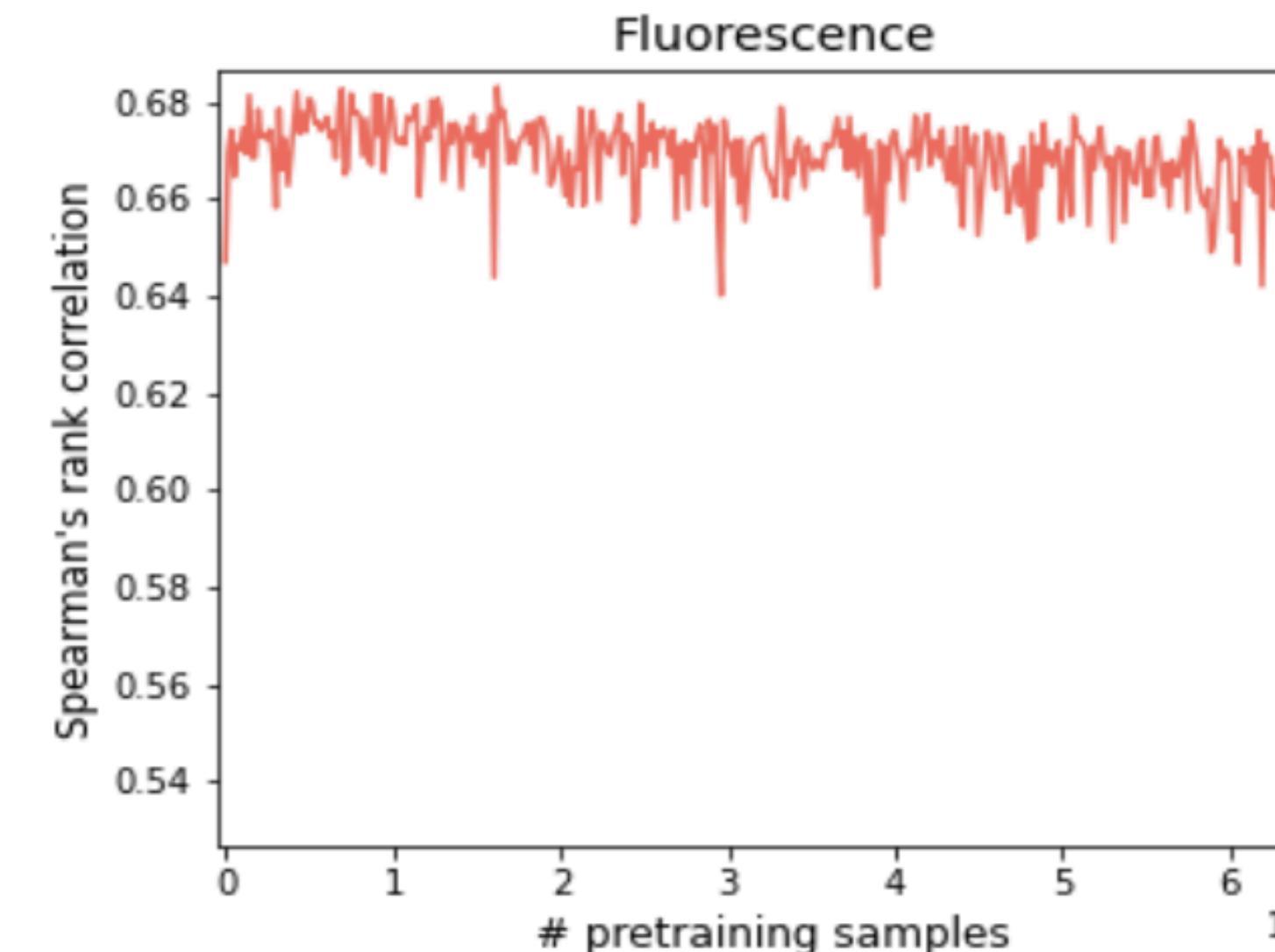
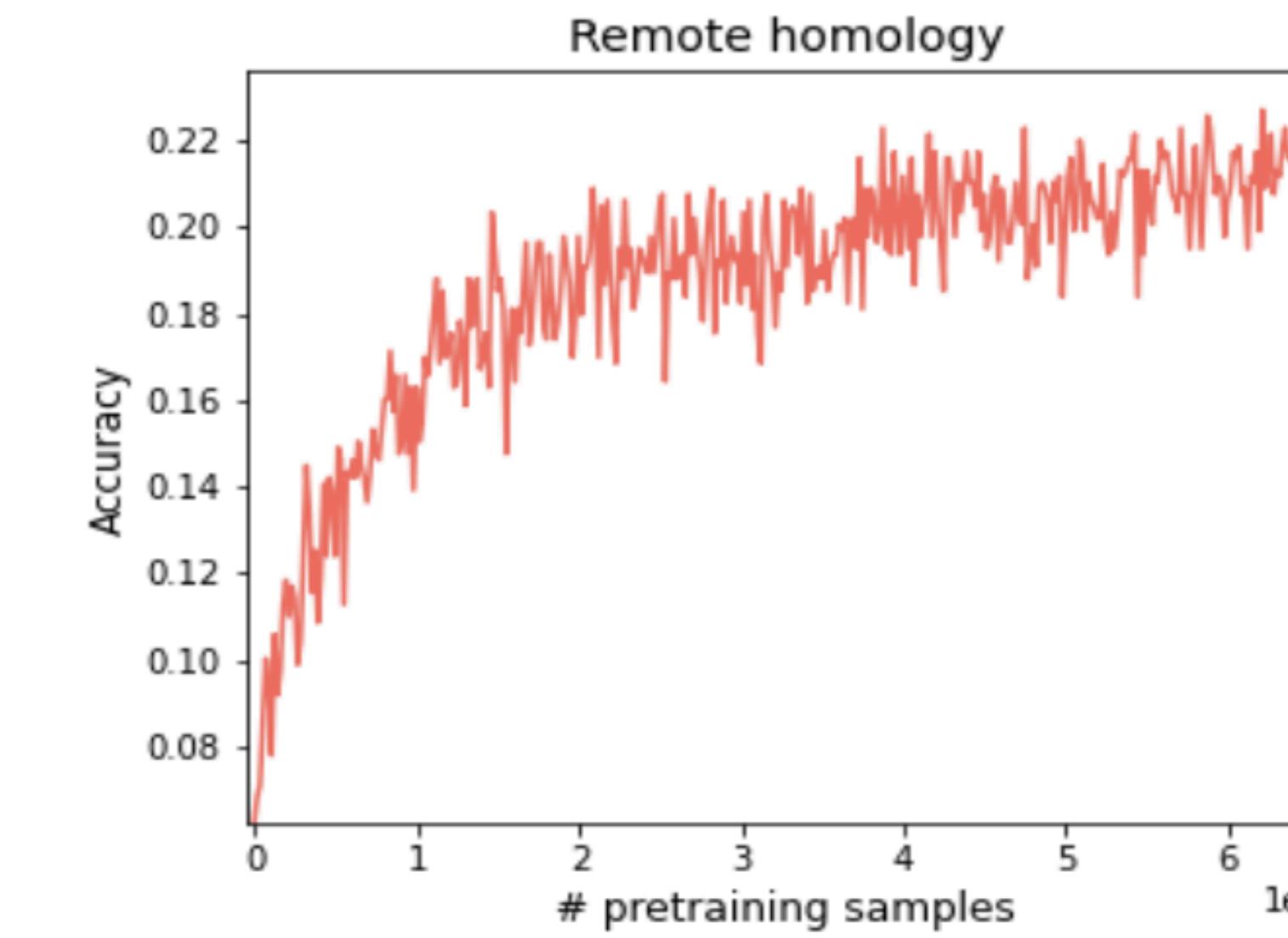
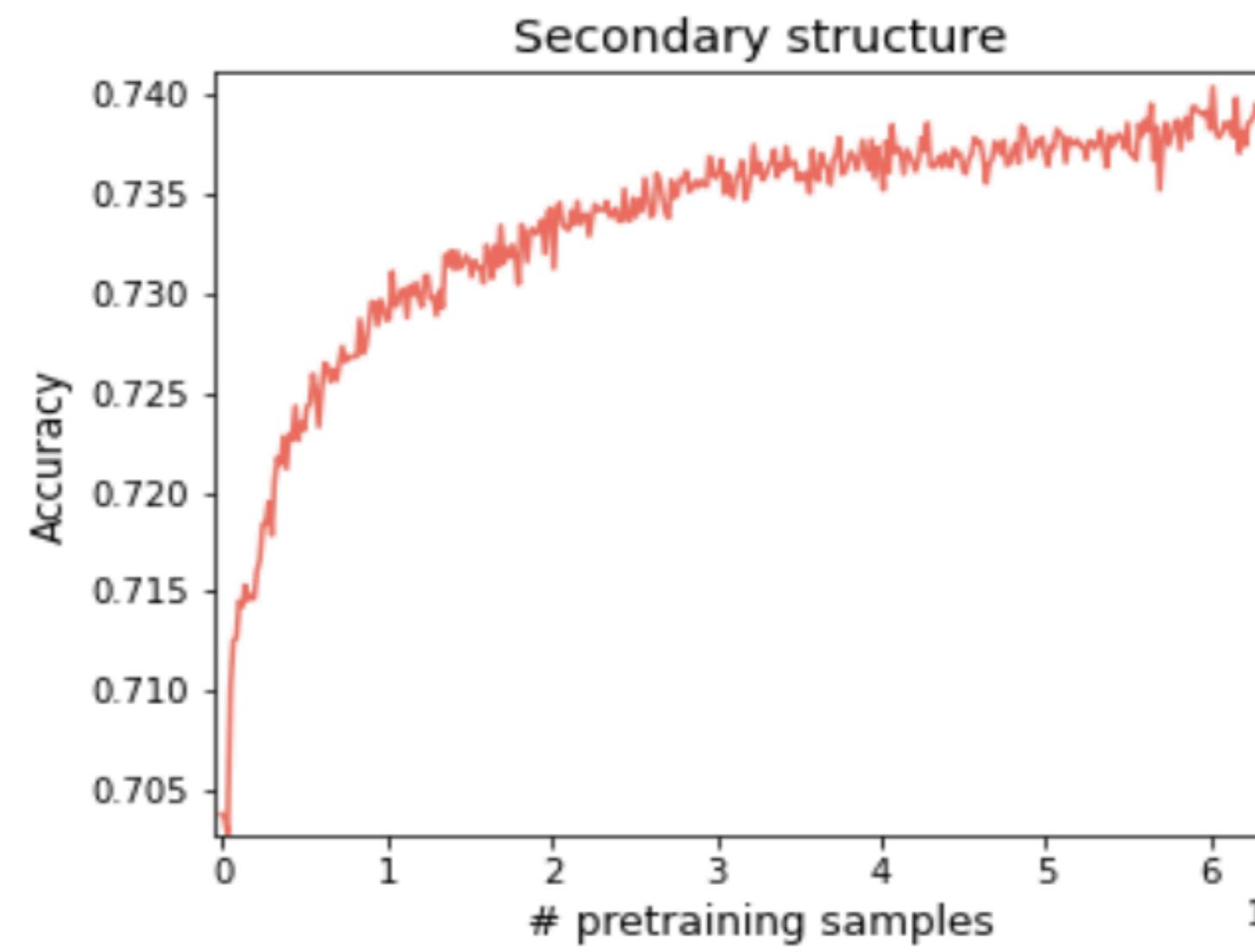
# When does pretraining help?



# When does pretraining help?



# When does pretraining help?



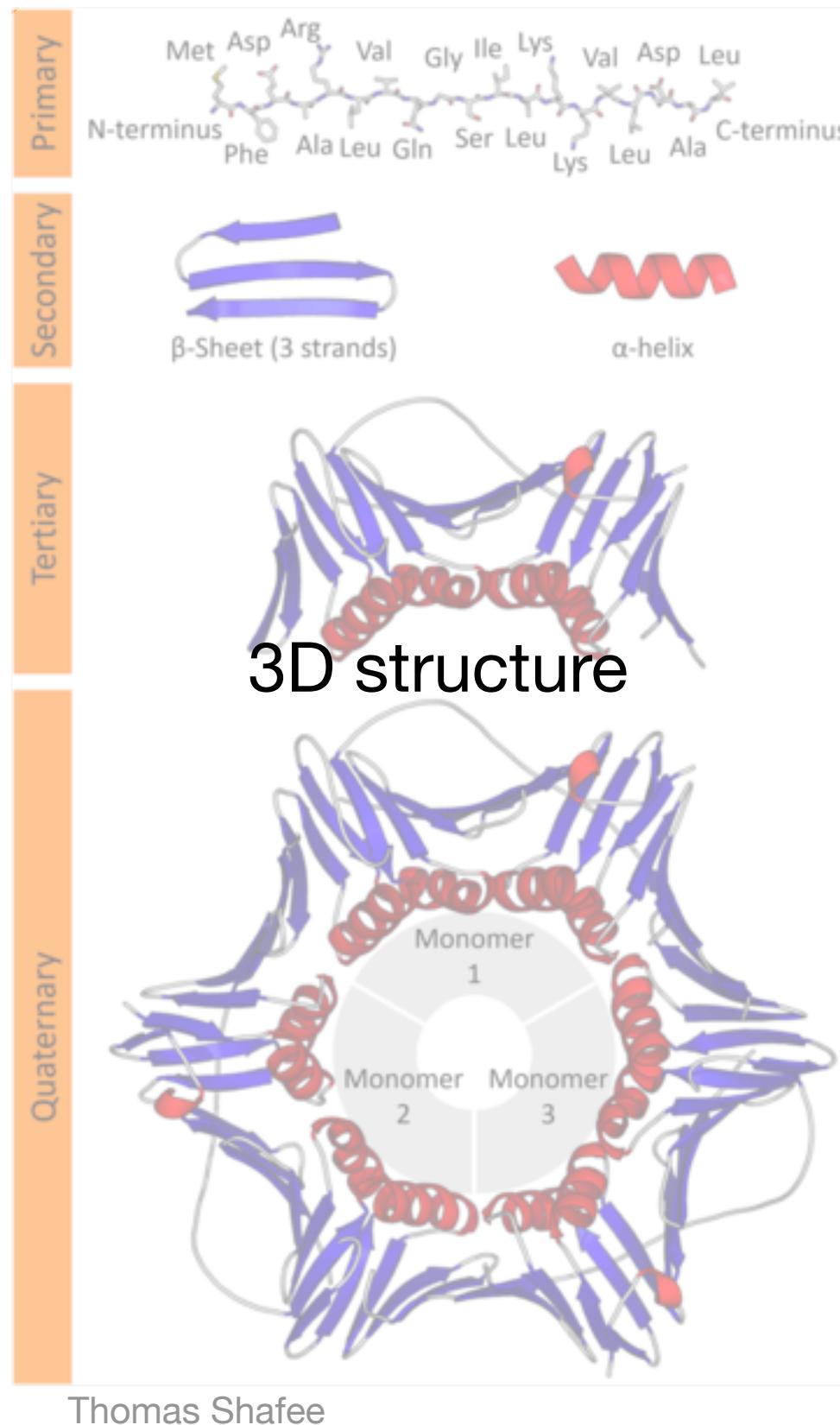
# Protein sequences are not language

# Protein sequences are not language

Differences between proteins and NL

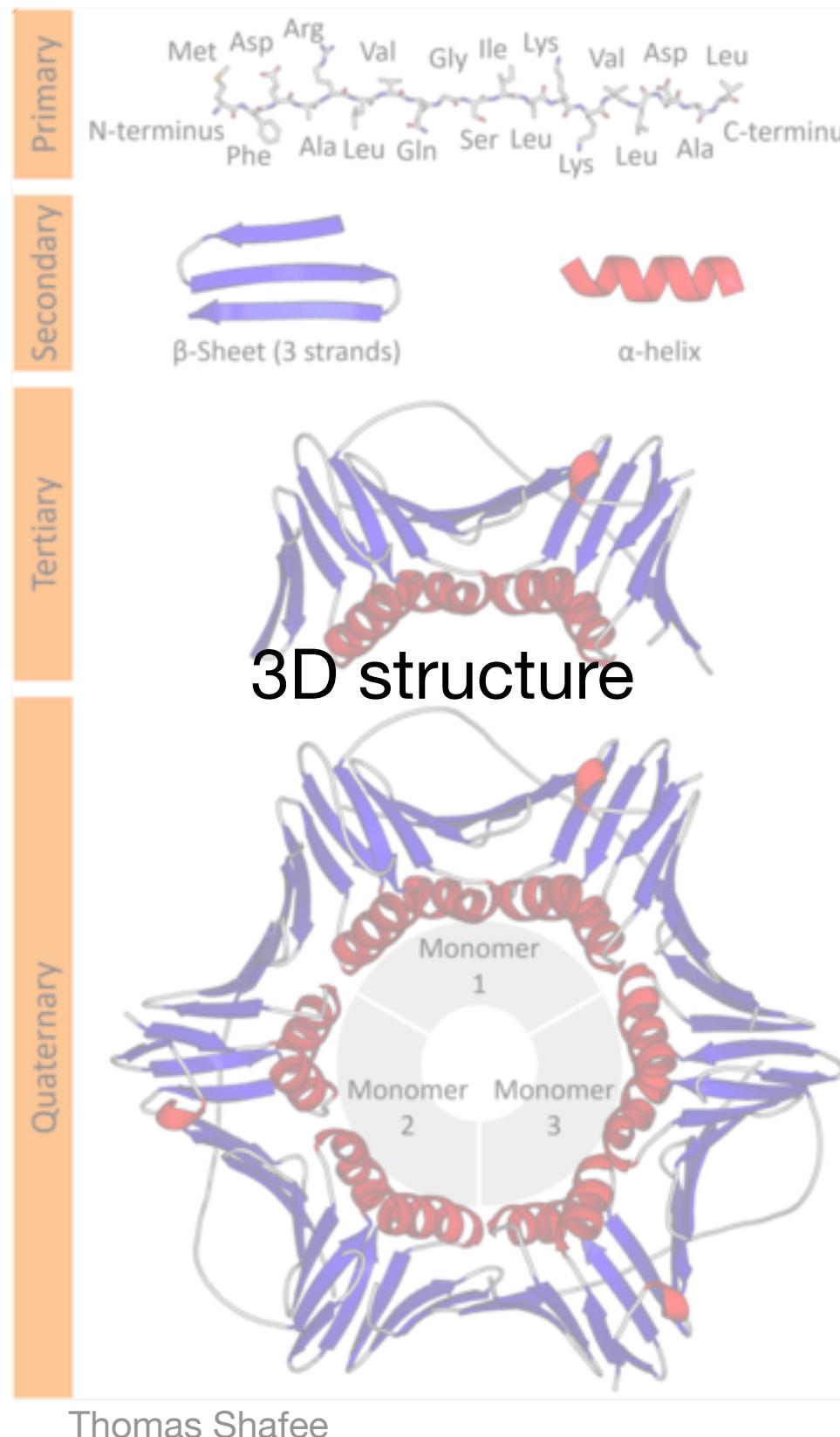
# Protein sequences are not language

Differences between proteins and NL



# Protein sequences are not language

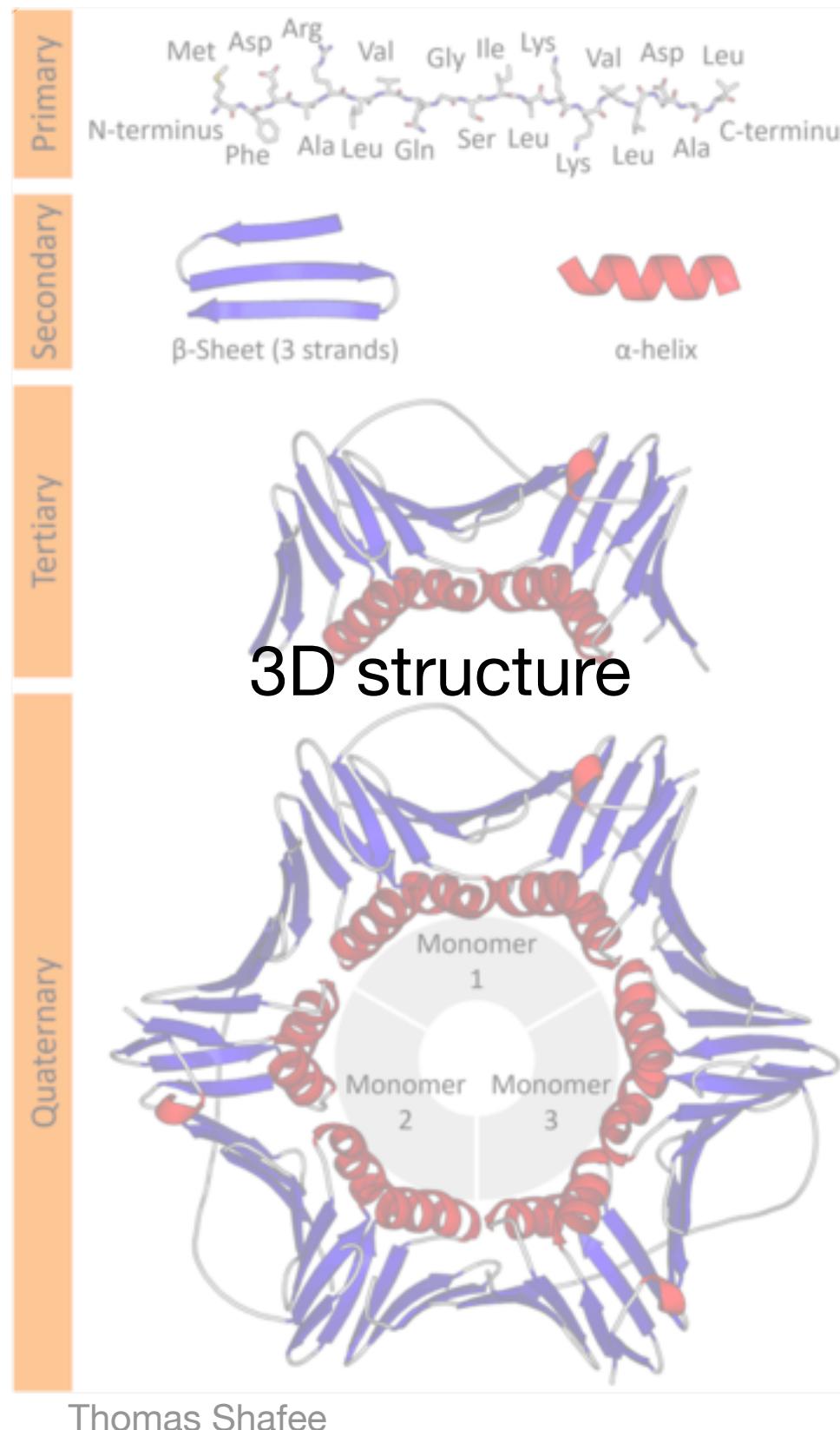
Differences between proteins and NL



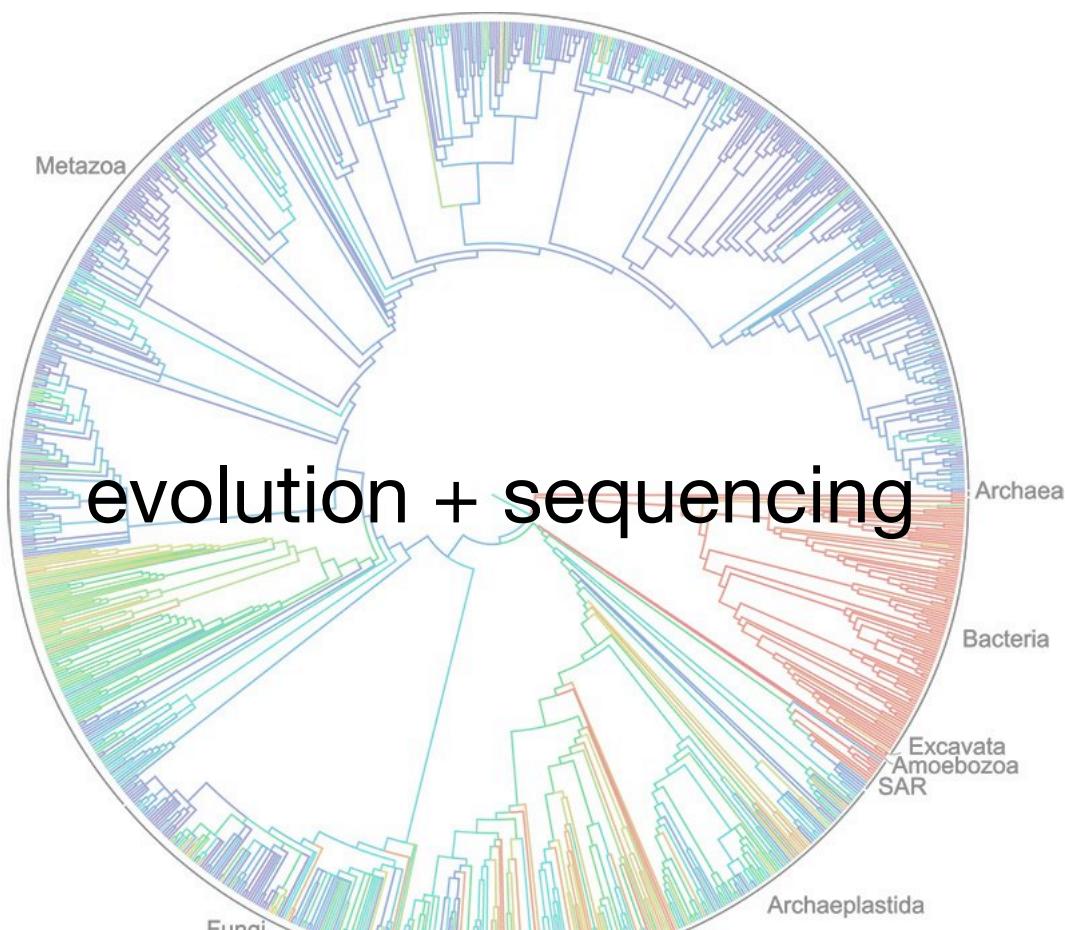
ACDEFGHIKLMNPQRSTVWY  
vs  
words and subwords

# Protein sequences are not language

Differences between proteins and NL

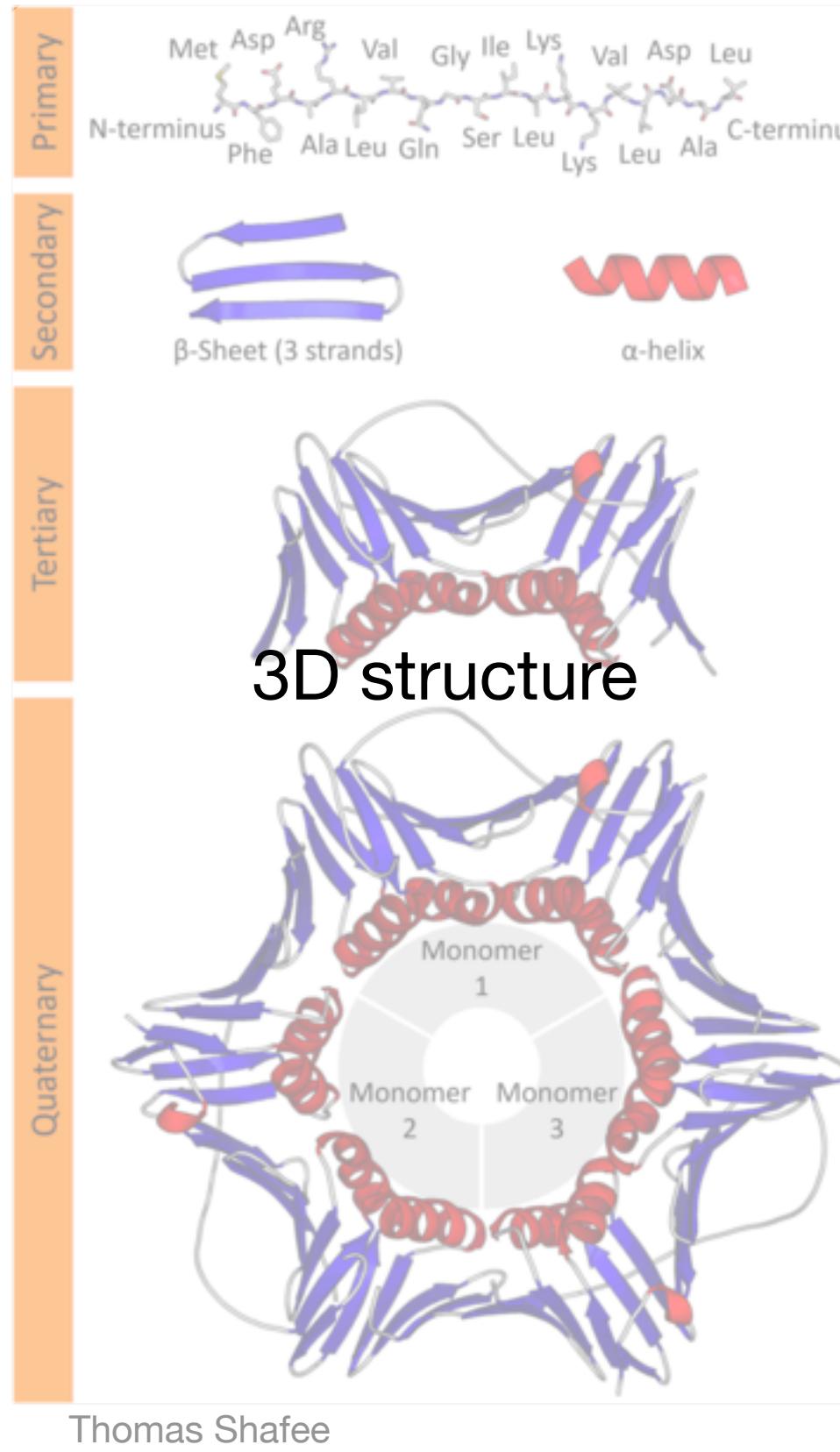


ACDEFGHIKLMNPQRSTVWY  
vs  
words and subwords

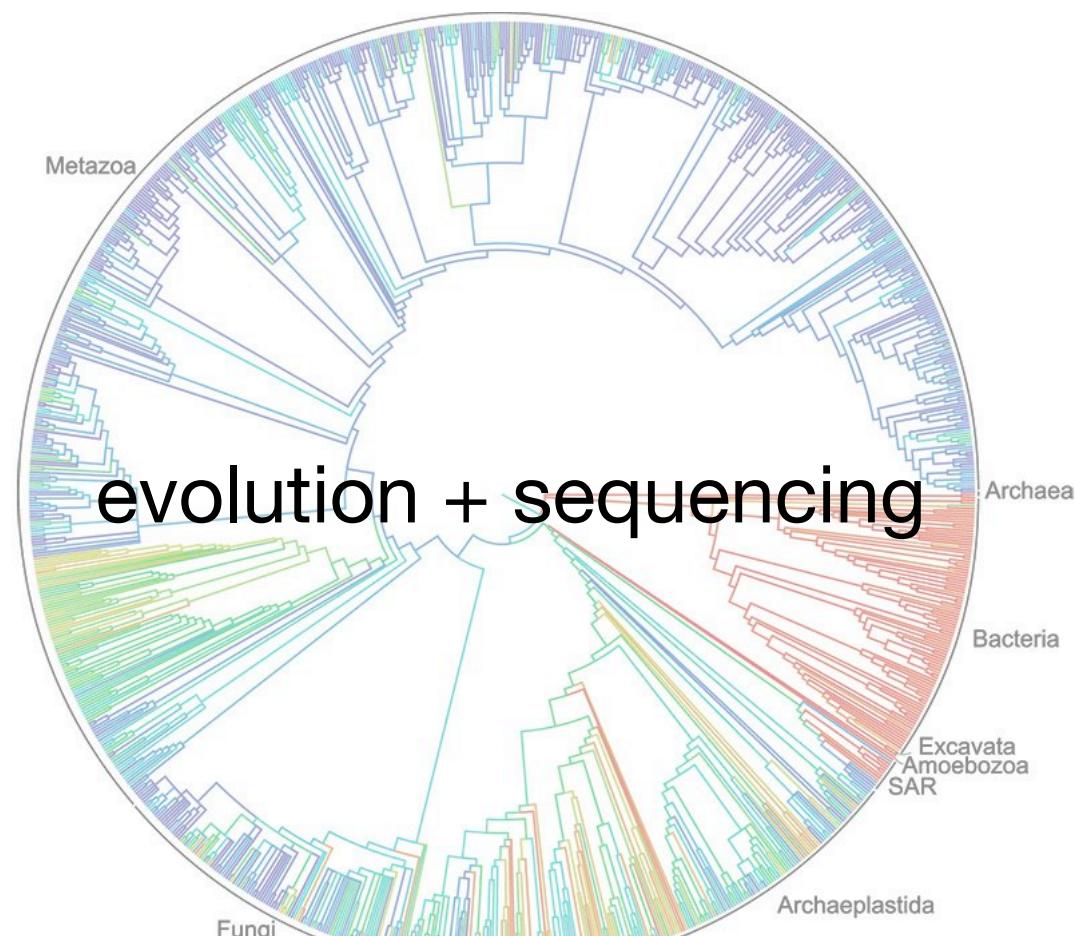


# Protein sequences are not language

## Differences between proteins and NL



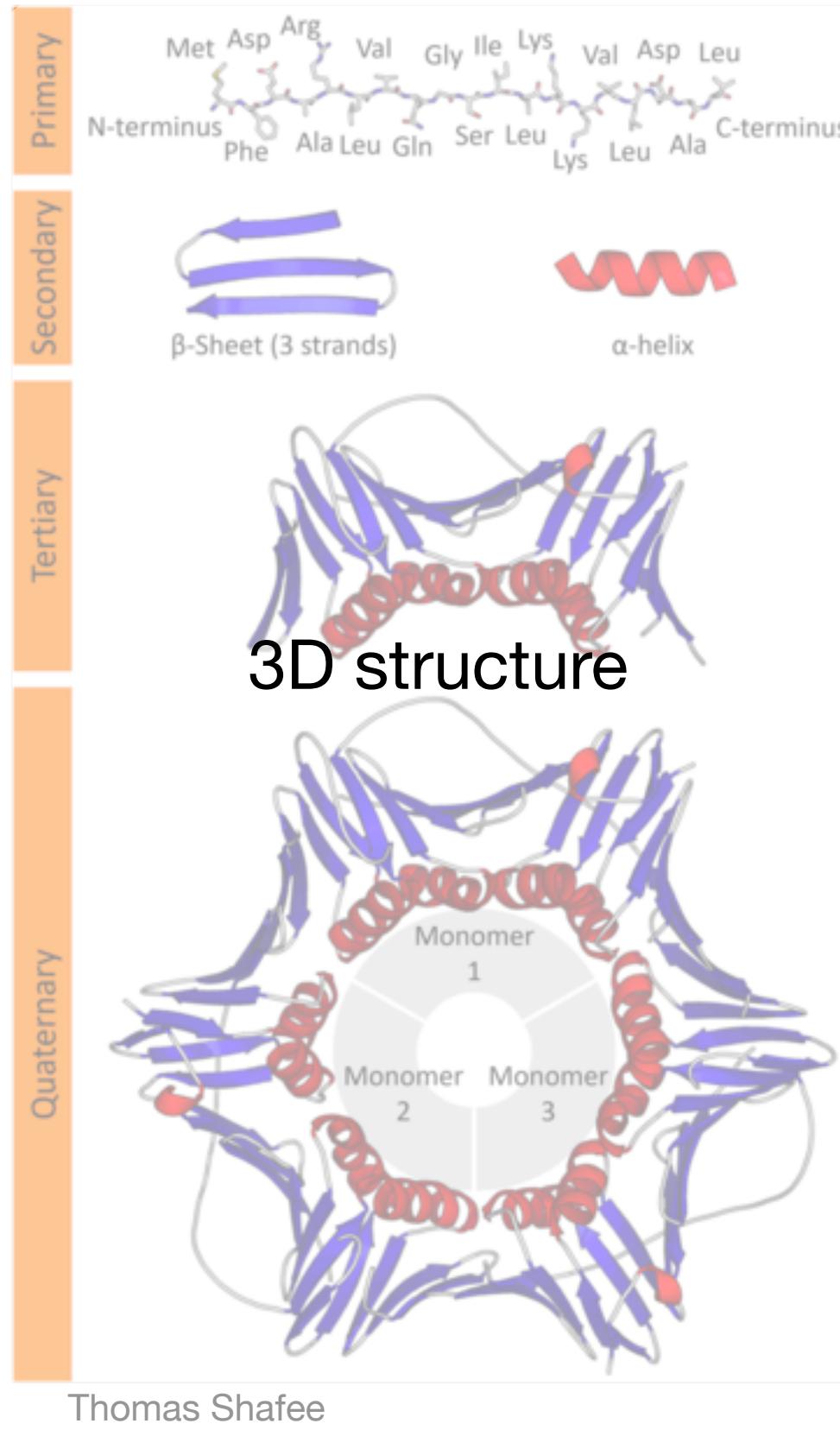
ACDEFGHIKLMNPQRSTVWY  
vs  
words and subwords



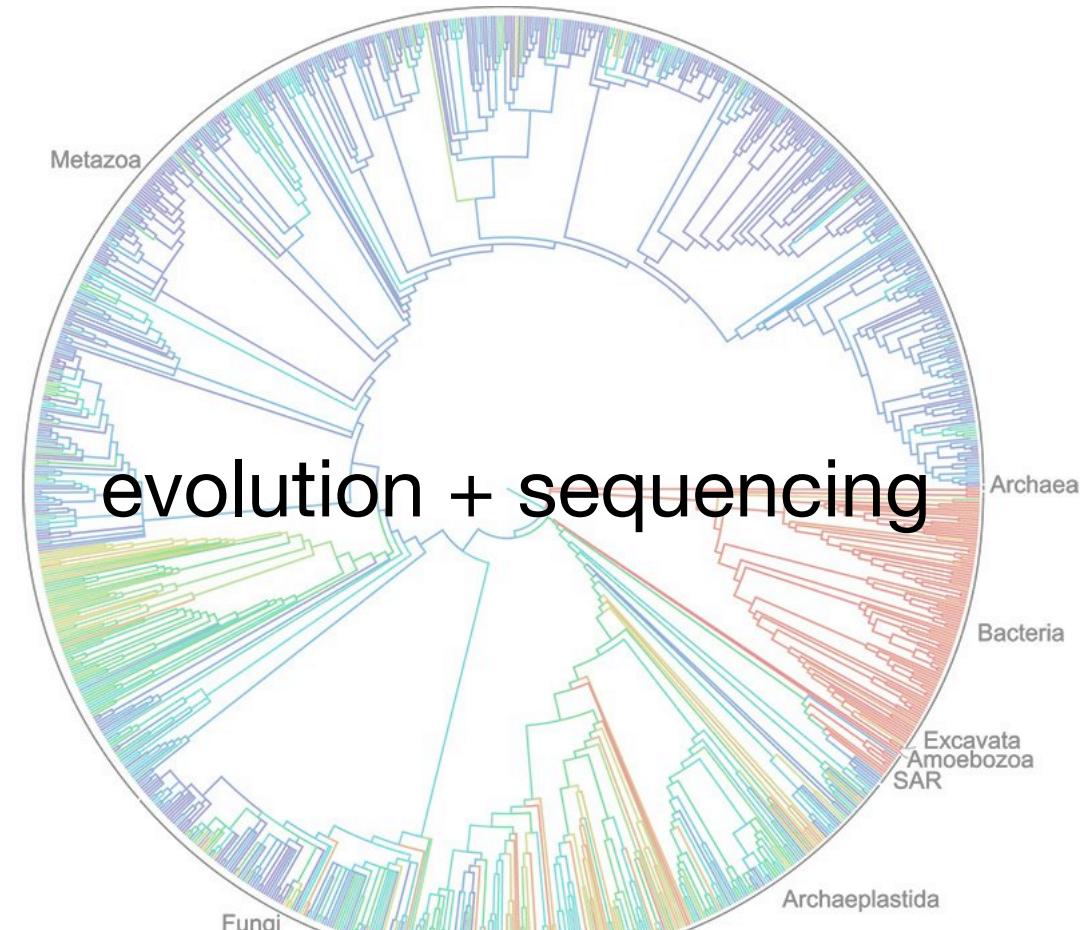
## Multi-modal models

# Protein sequences are not language

## Differences between proteins and NL



ACDEFGHIKLMNPQRSTVWY  
vs  
words and subwords



evolution + sequencing

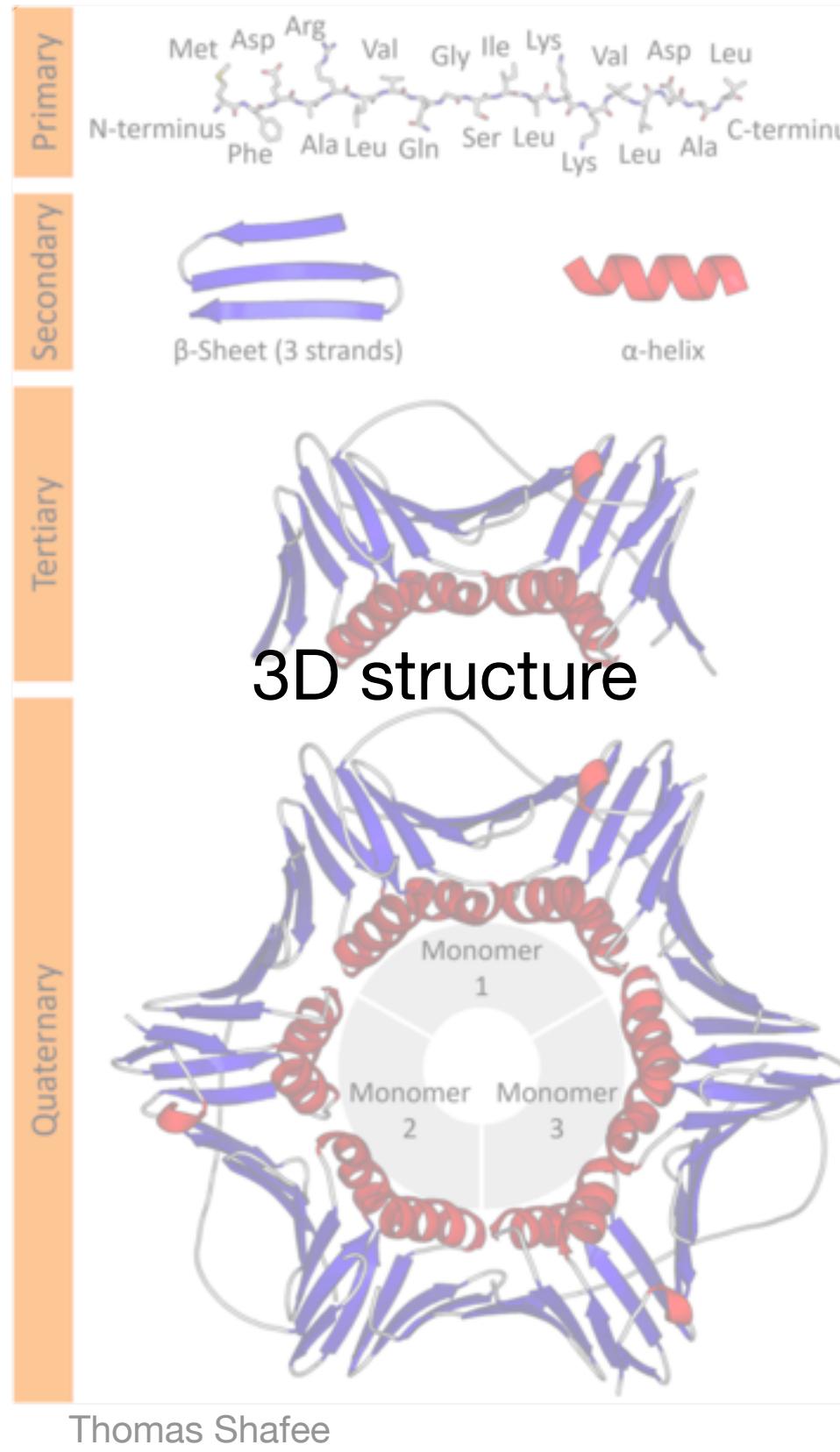
## Multi-modal models

sequence MGHTQWI ...

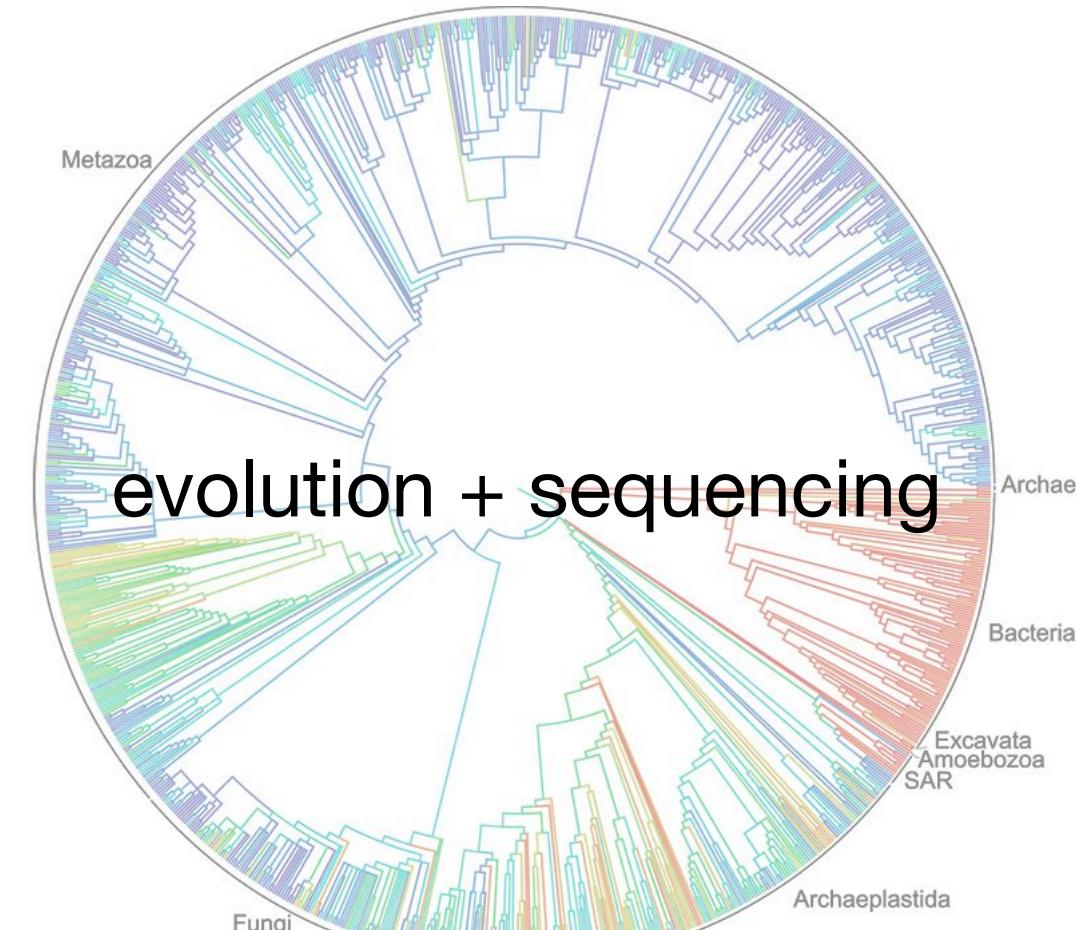
model

# Protein sequences are not language

## Differences between proteins and NL

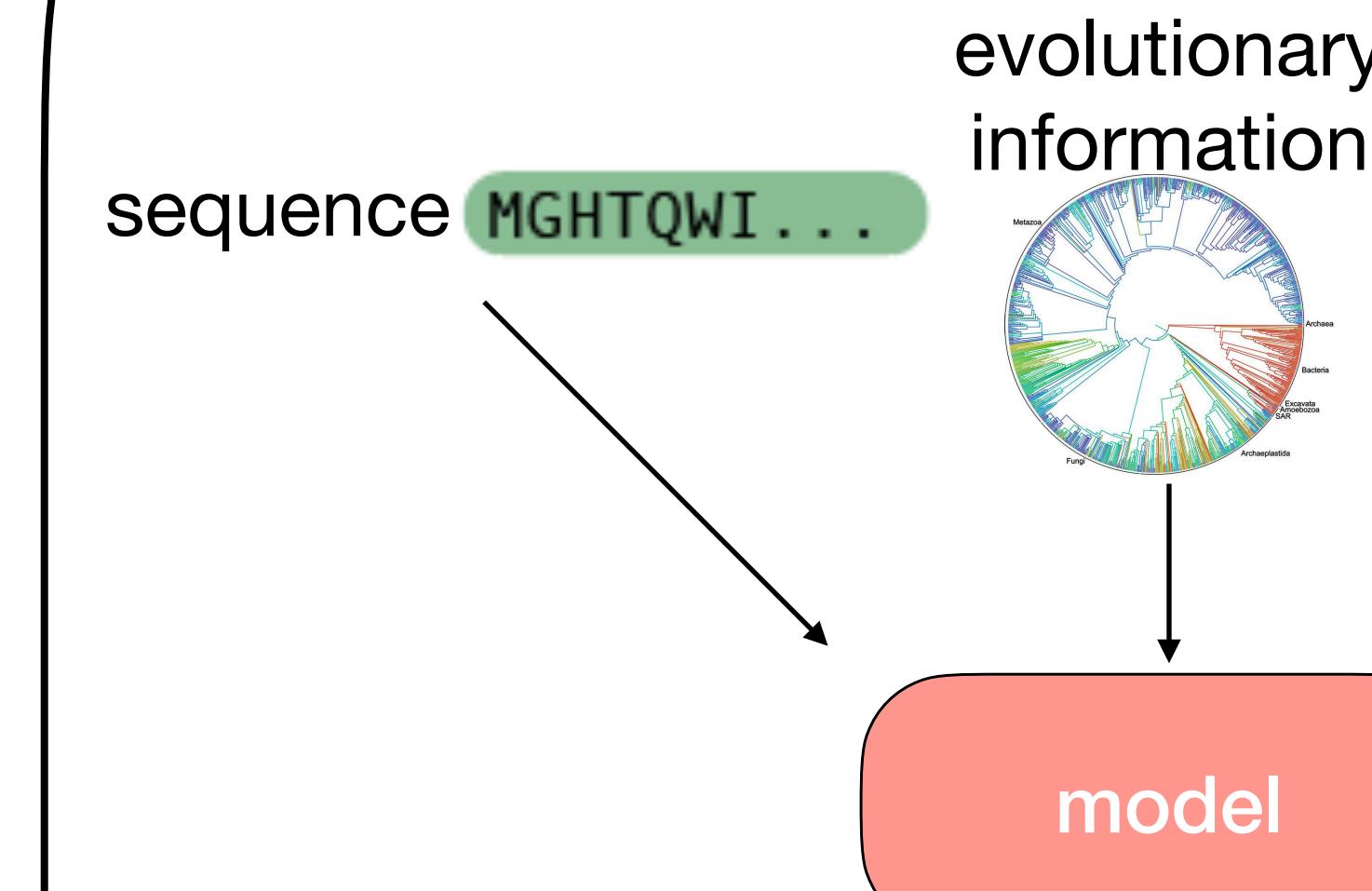


ACDEFGHIKLMNPQRSTVWY  
vs  
words and subwords



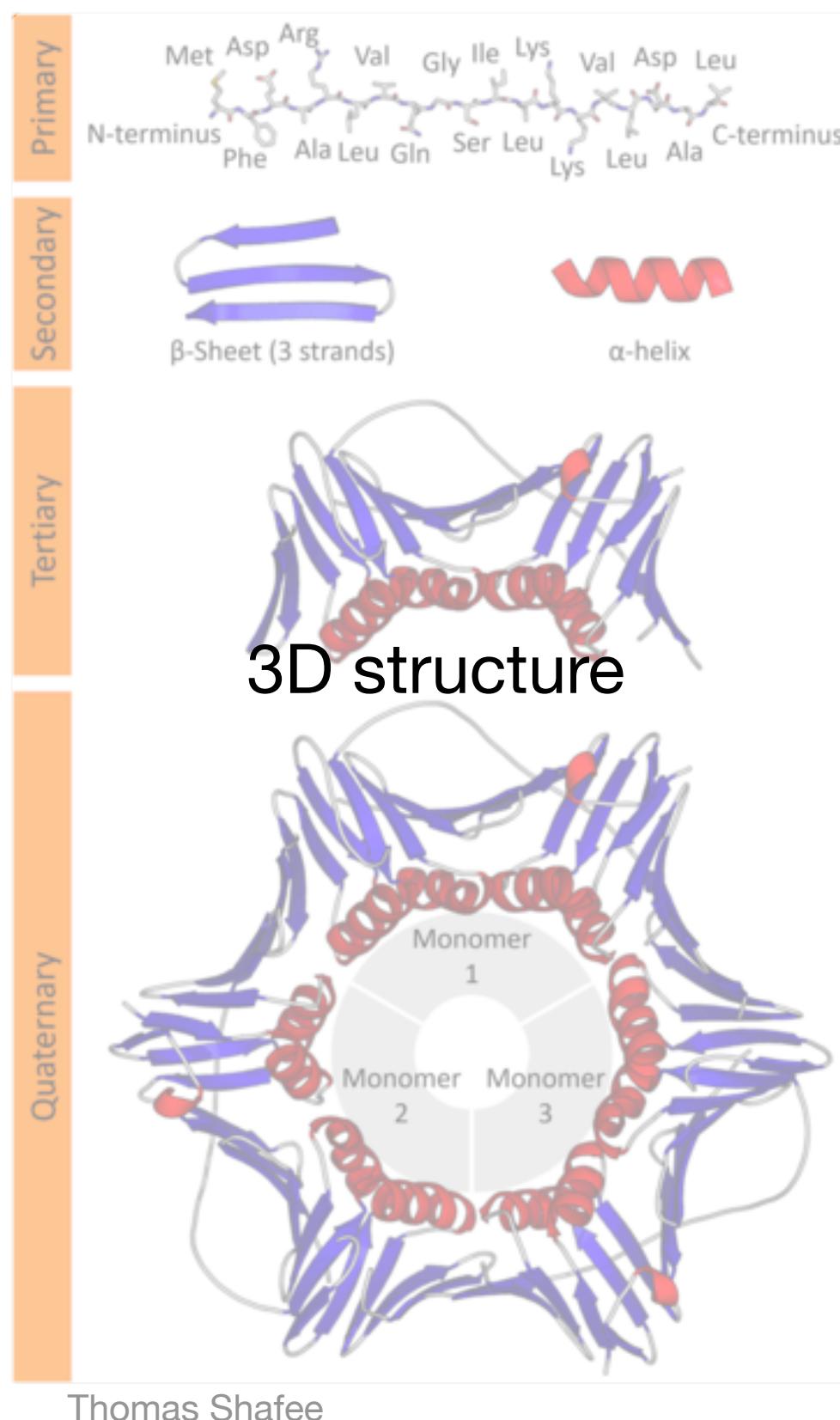
evolution + sequencing

## Multi-modal models

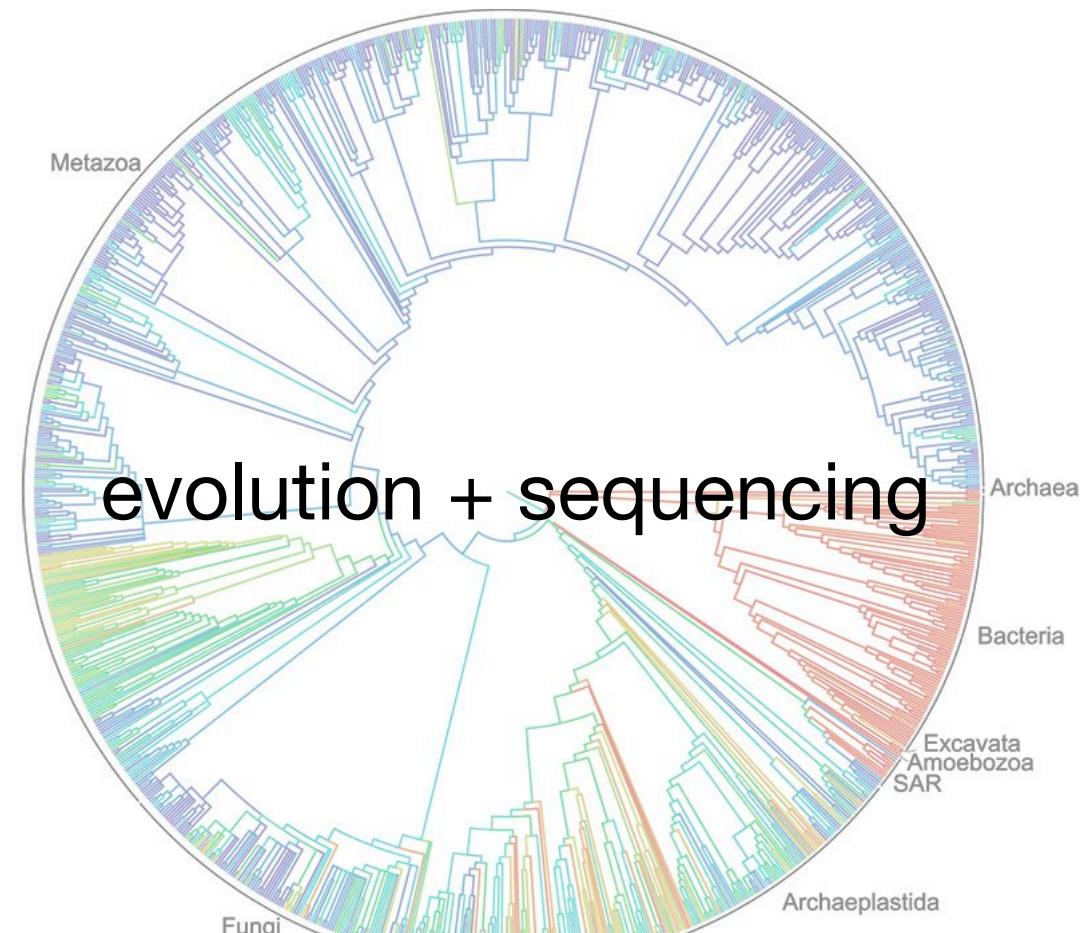


# Protein sequences are not language

## Differences between proteins and NL

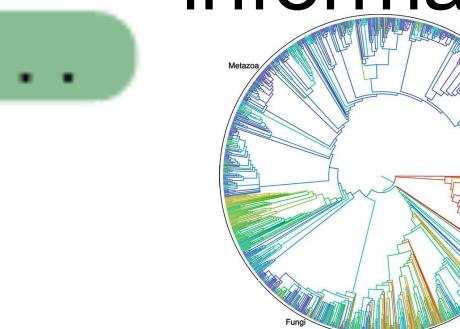


ACDEFGHIKLMNPQRSTVWY  
vs  
words and subwords

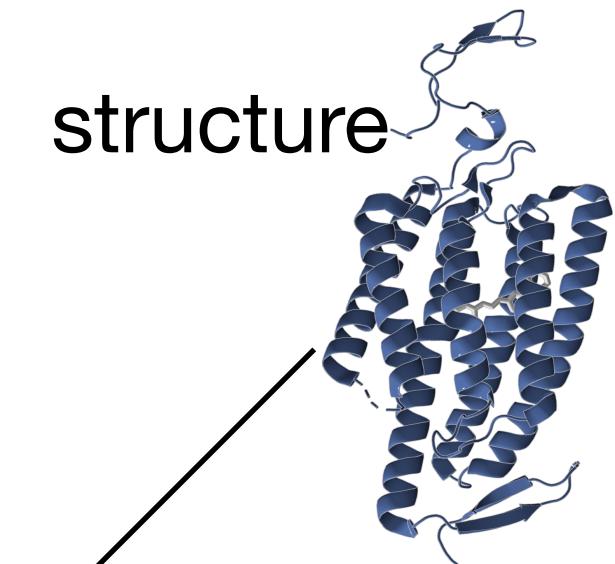


## Multi-modal models

sequence MGHTQWI ...  
evolutionary information

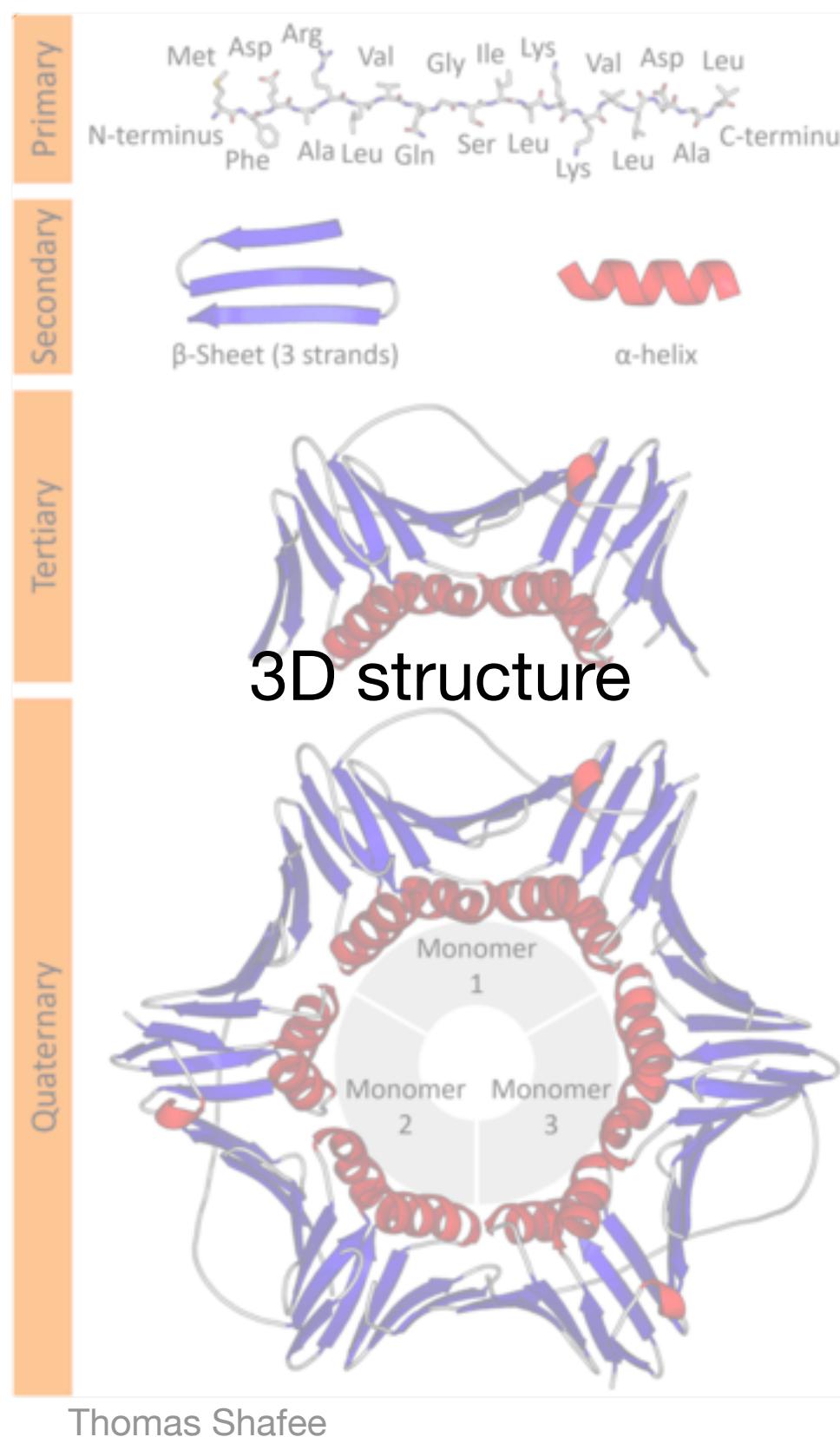


model

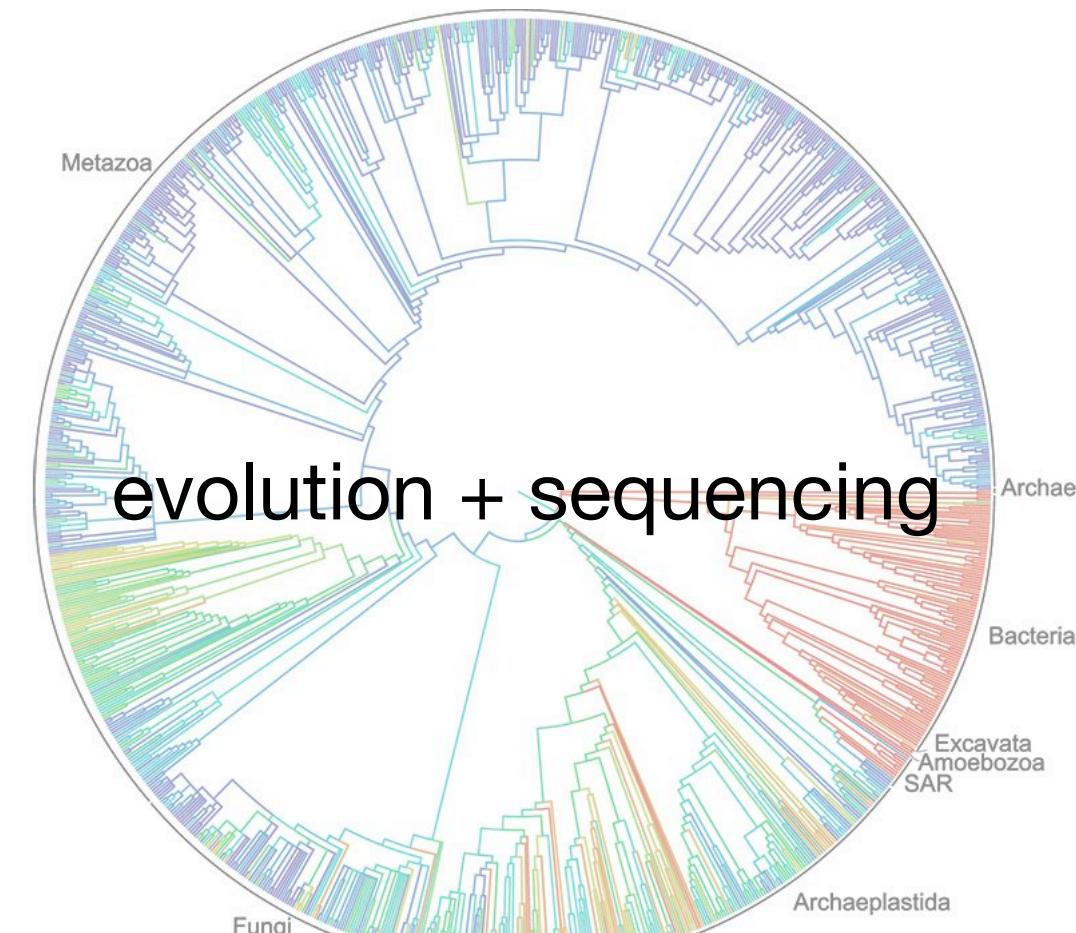


# Protein sequences are not language

## Differences between proteins and NL



ACDEFGHIKLMNPQRSTVWY  
vs  
words and subwords

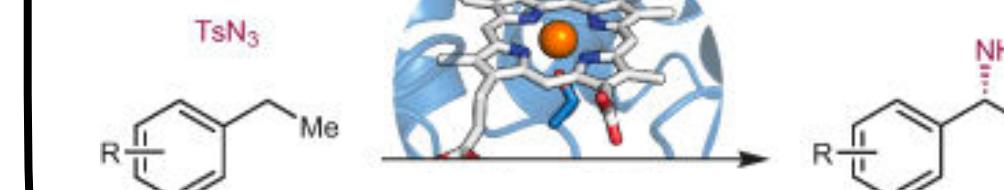


## Multi-modal models

sequence MGHTQWI ...  
evolutionary information

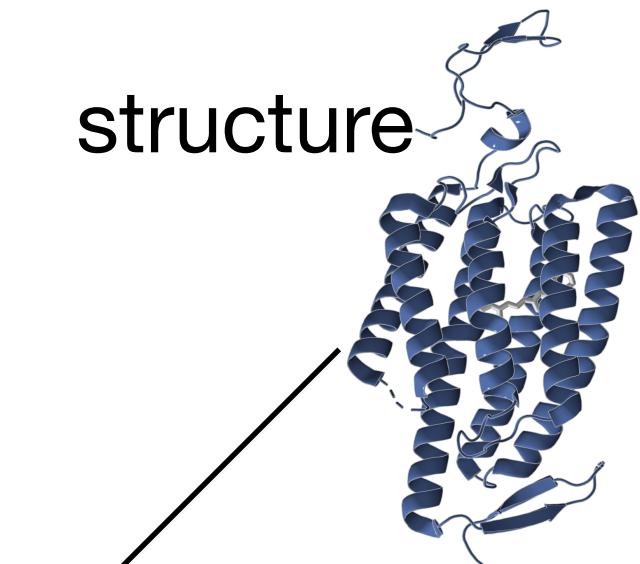
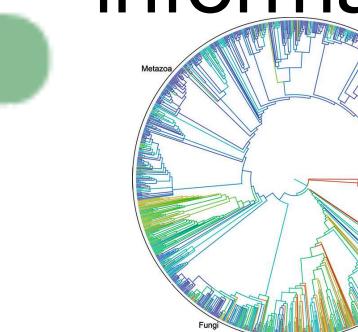
sequence MGHTQWI ...

chemistry



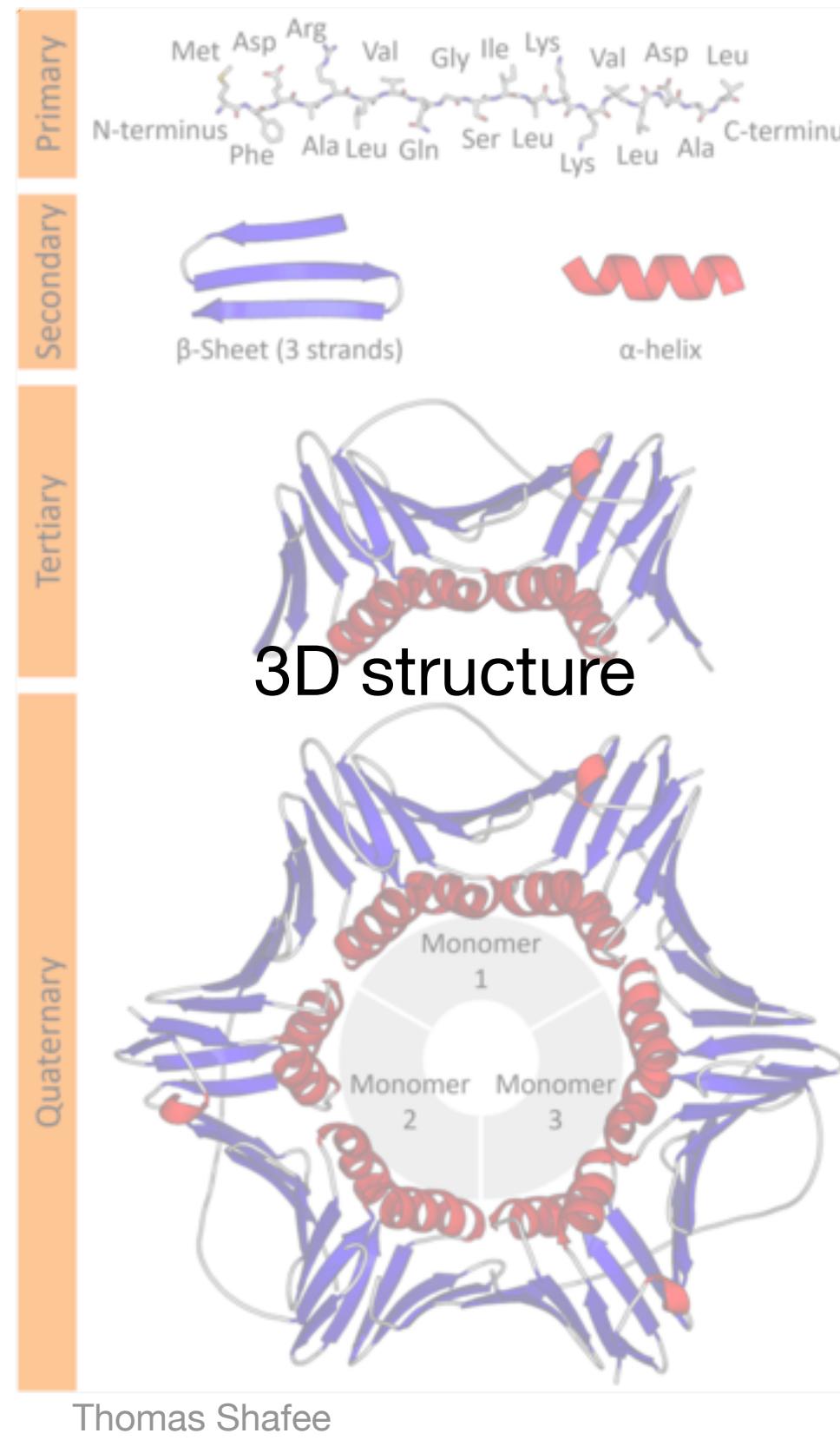
model

evolutionary  
information

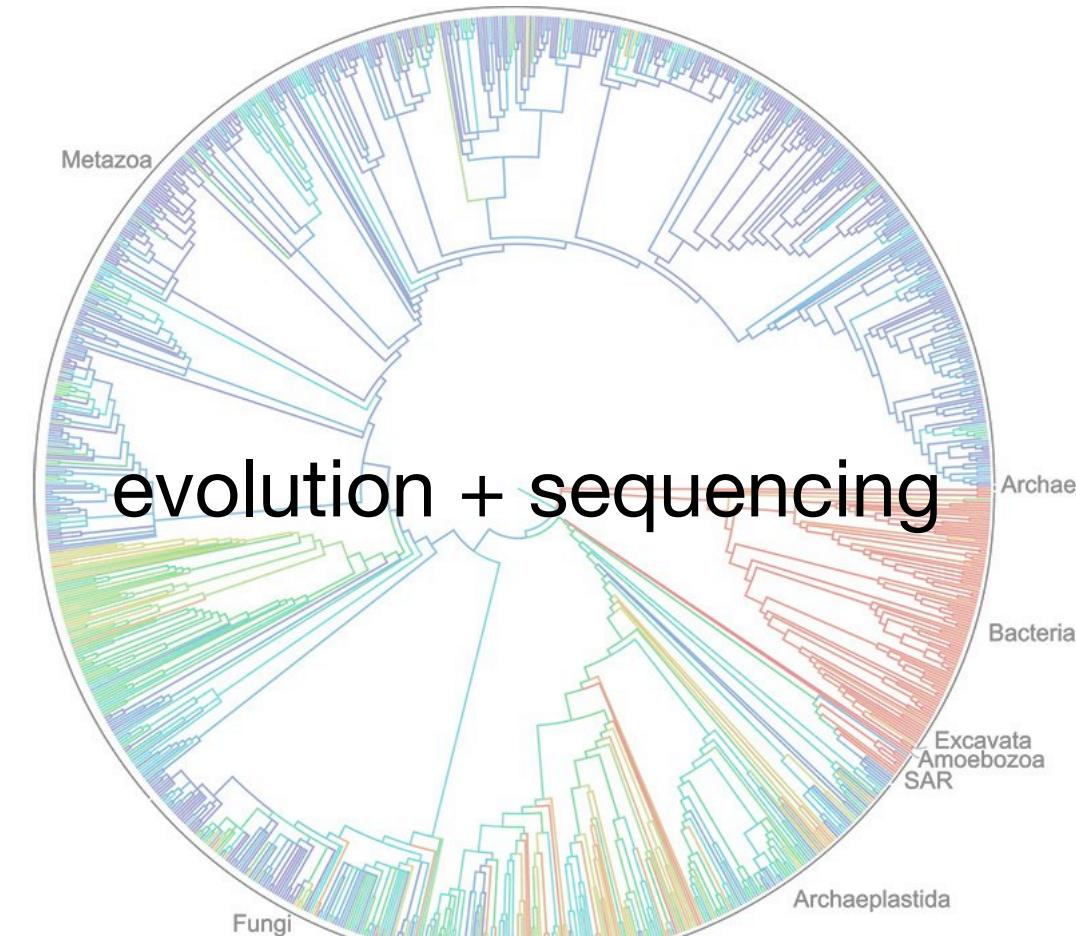


# Protein sequences are not language

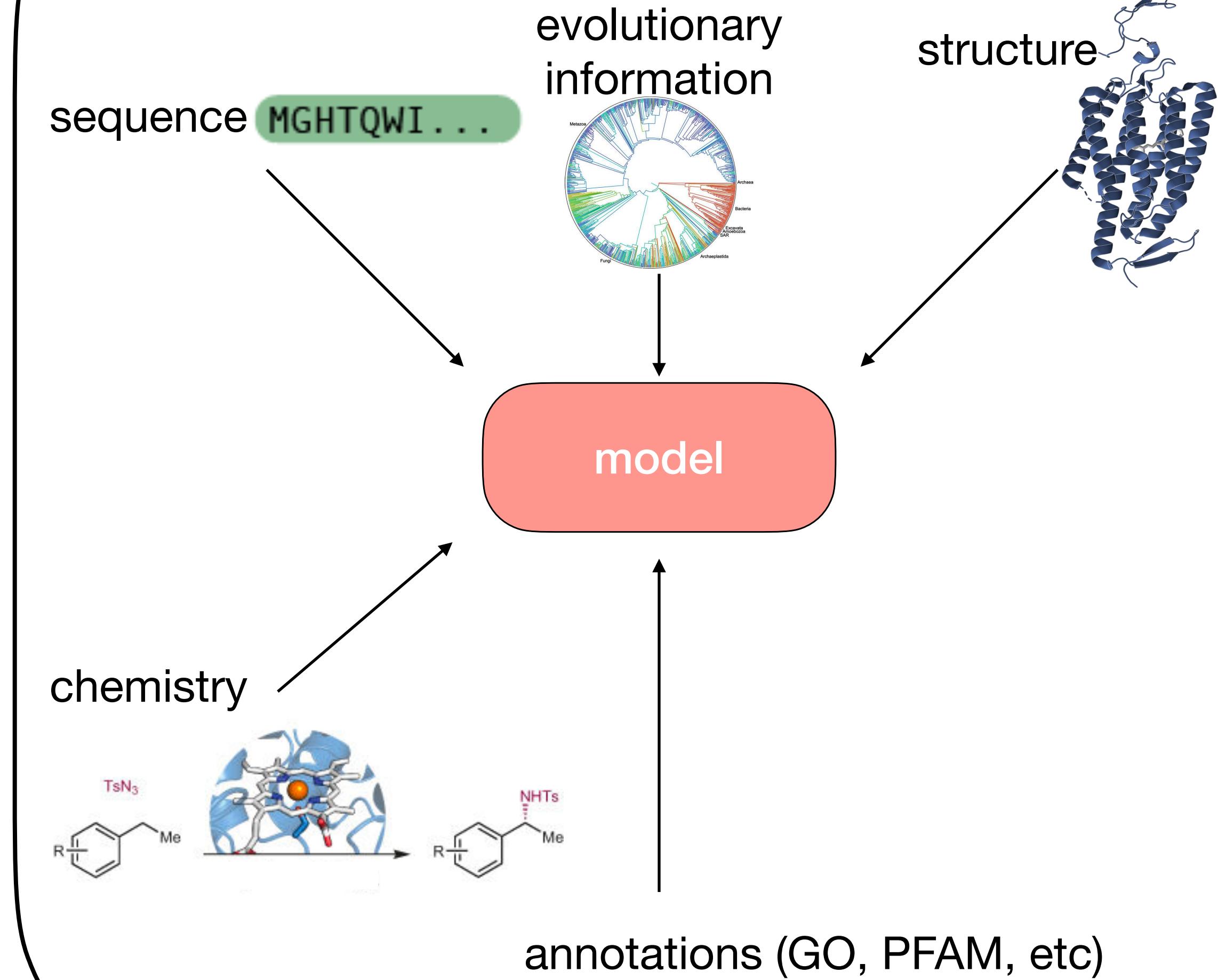
## Differences between proteins and NL



ACDEFGHIKLMNPQRSTVWY  
vs  
words and subwords

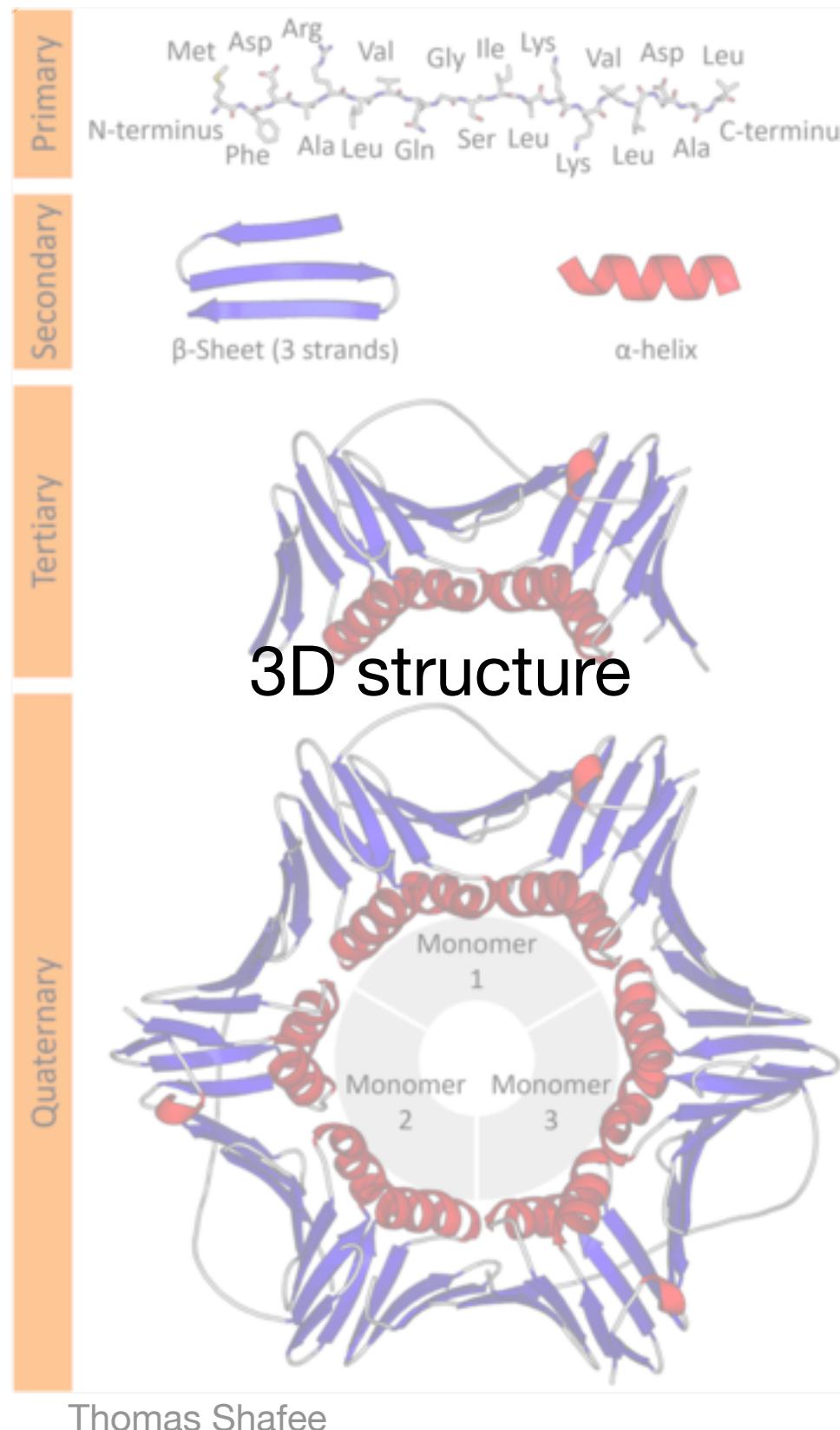


## Multi-modal models

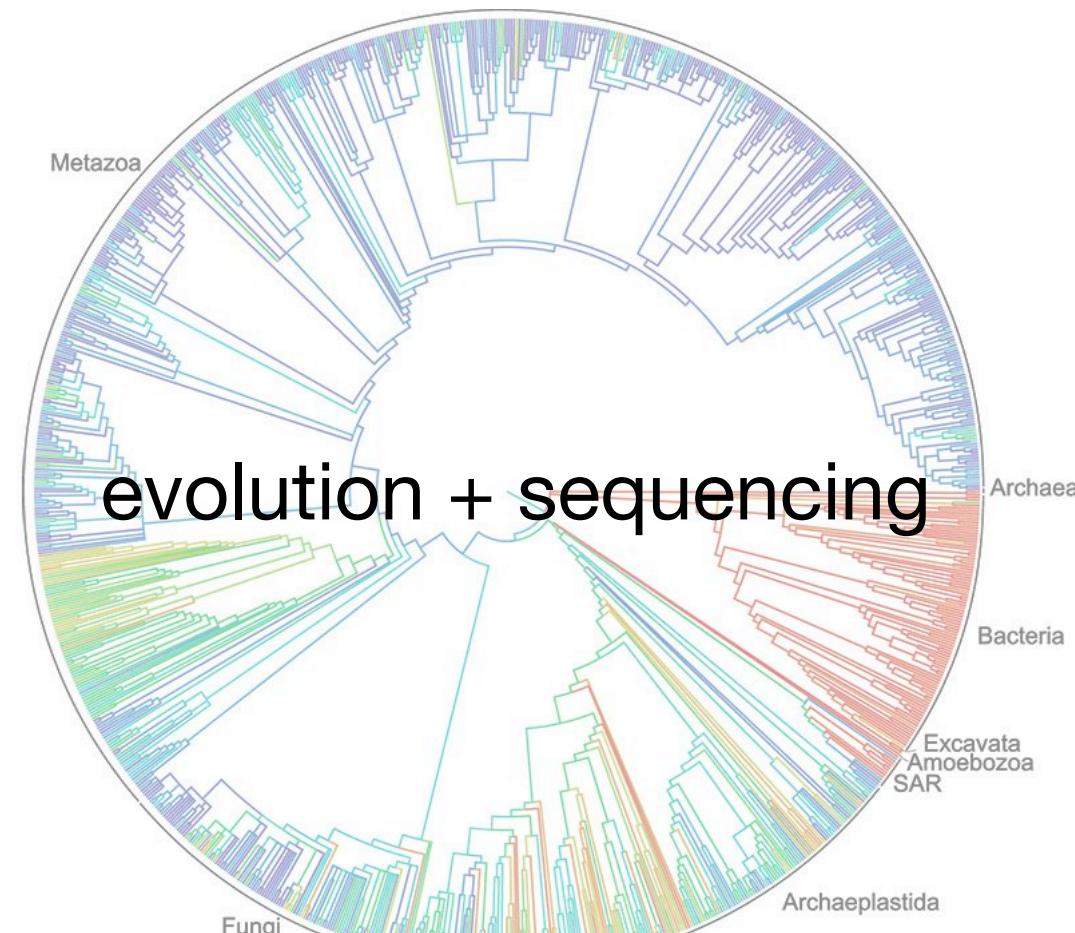


# Protein sequences are not language

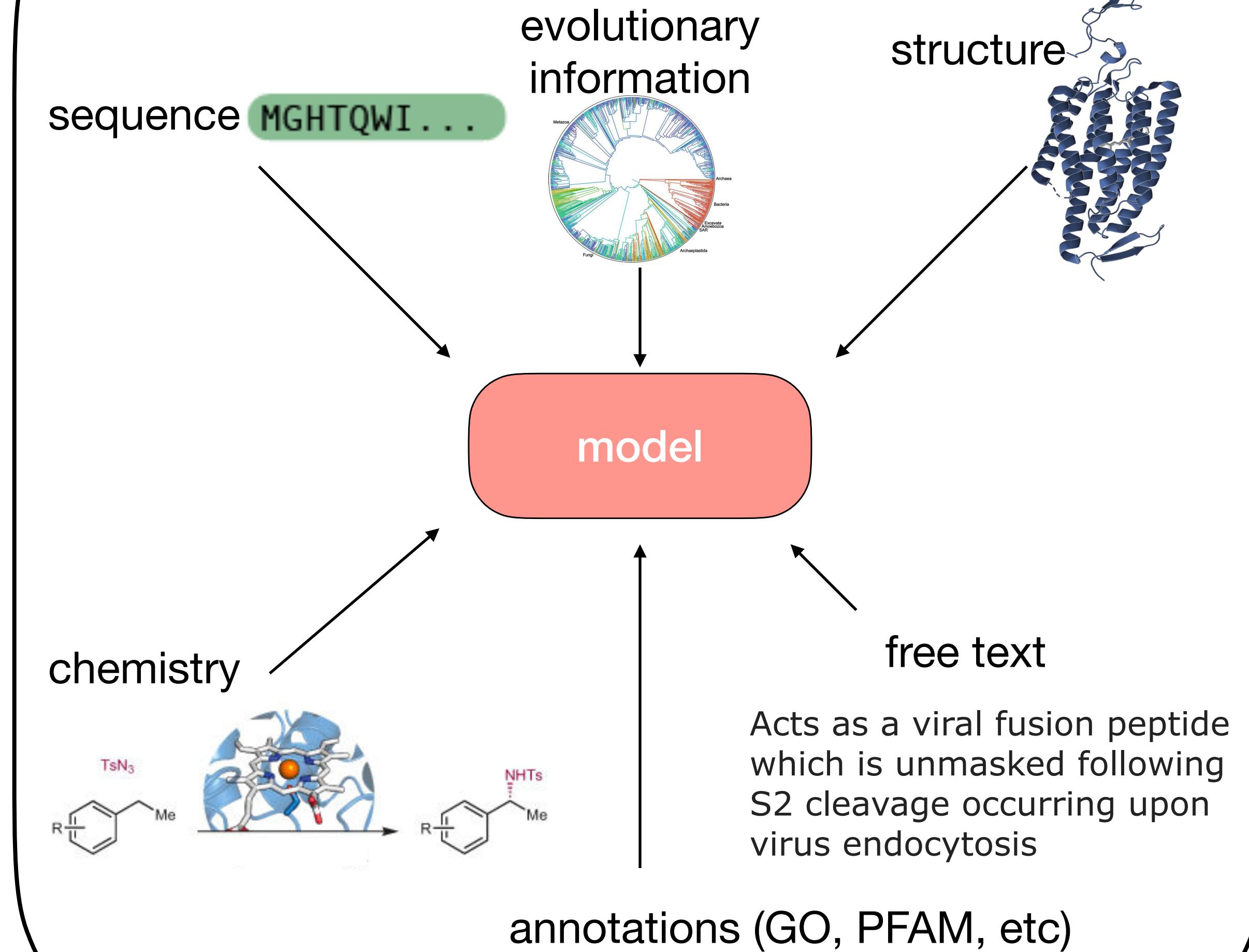
## Differences between proteins and NL



ACDEFGHIKLMNPQRSTVWY  
vs  
words and subwords

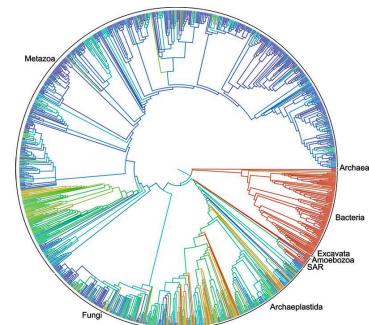


## Multi-modal models



# Multiple sequence alignments encode structural information

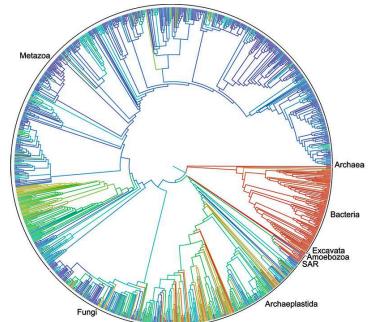
evolutionary  
information



# Multiple sequence alignments encode structural information

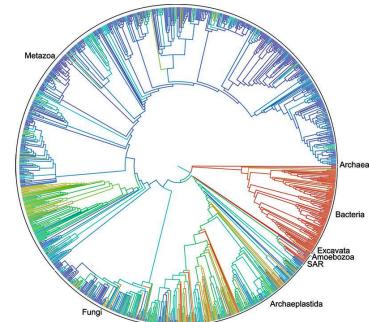
Alignments contain information  
about contacts

evolutionary  
information



# Multiple sequence alignments encode structural information

evolutionary  
information

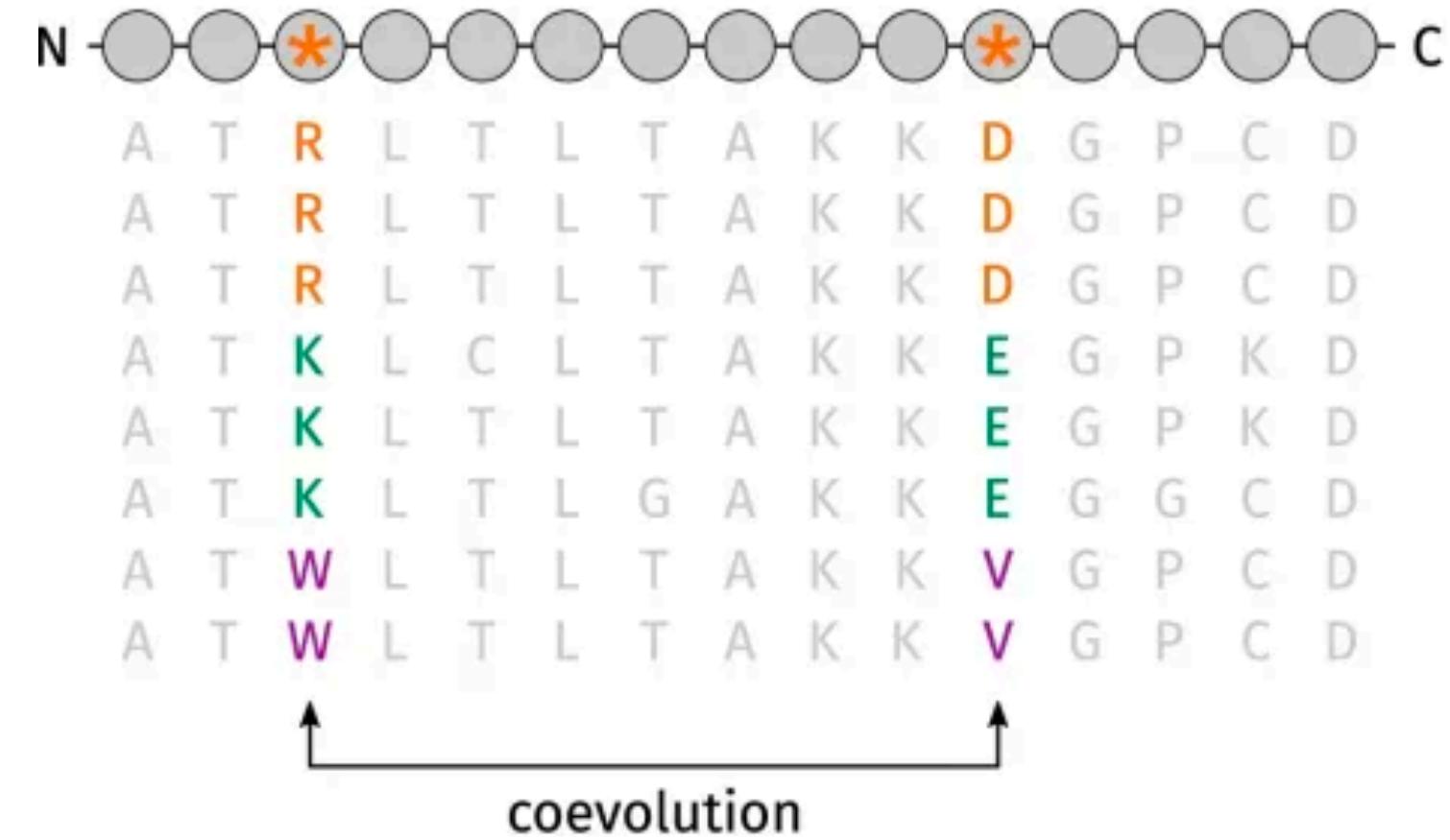


Alignments contain information  
about contacts

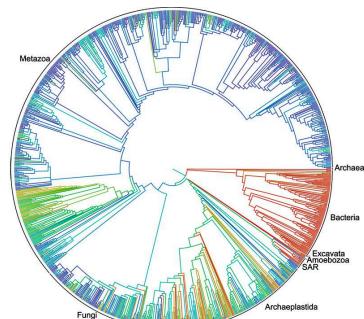


# Multiple sequence alignments encode structural information

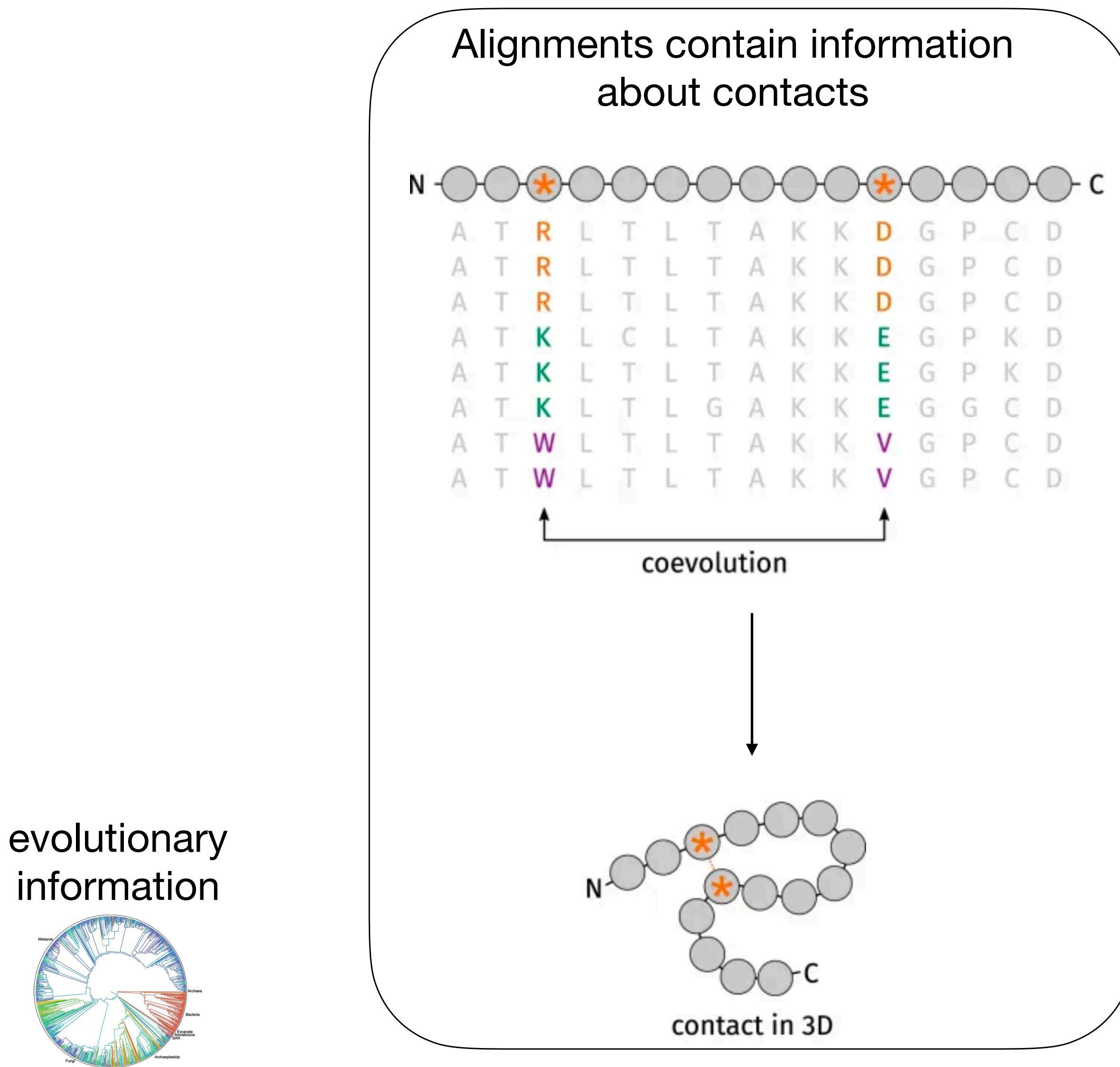
Alignments contain information about contacts



evolutionary information

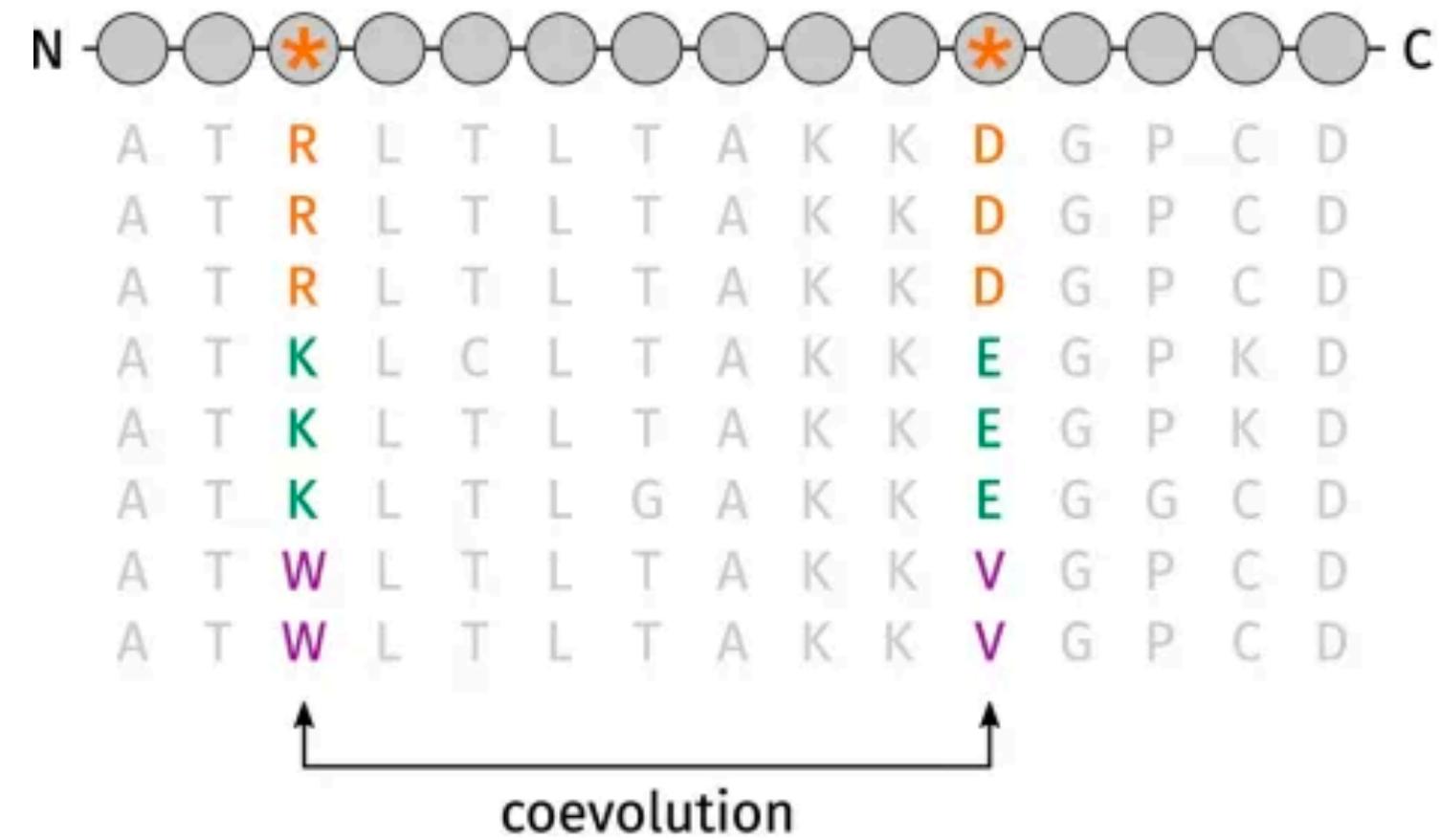


# Multiple sequence alignments encode structural information

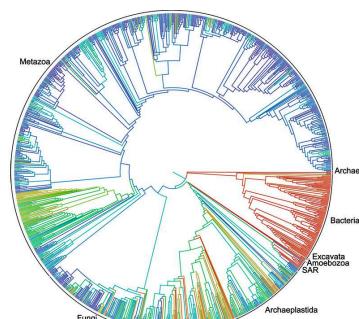


# Multiple sequence alignments encode structural information

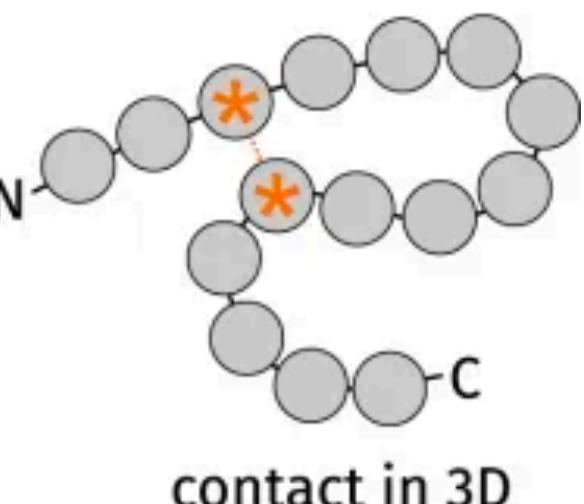
Alignments contain information about contacts



evolutionary information

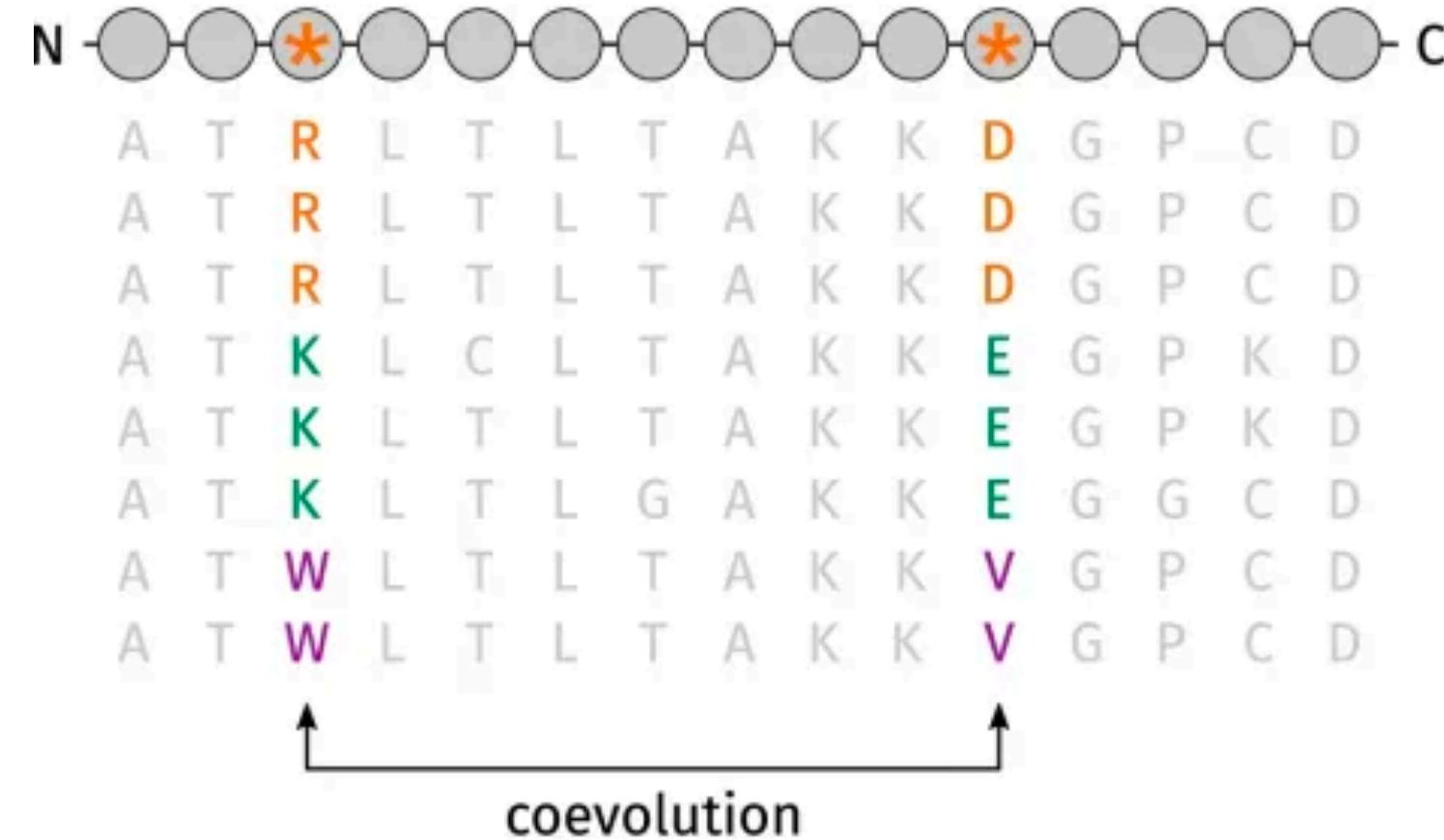


Contacts constrain fold



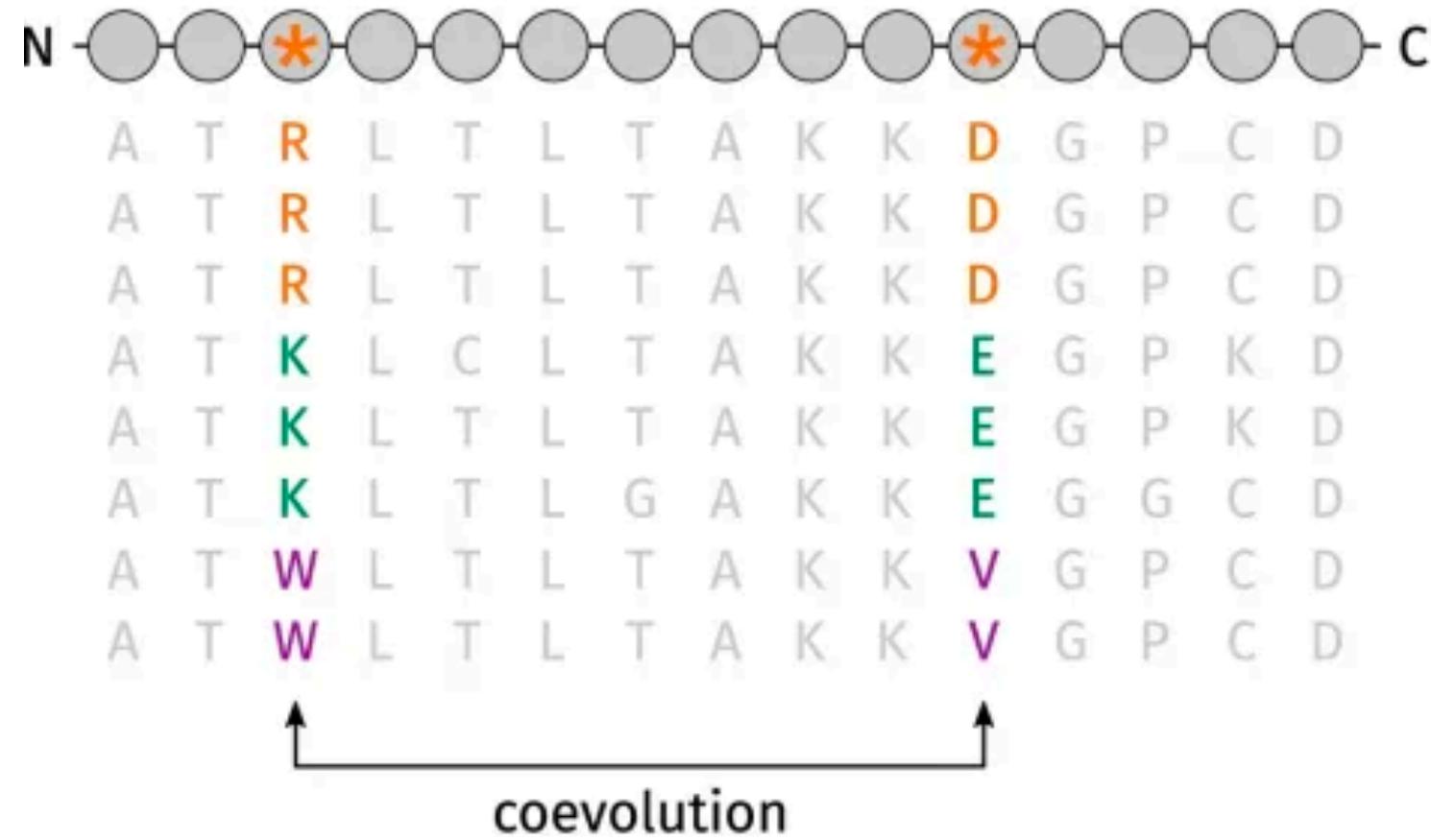
# Multiple sequence alignments encode structural information

Alignments contain information about contacts

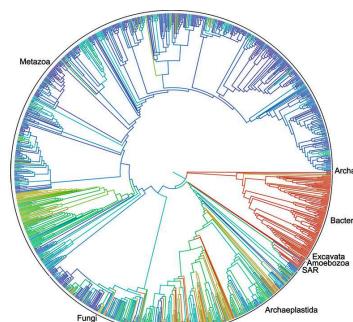


# Multiple sequence alignments encode structural information

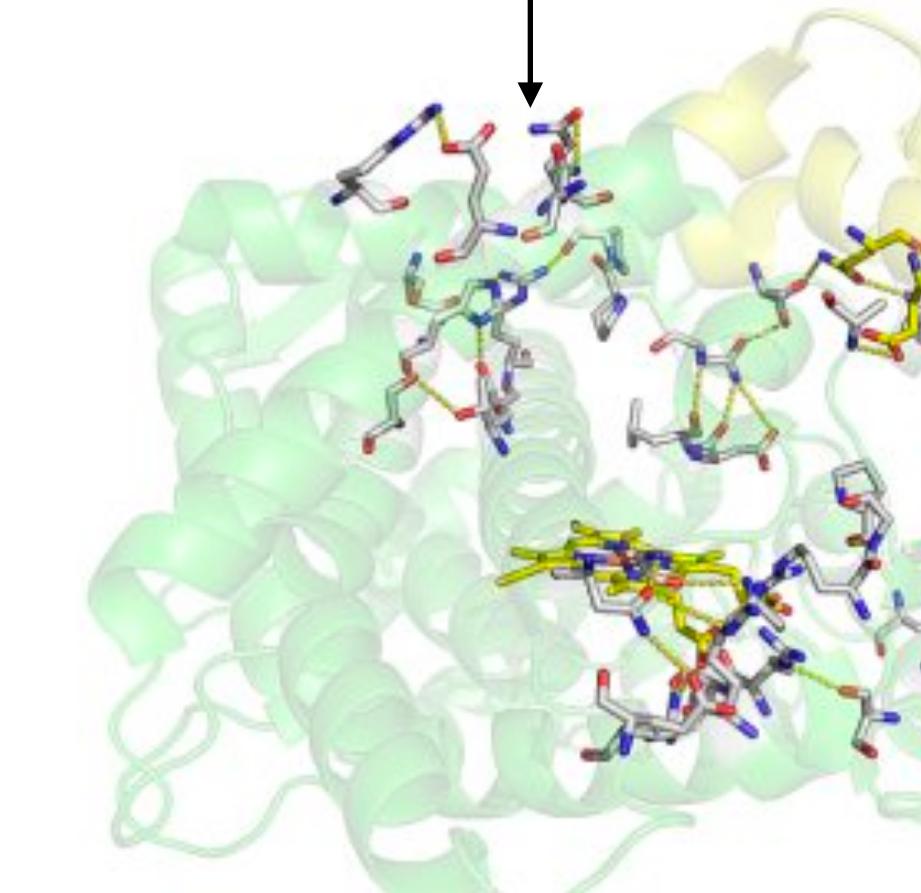
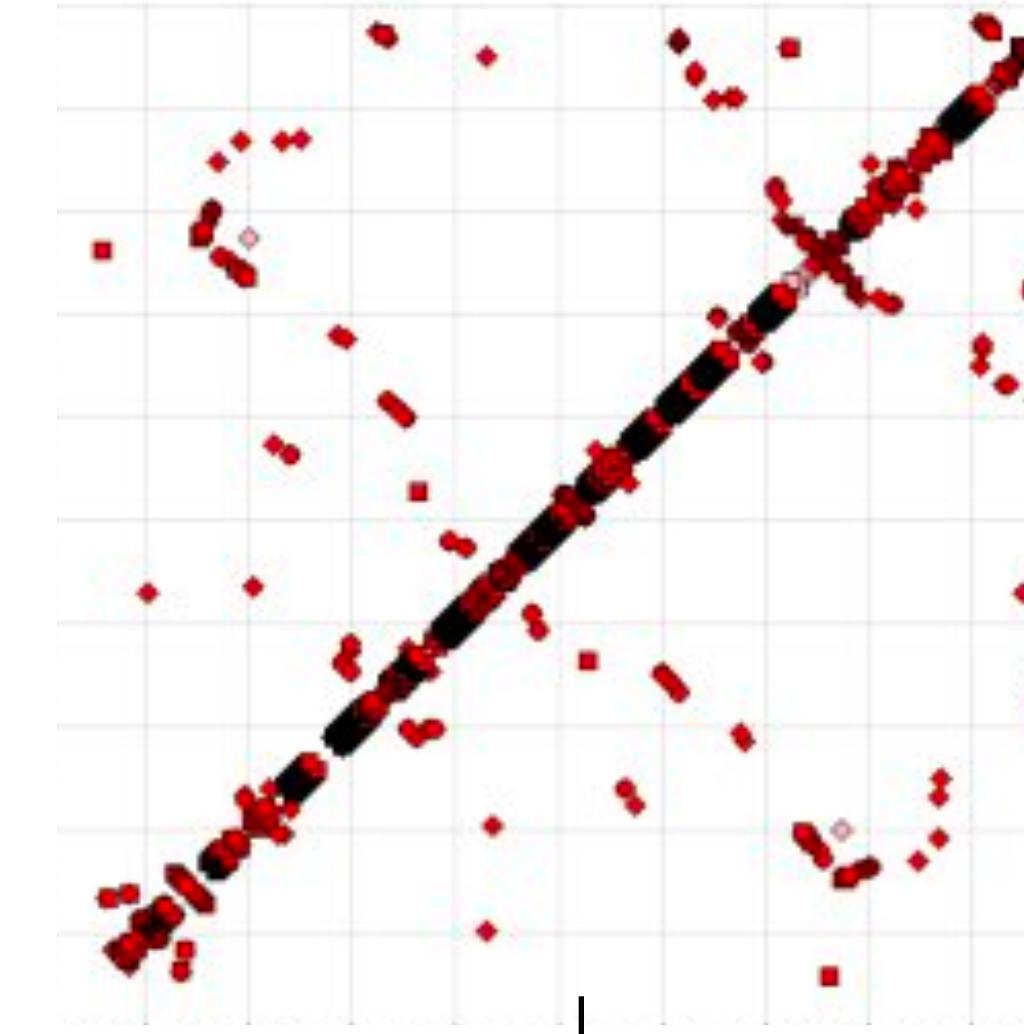
Alignments contain information  
about contacts



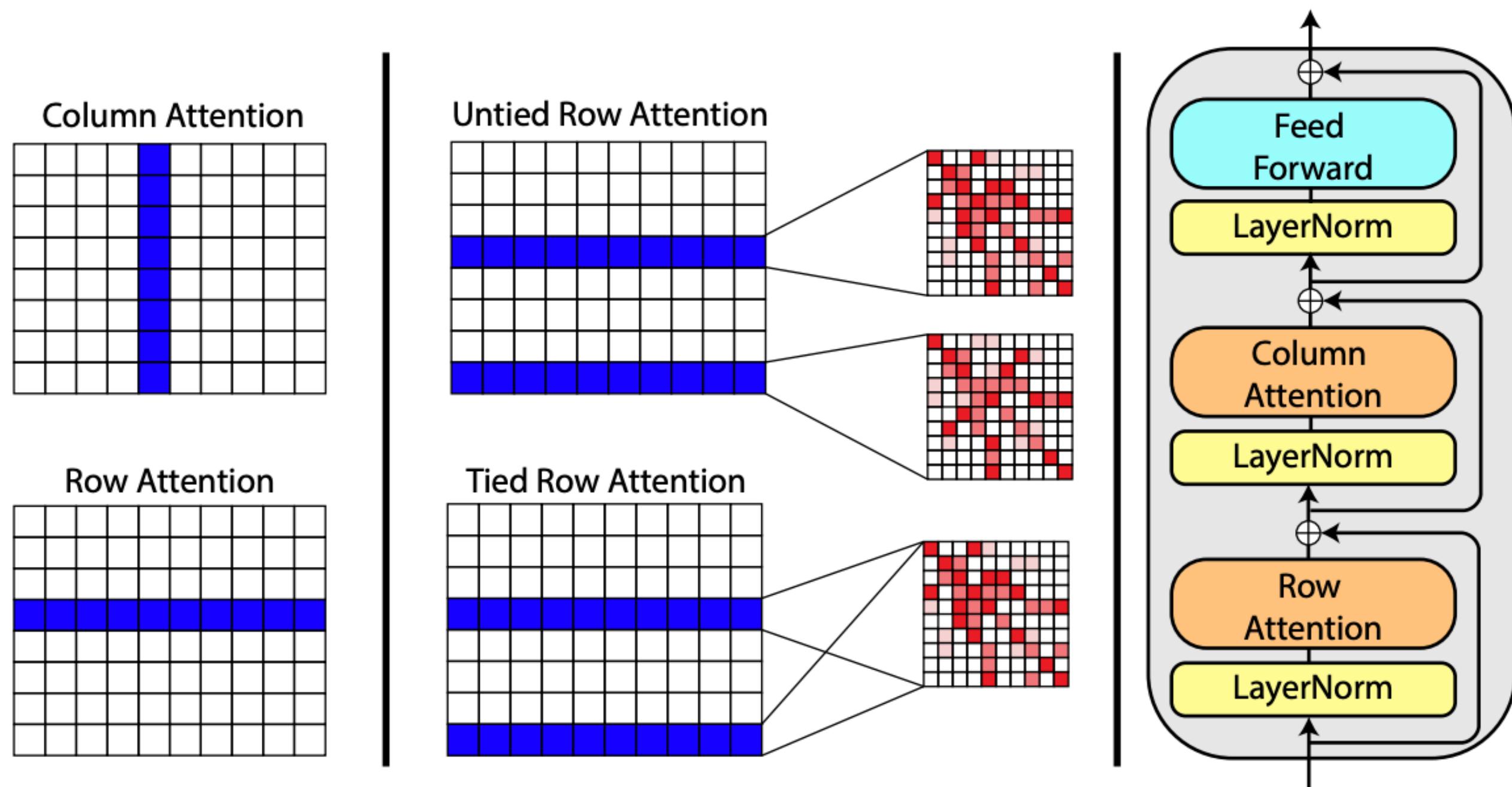
# evolutionary information



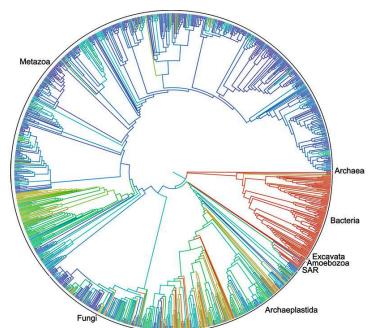
## Contacts constrain fold



# MSA-Transformer extracts information from MSAs

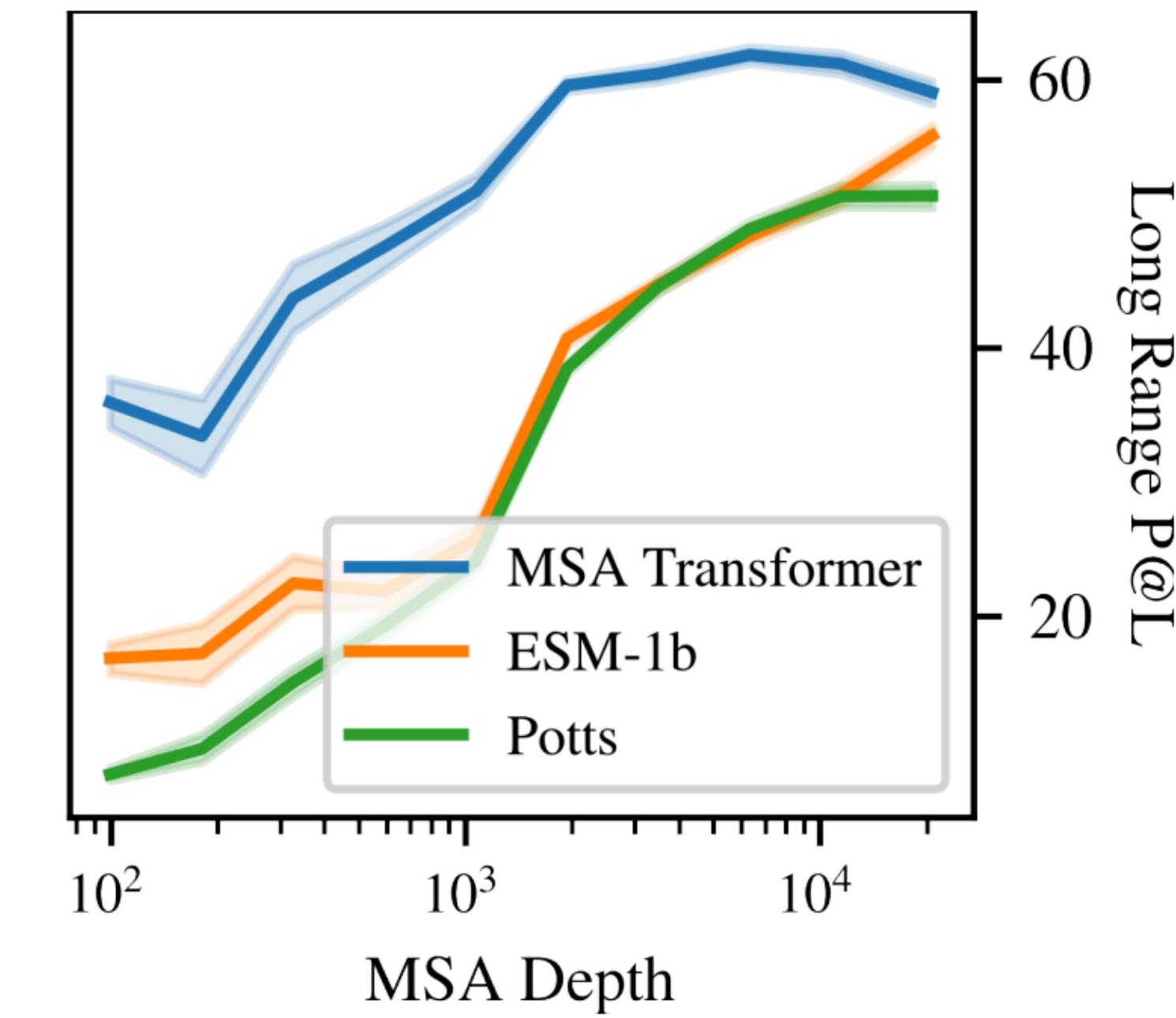
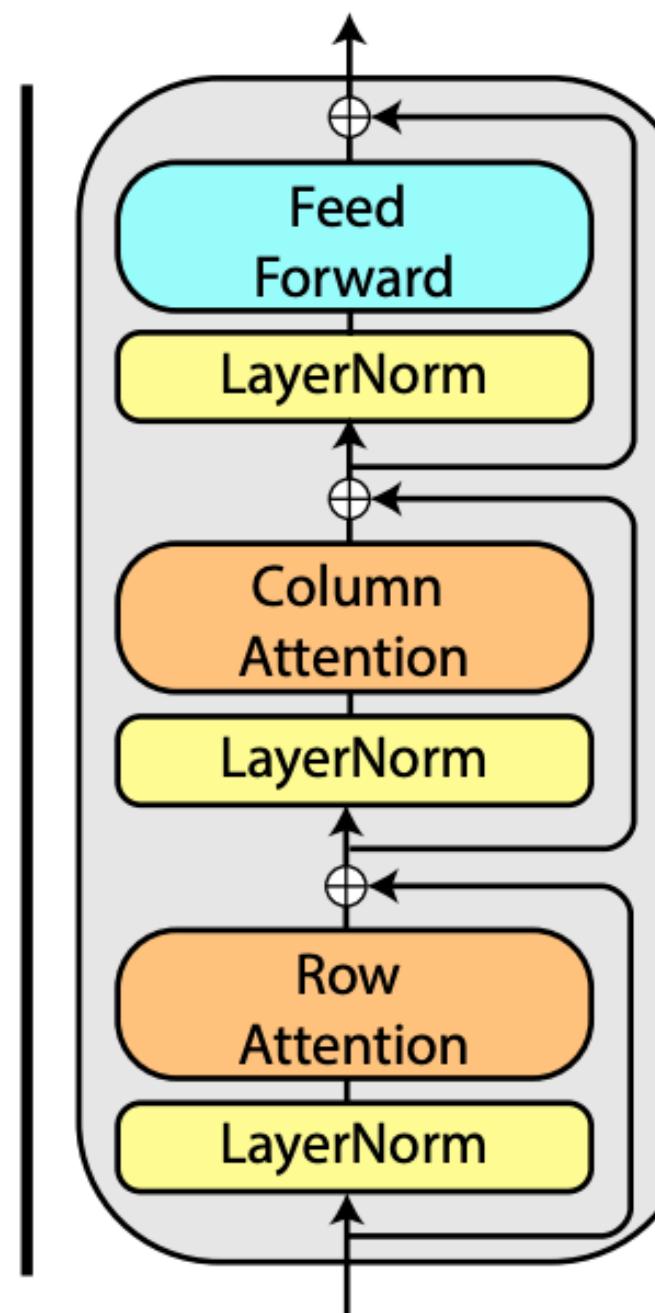
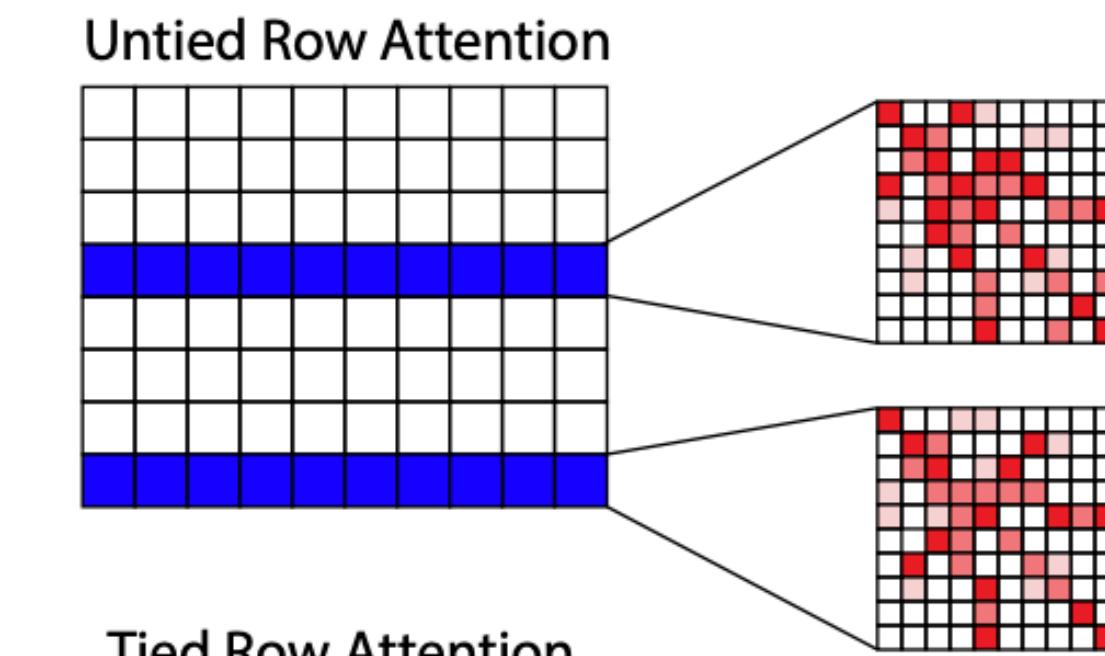
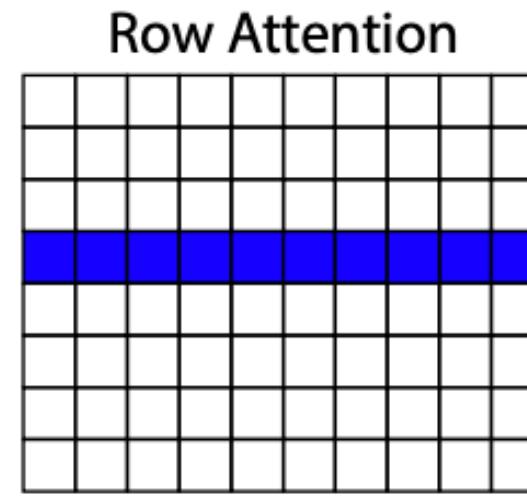
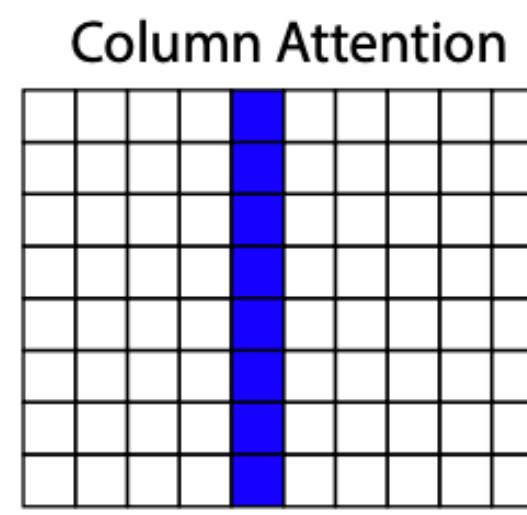


evolutionary  
information

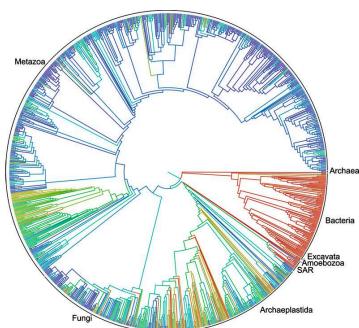


Rao et al. 2021

# MSA-Transformer extracts information from MSAs

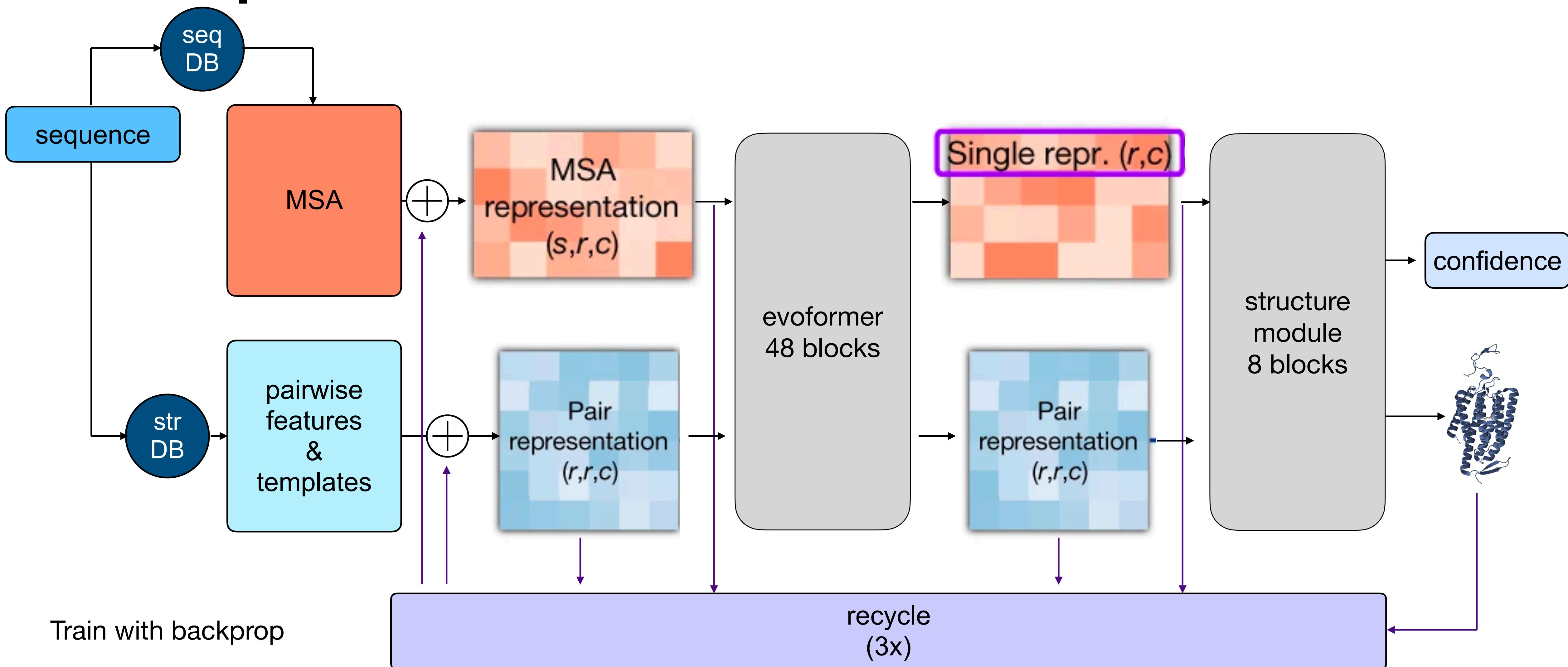


evolutionary  
information



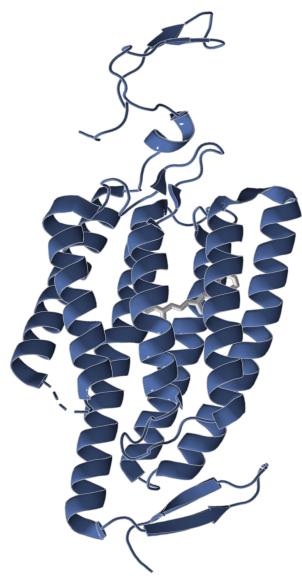
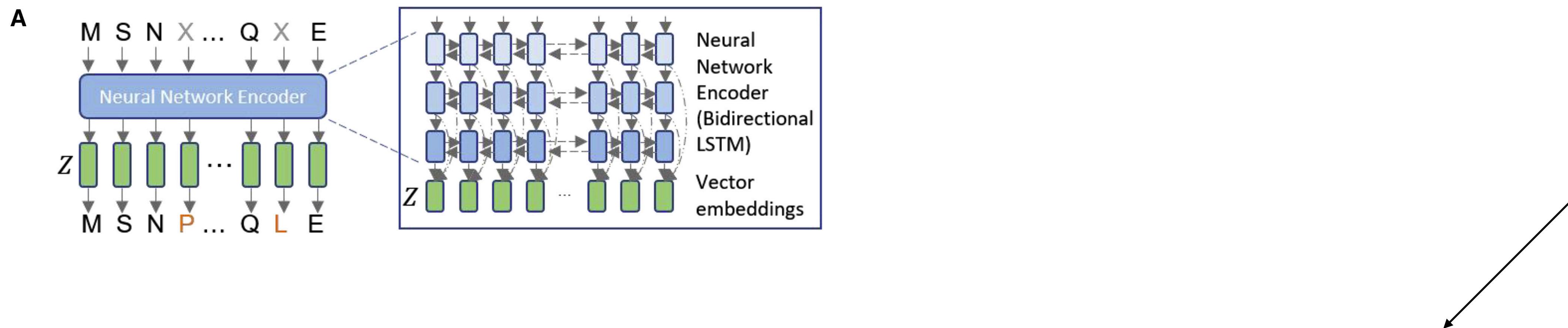
Rao et al. 2021

# AlphaFold uses MSAs and transformers

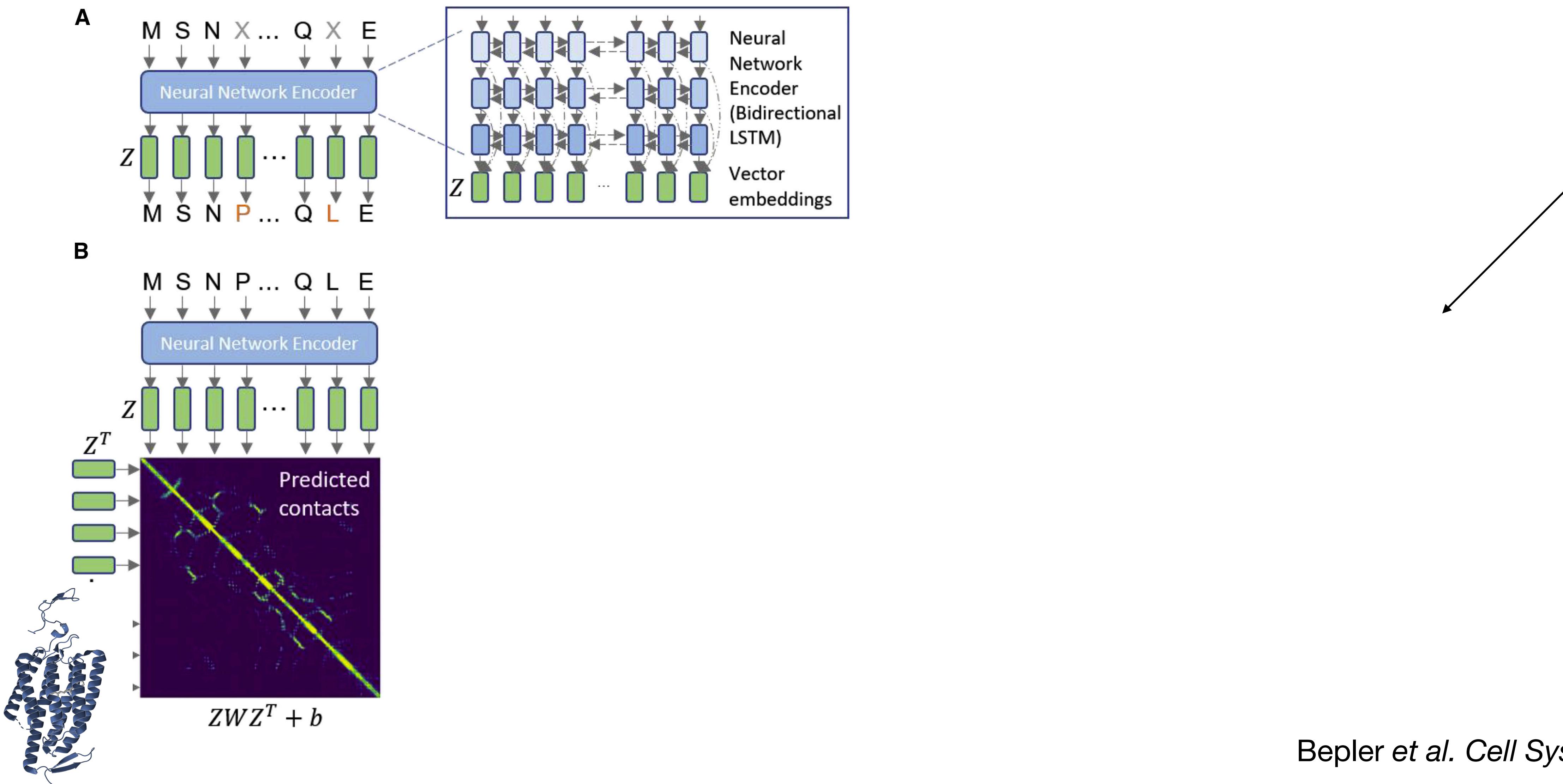


16 TPUv3s for several weeks

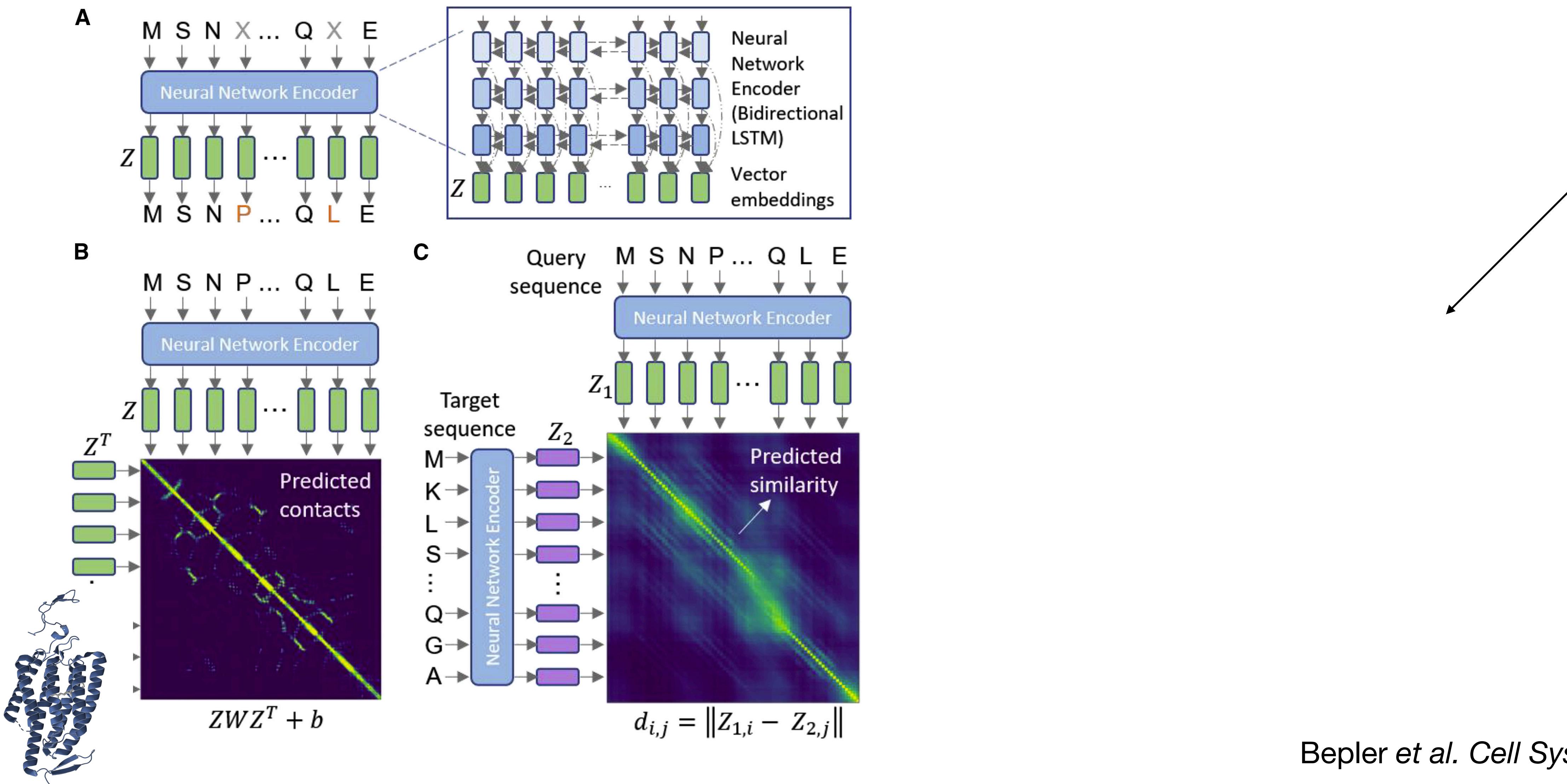
# Can use structure as a pretext task



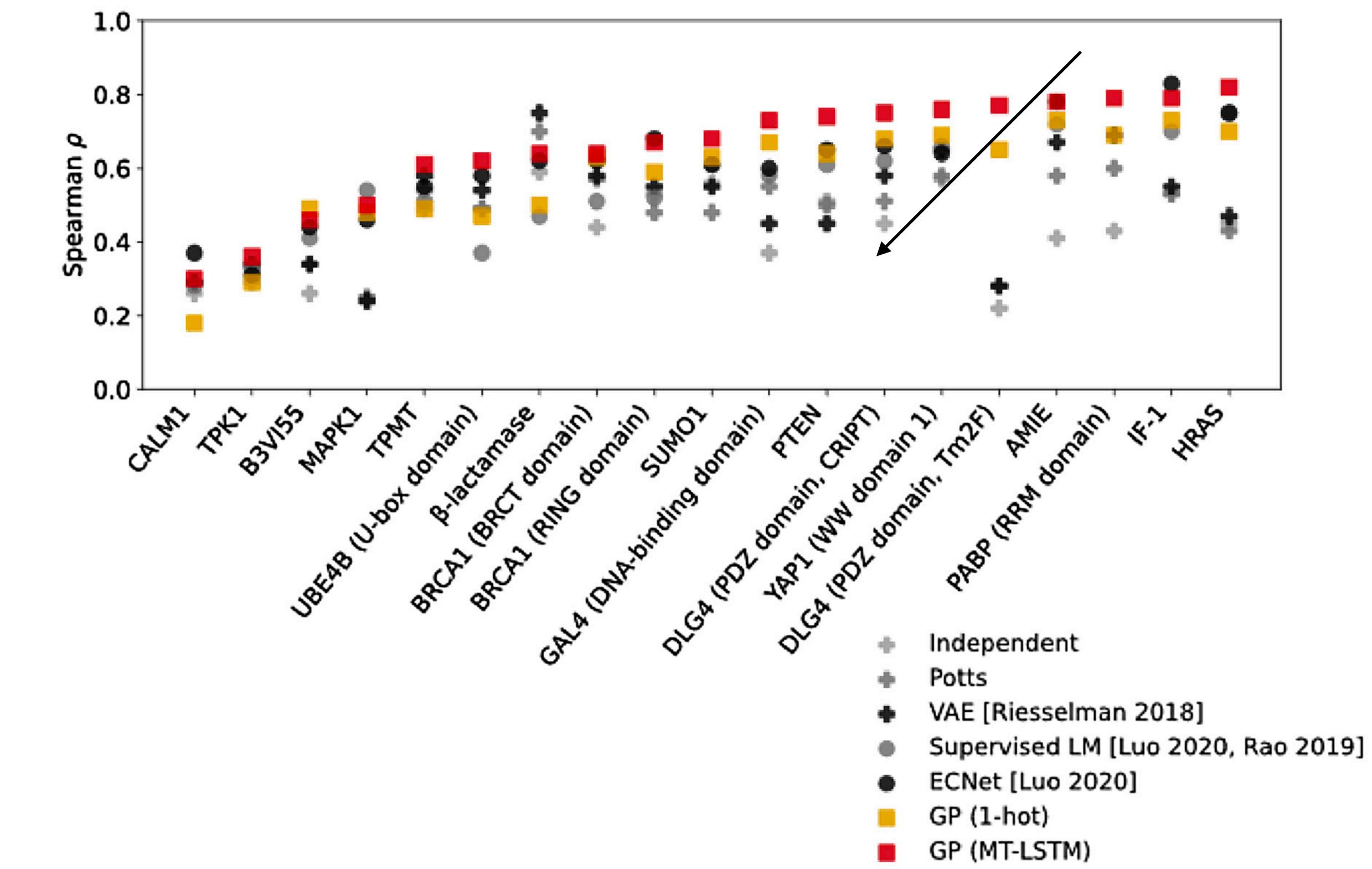
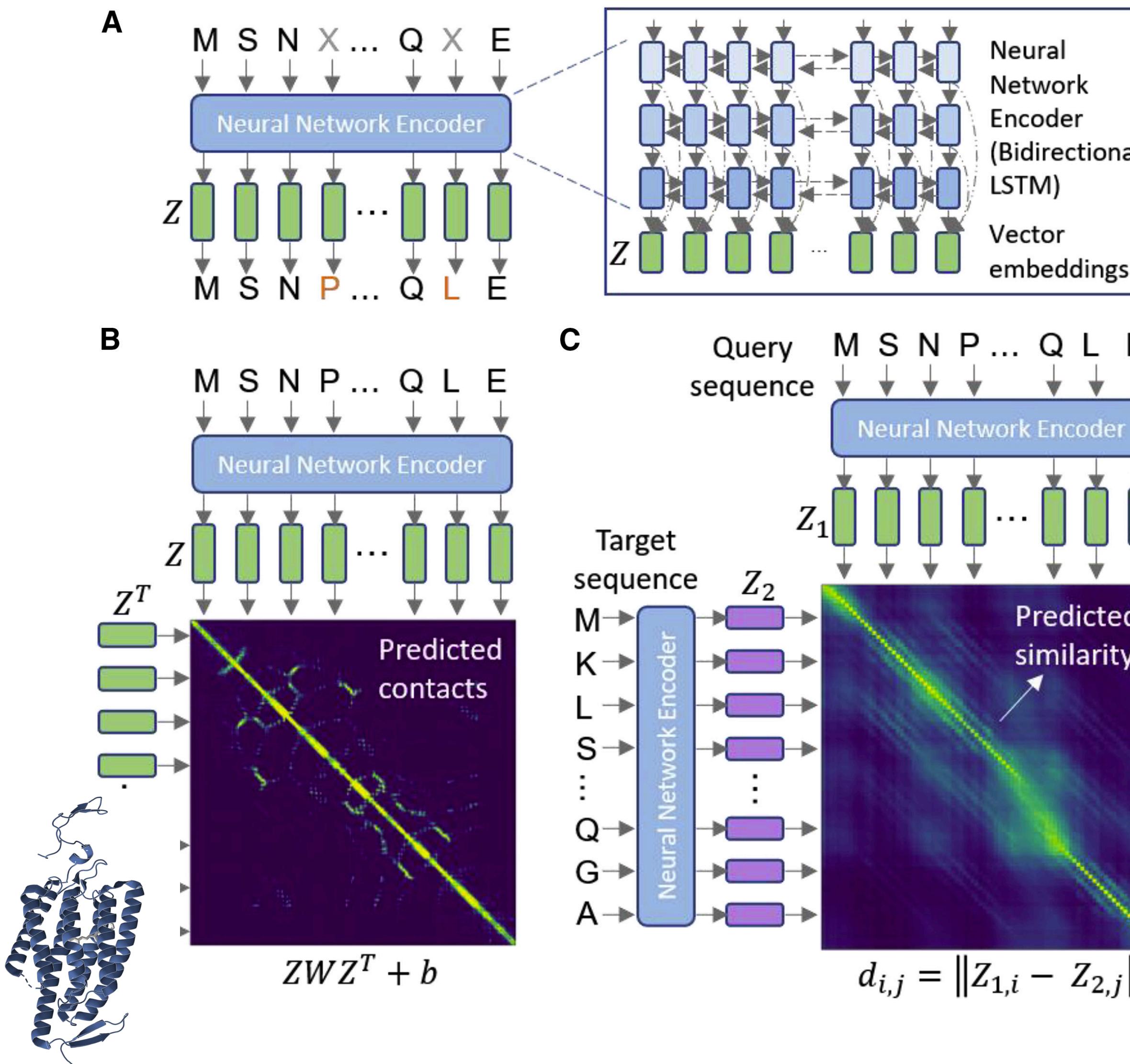
# Can use structure as a pretext task



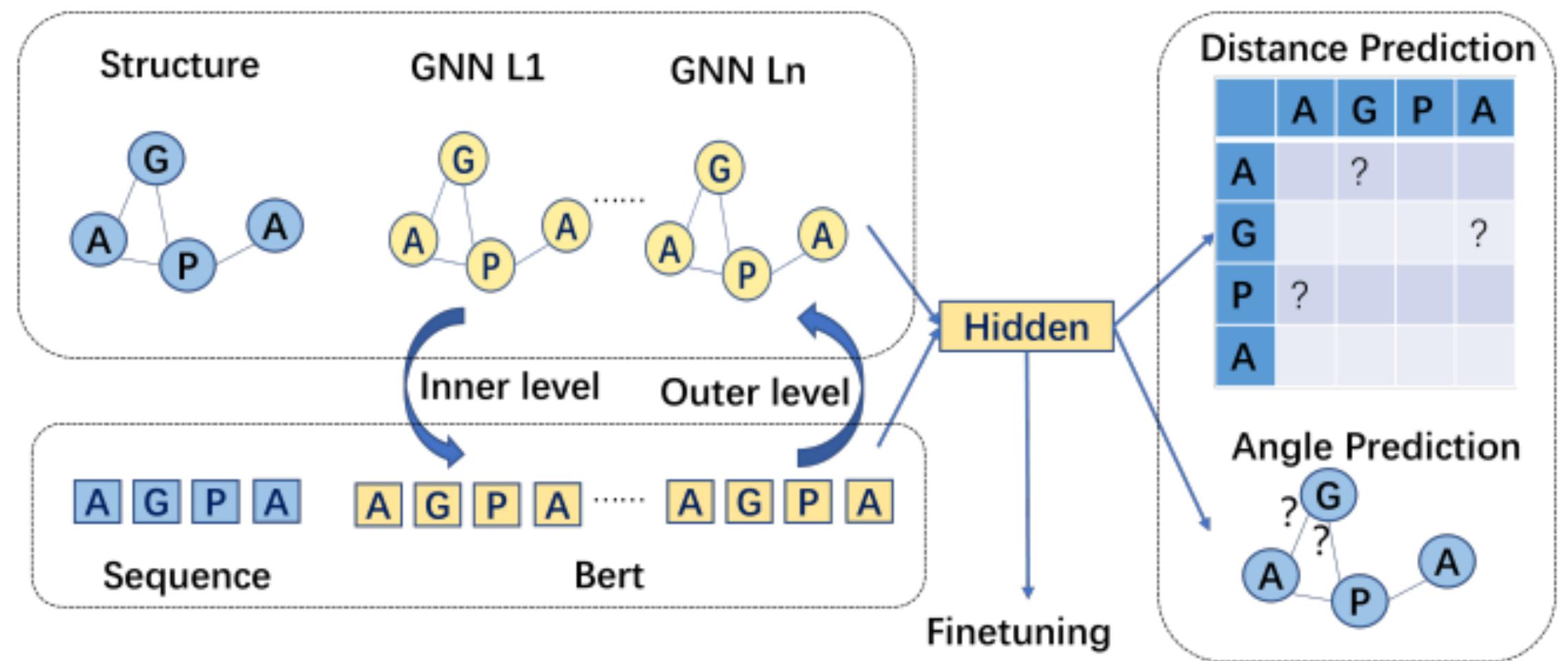
# Can use structure as a pretext task



# Can use structure as a pretext task

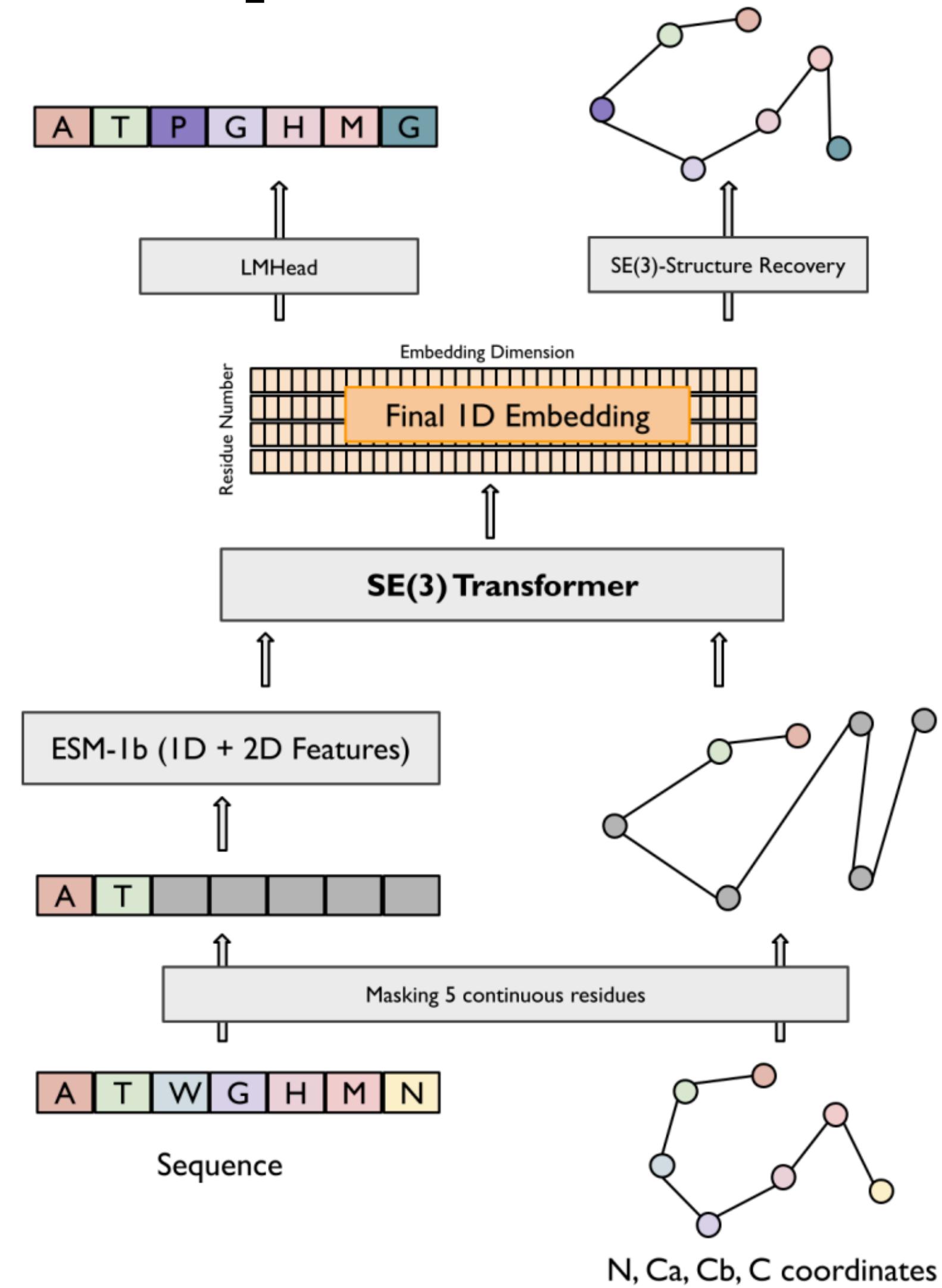
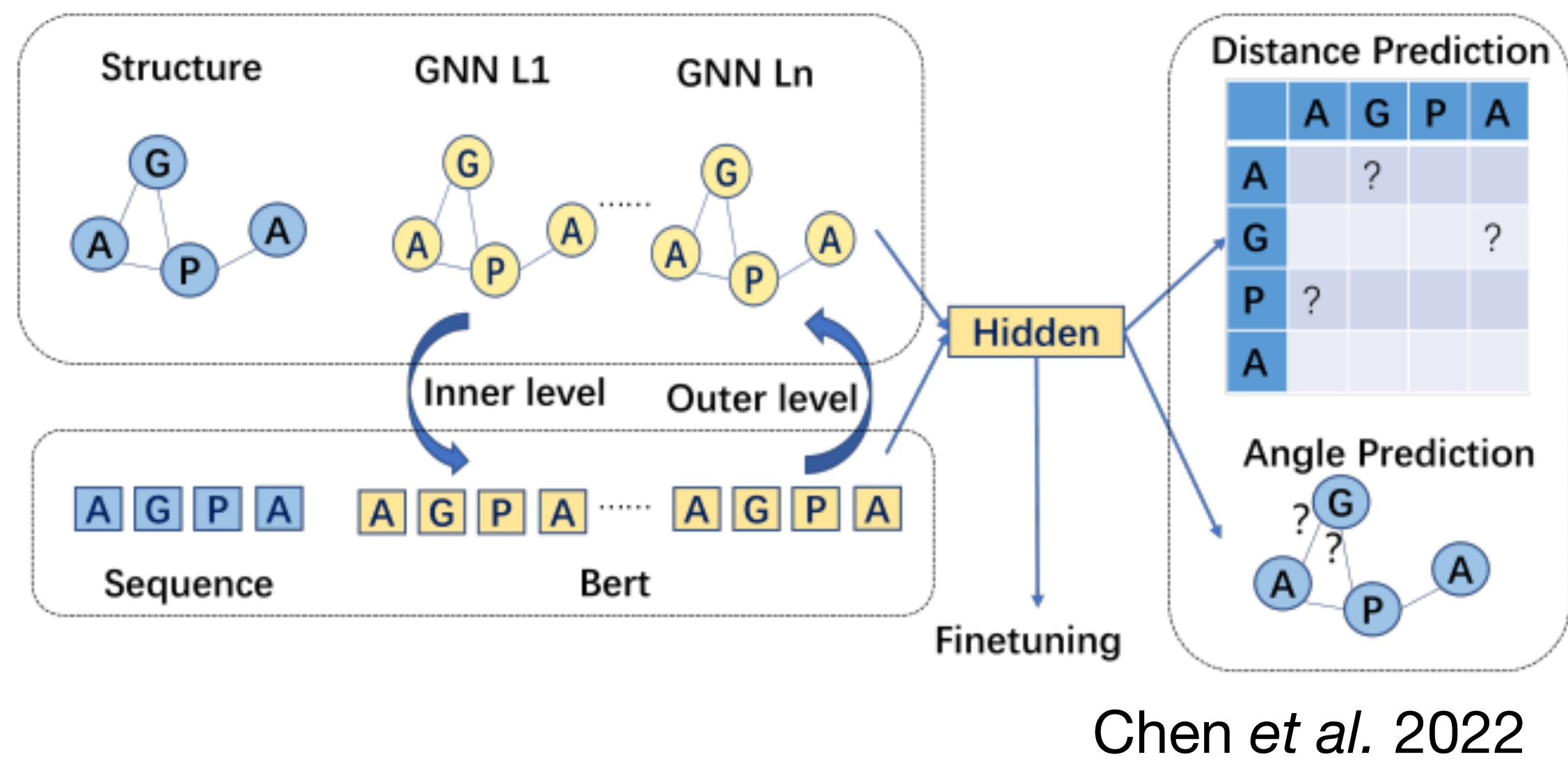


# Can use structure as a pretext task



Chen et al. 2022

# Can use structure as a pretext task



# Conditioning on structure improves sequence reconstruction

# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

# Conditioning on structure improves sequence reconstruction

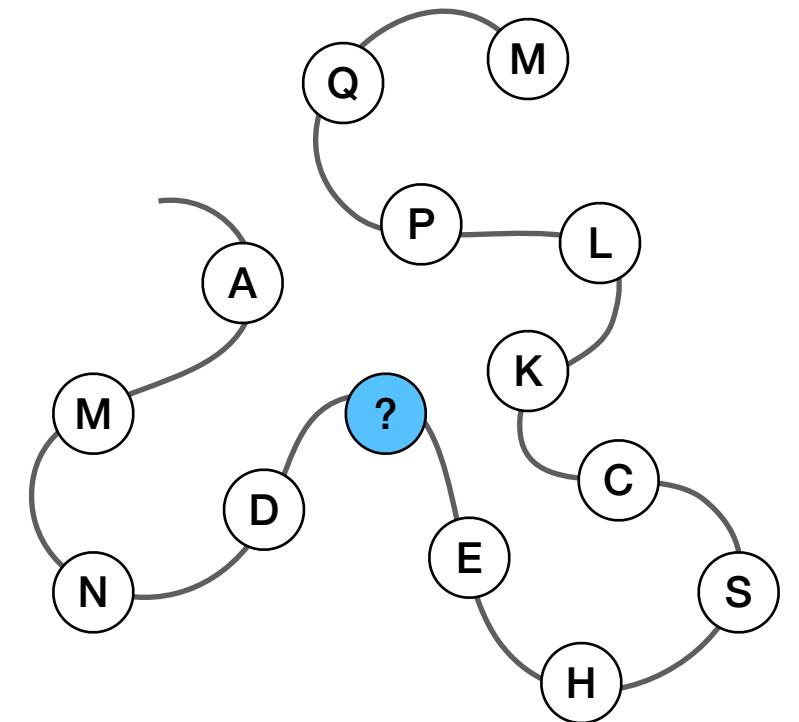
Structure-conditioned masked language model

Goal: Predict  $P(s_{masked} | \text{structure}, s_{unmasked})$

# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

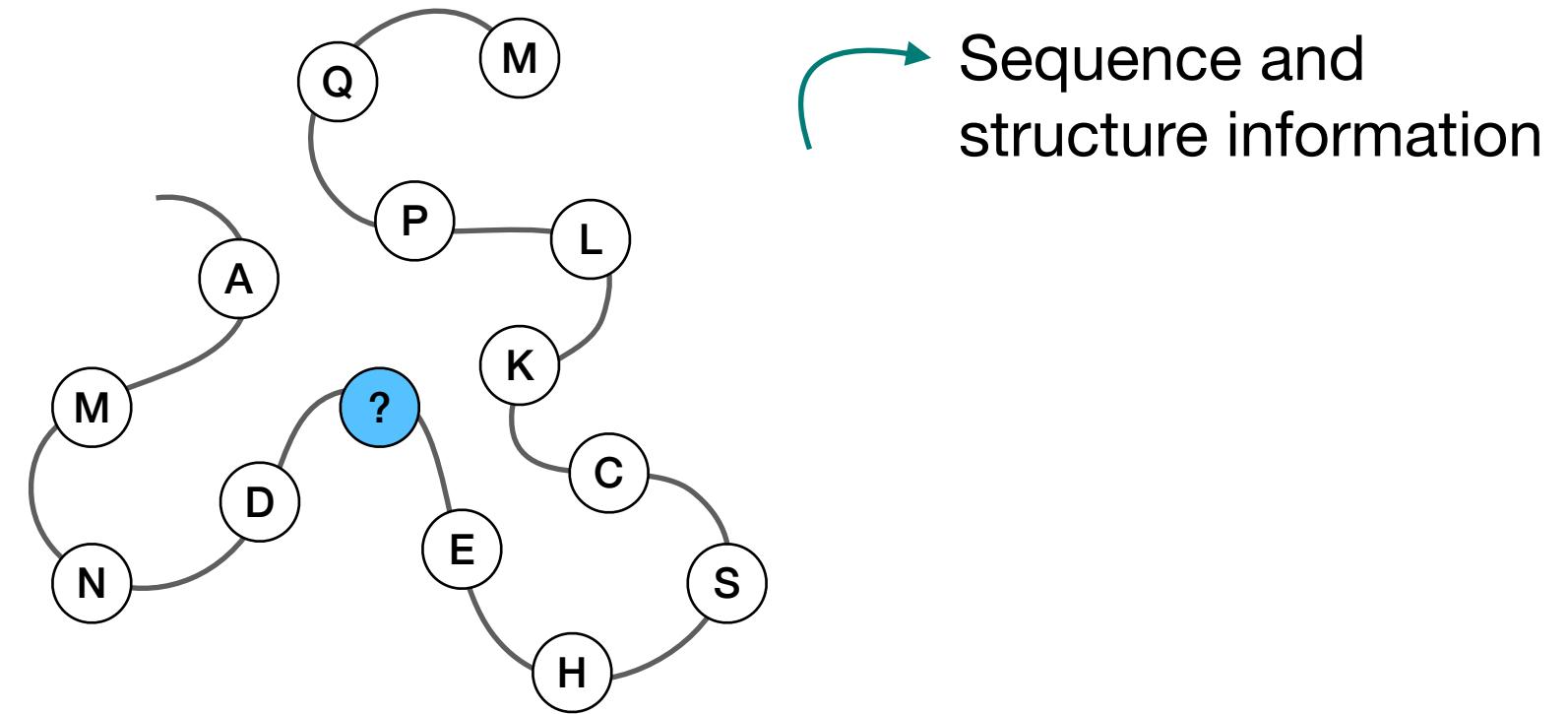
Goal: Predict  $P(s_{masked} | \text{structure}, s_{unmasked})$



# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

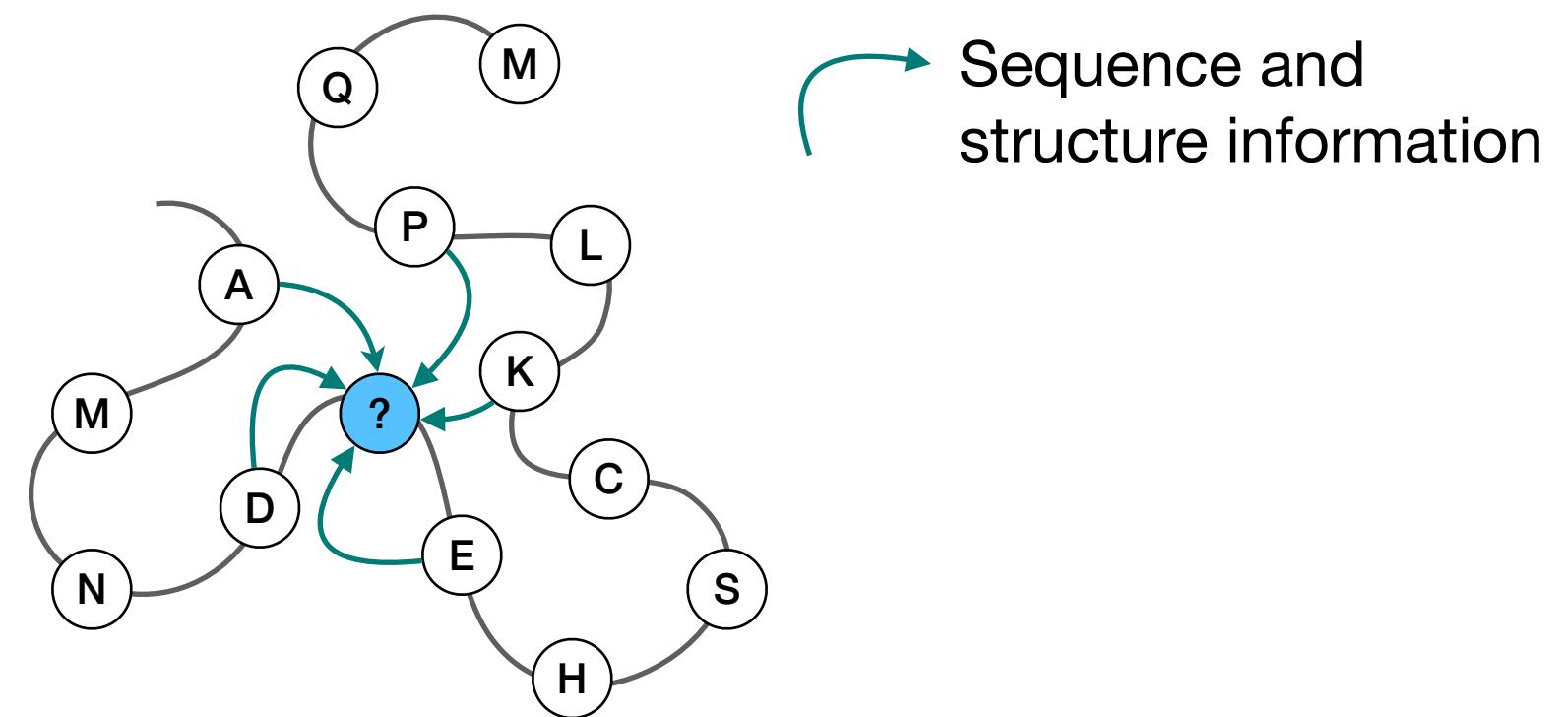
Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$



# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

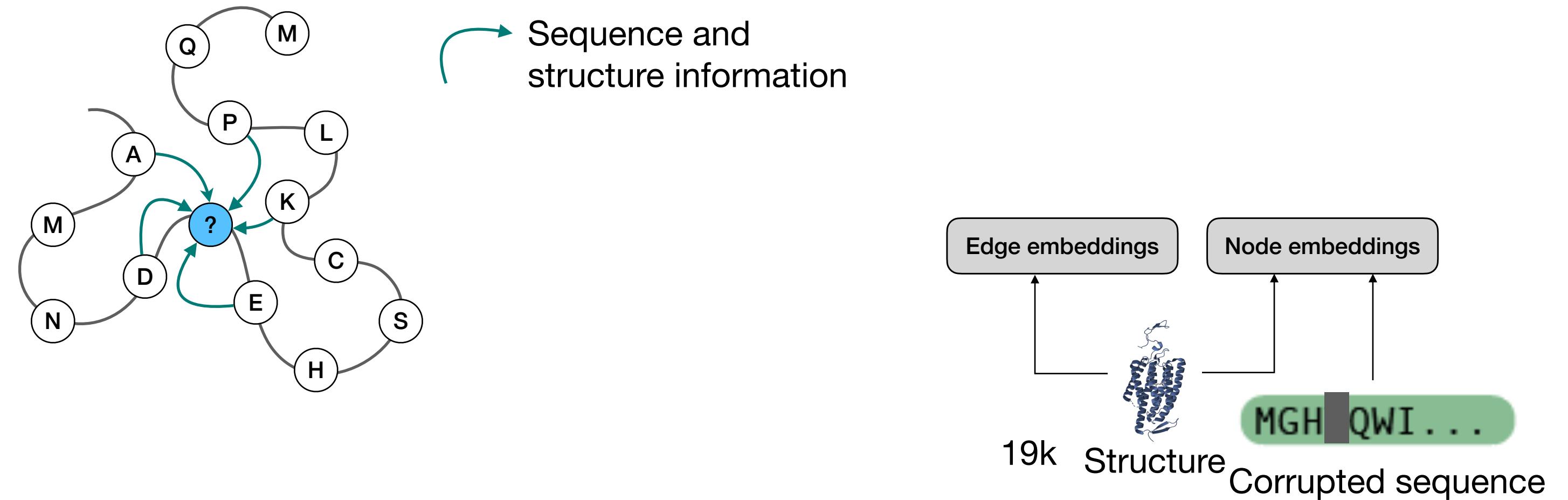
Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$



# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

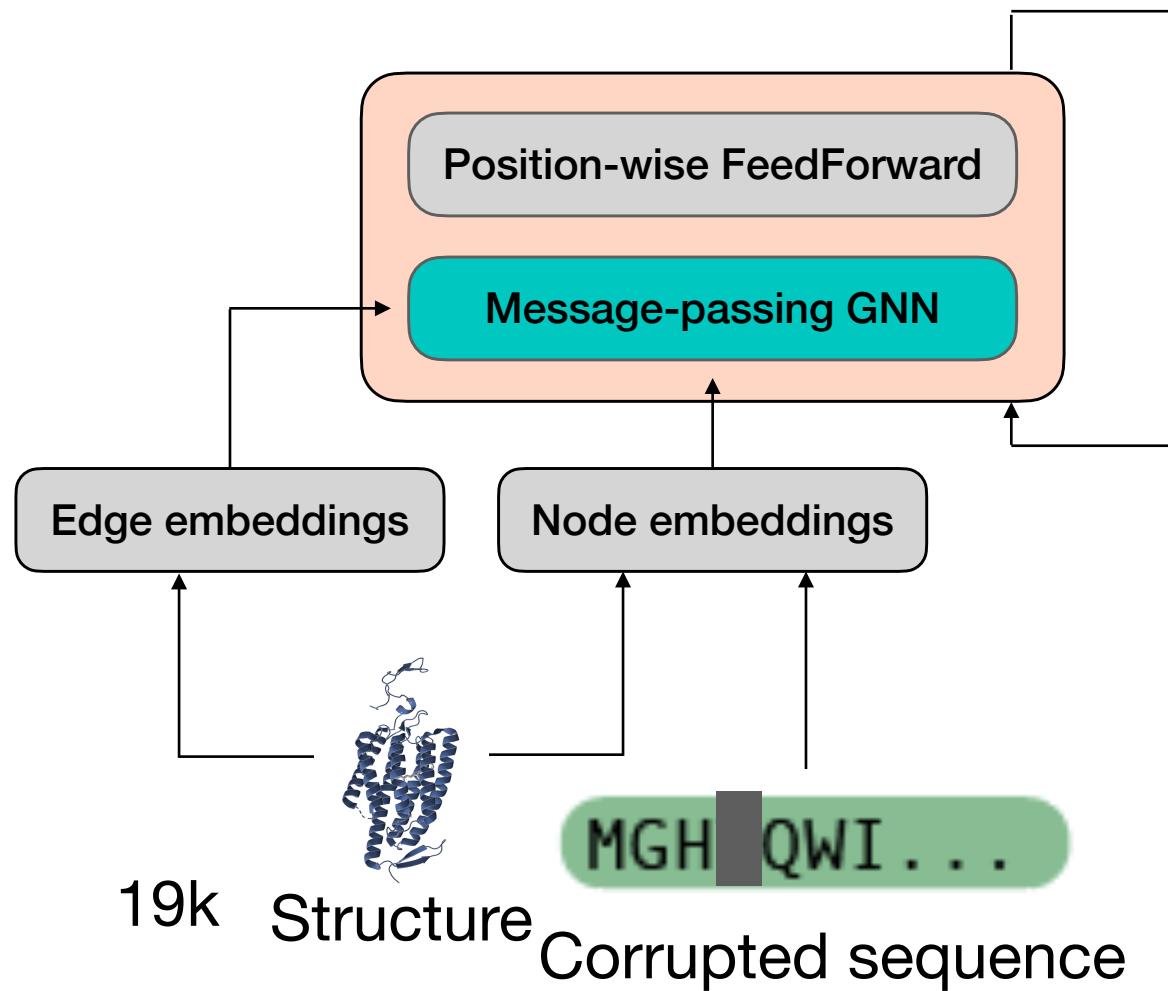
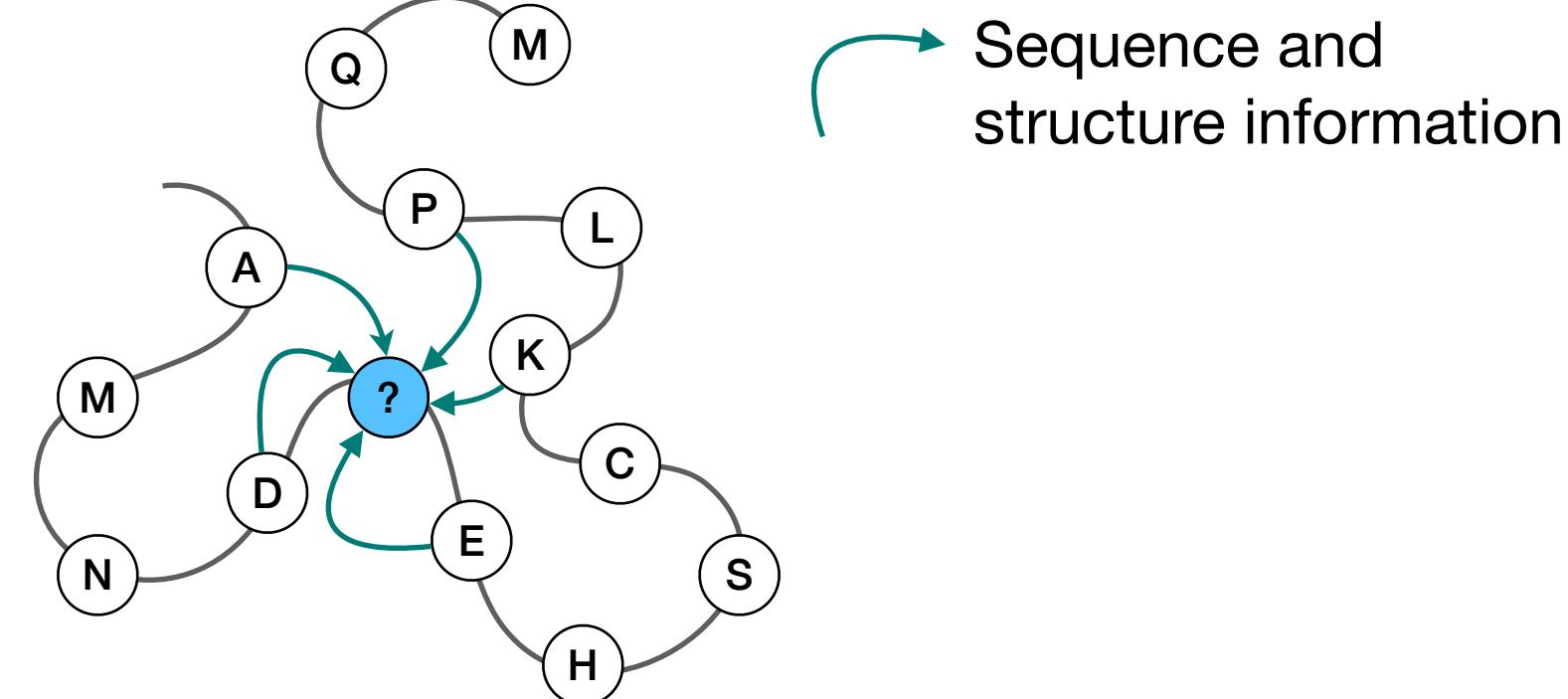
Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$



# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

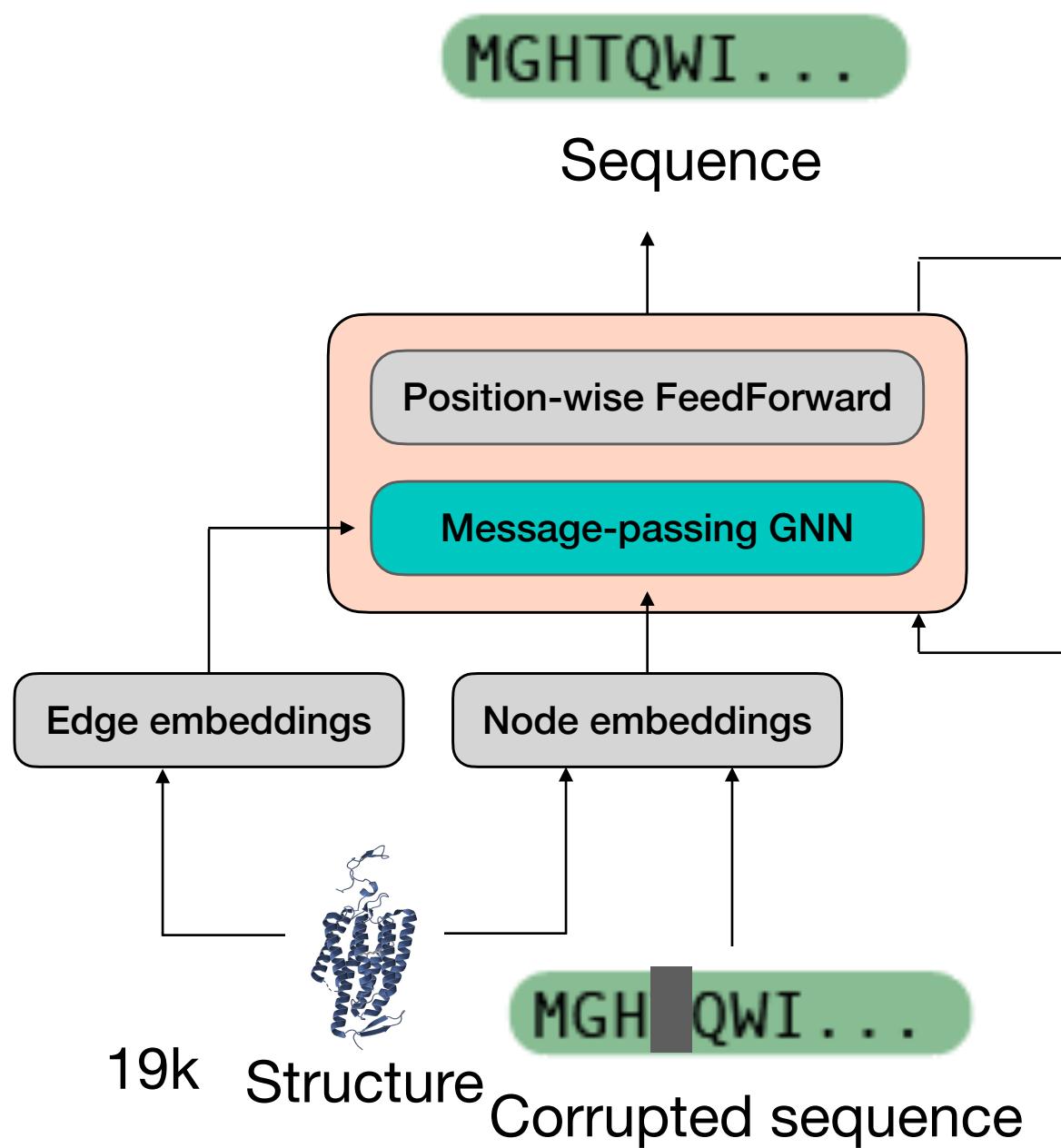
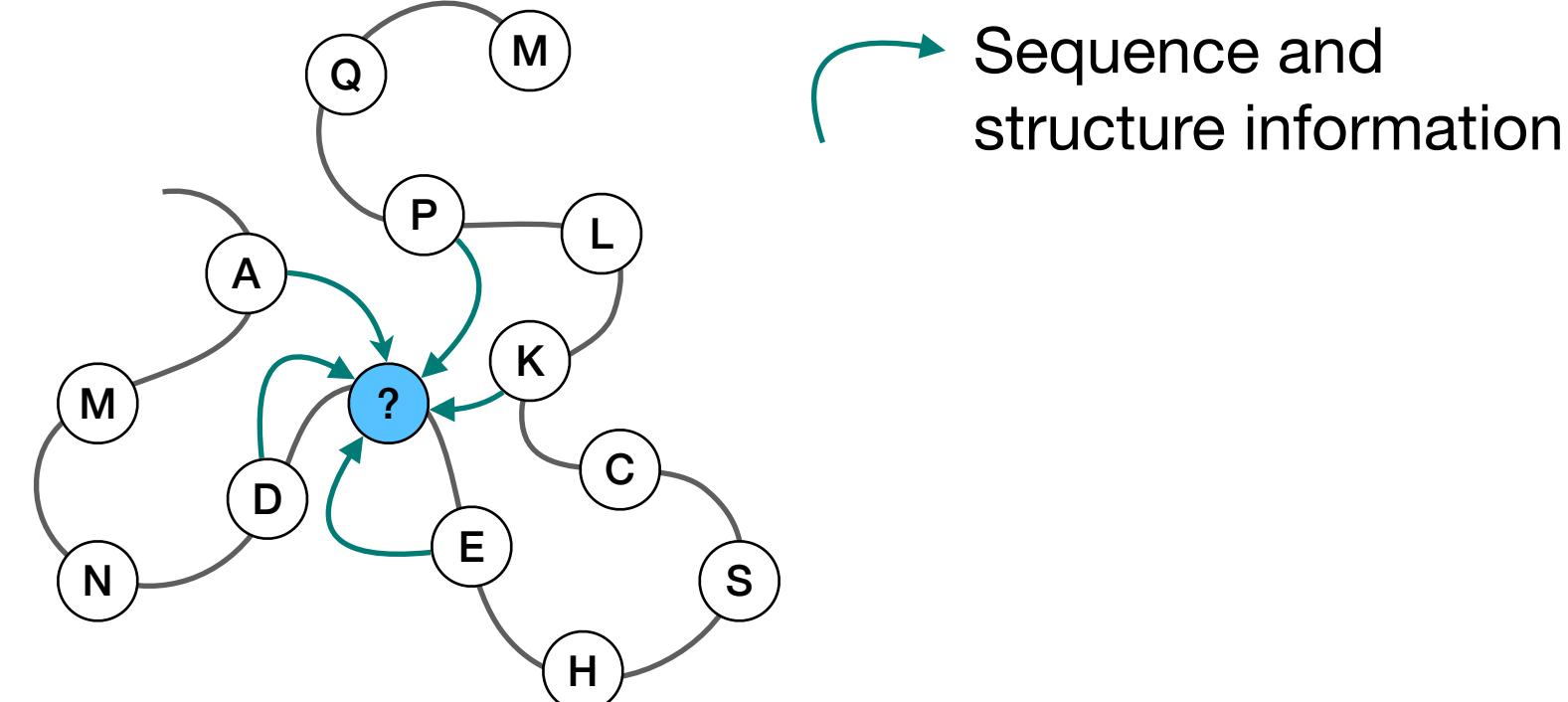
Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$



# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

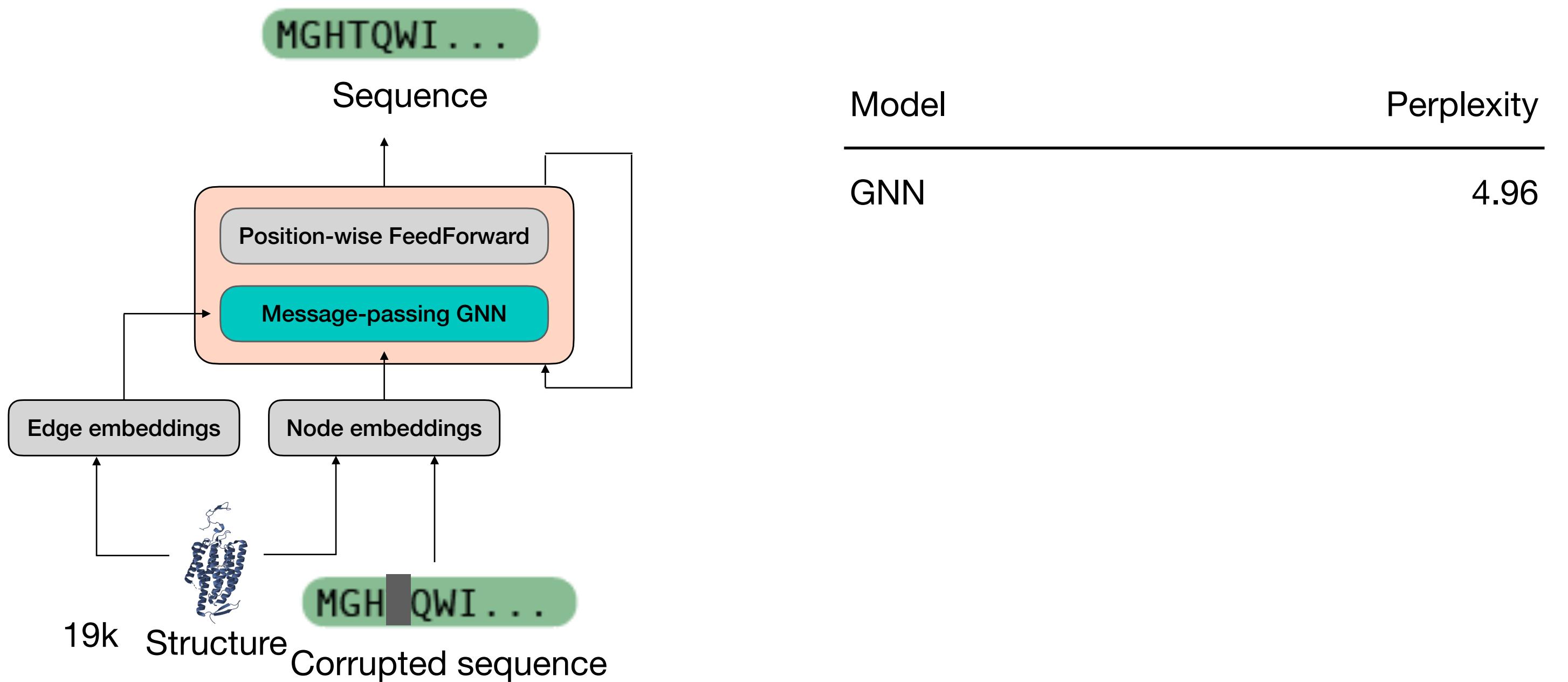
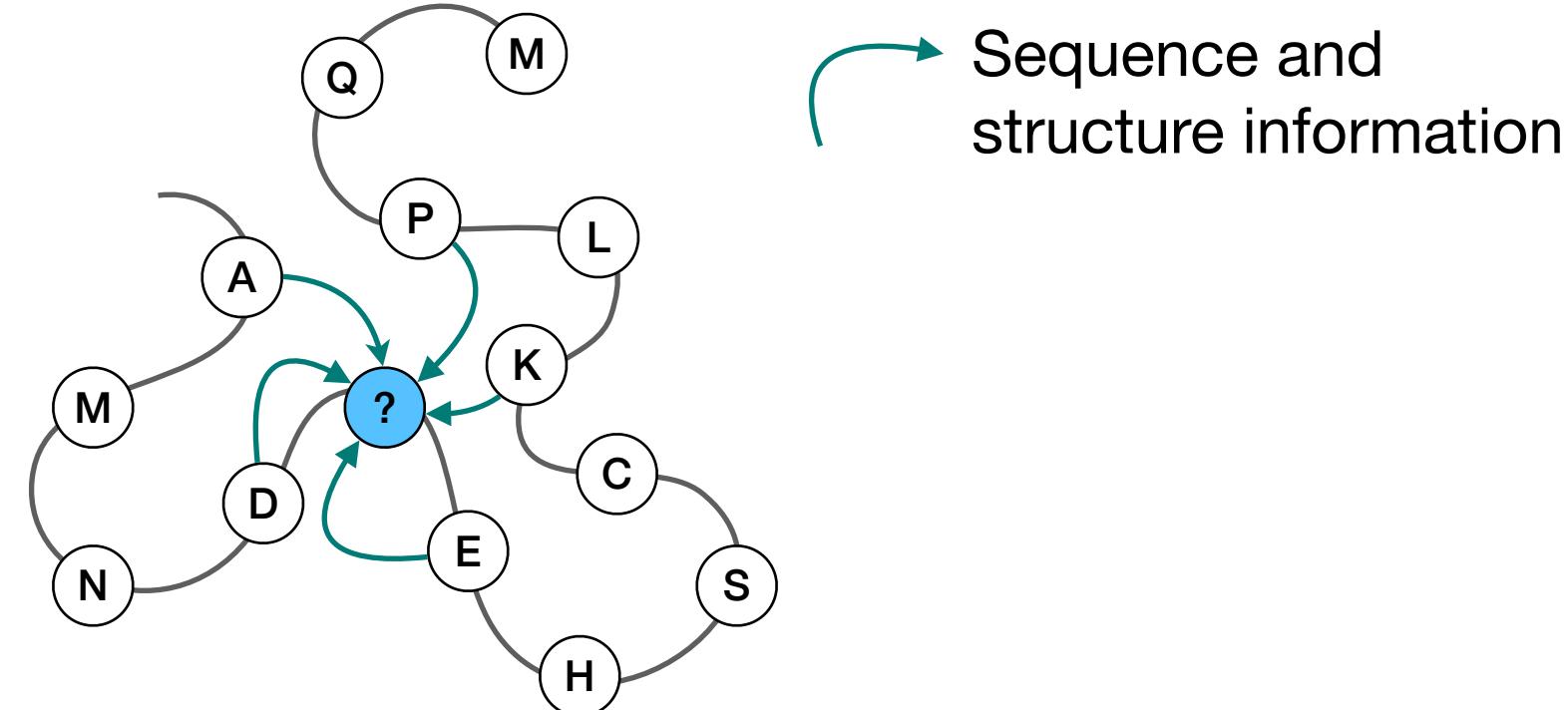
Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$



# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

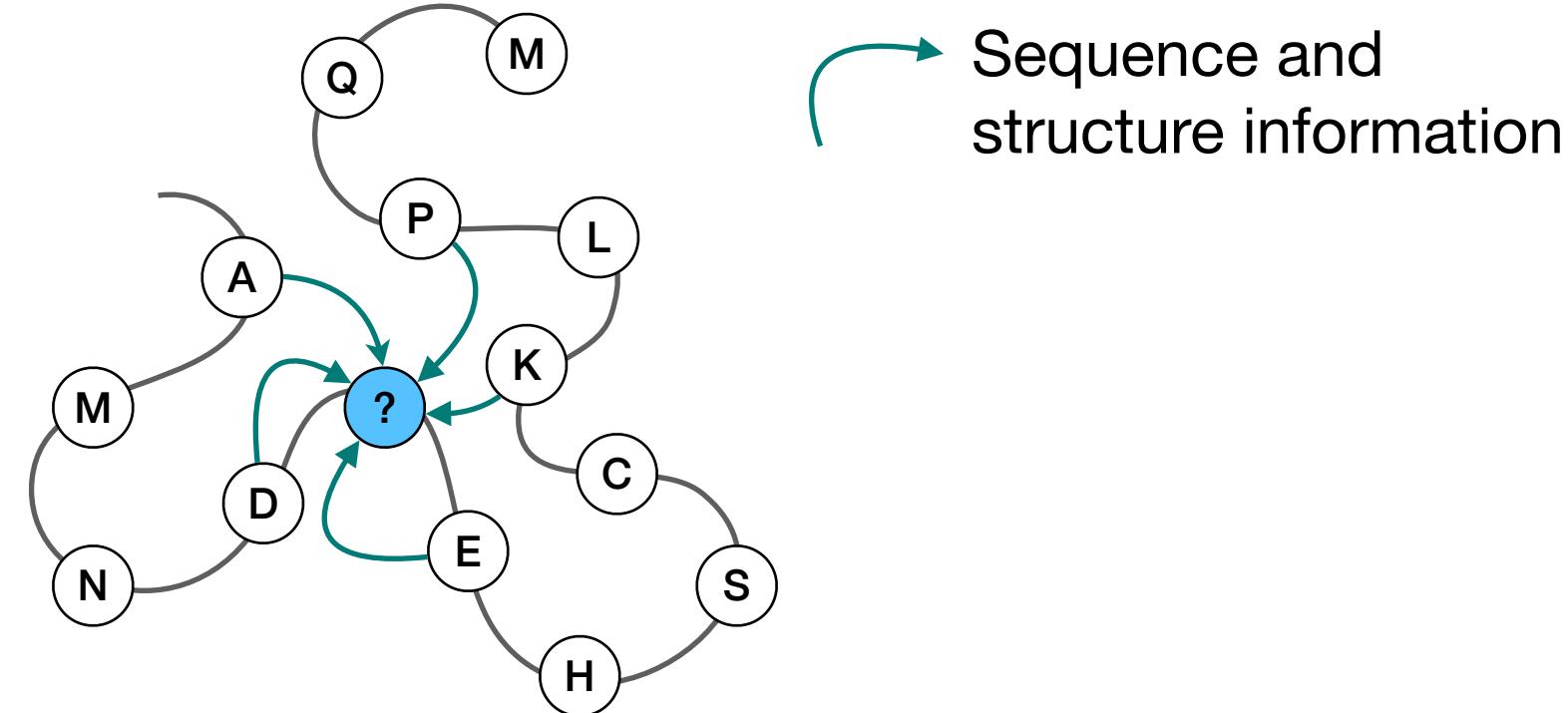
Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$



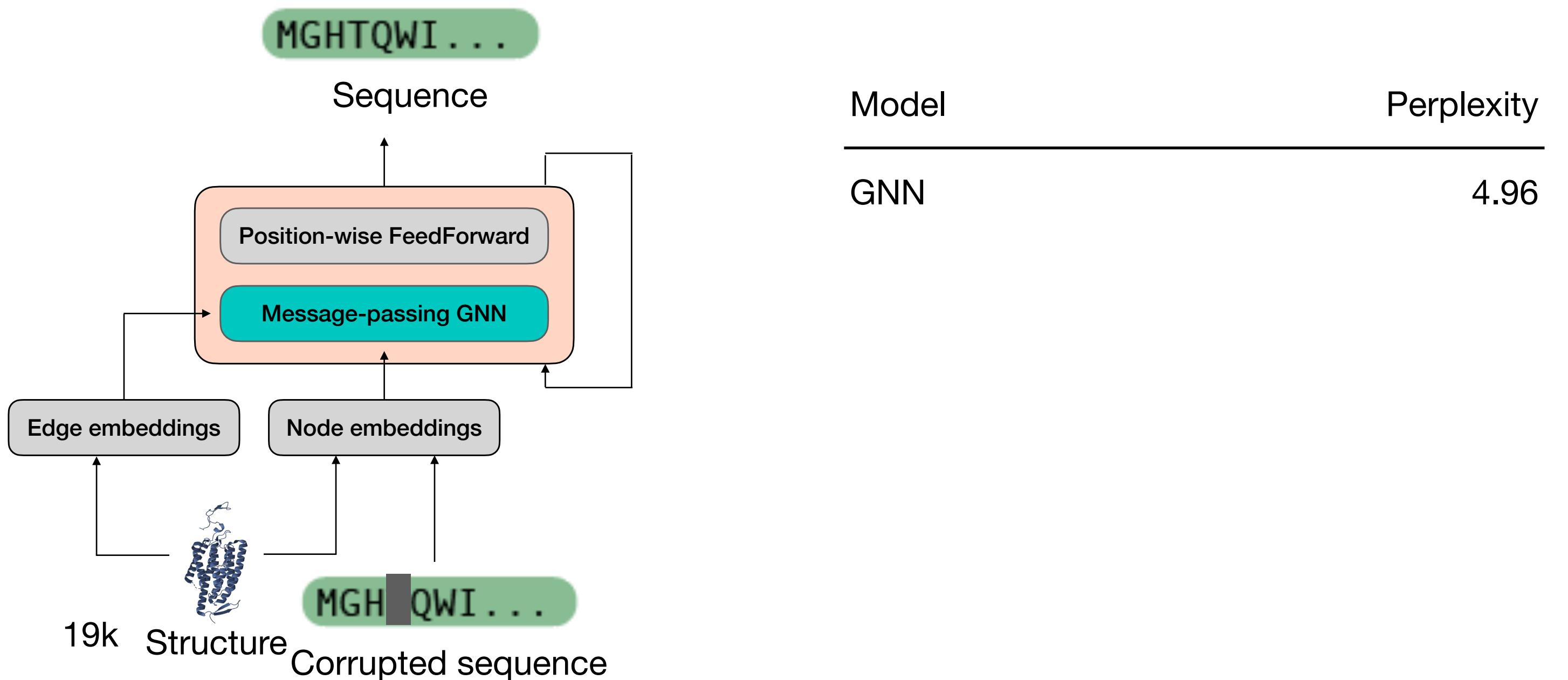
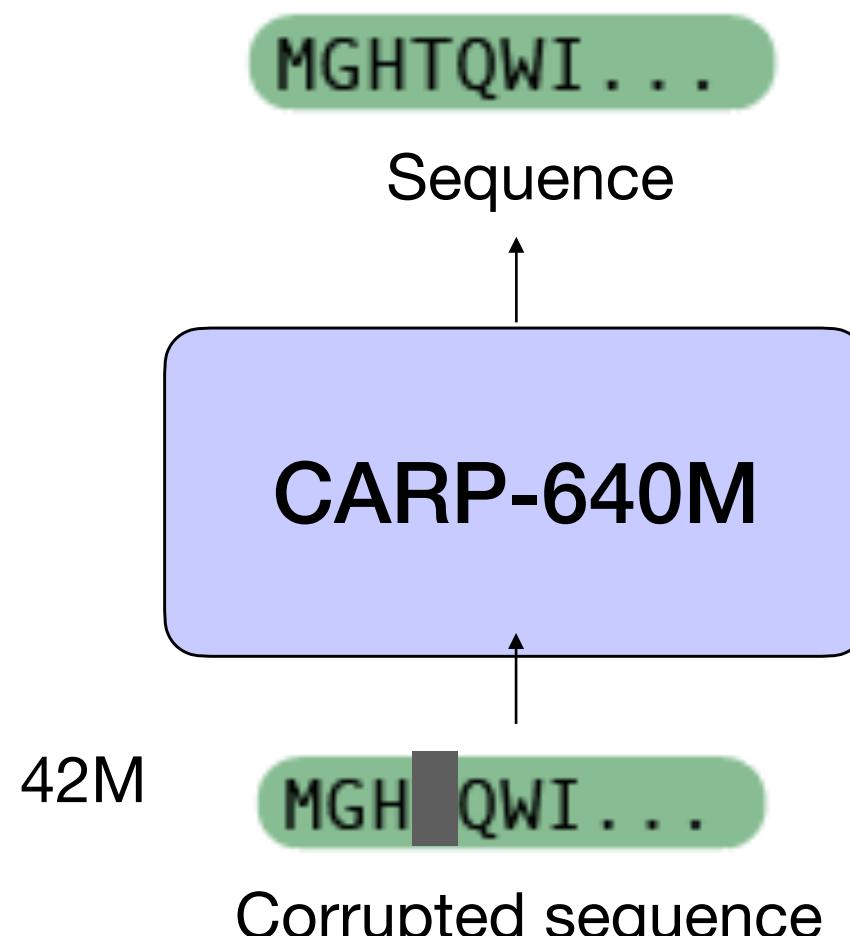
# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$



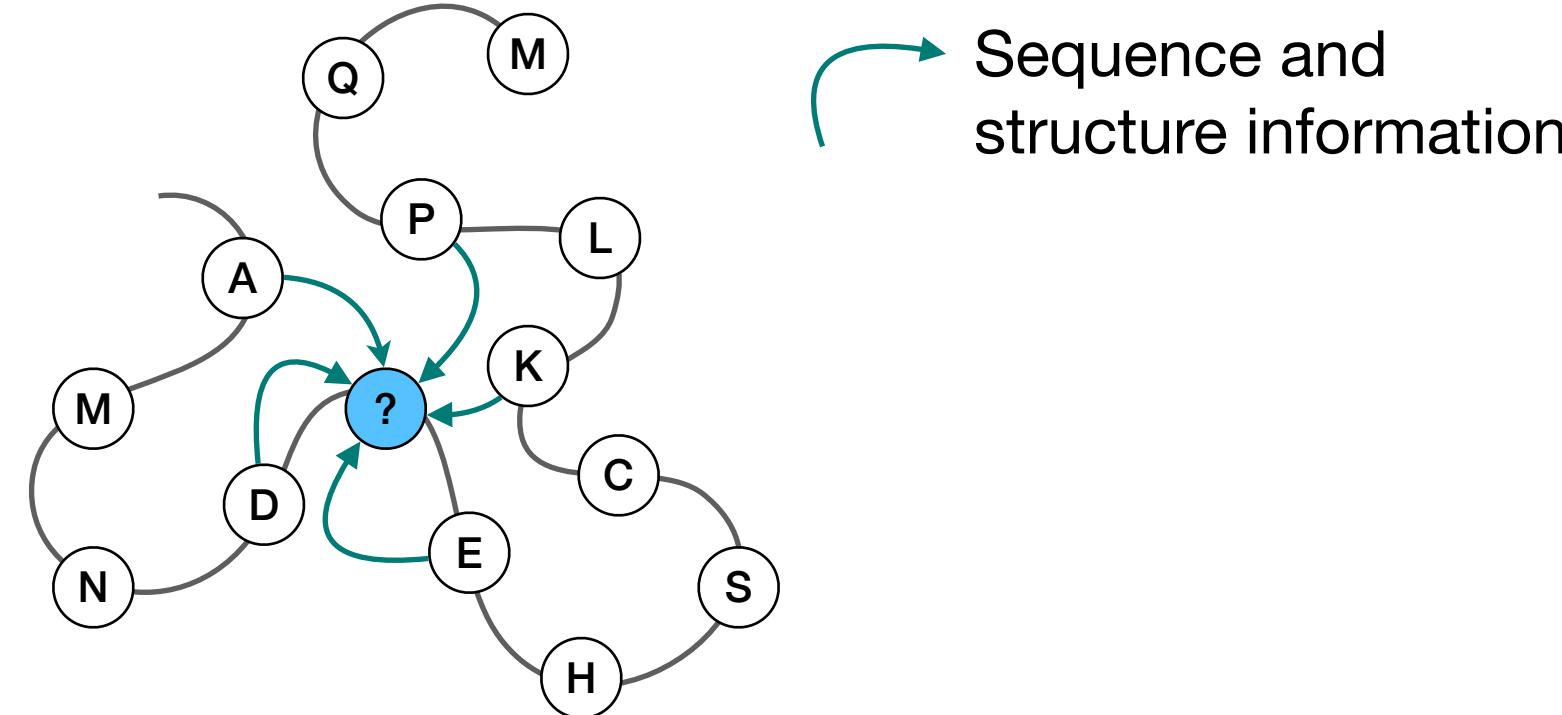
Masked language model



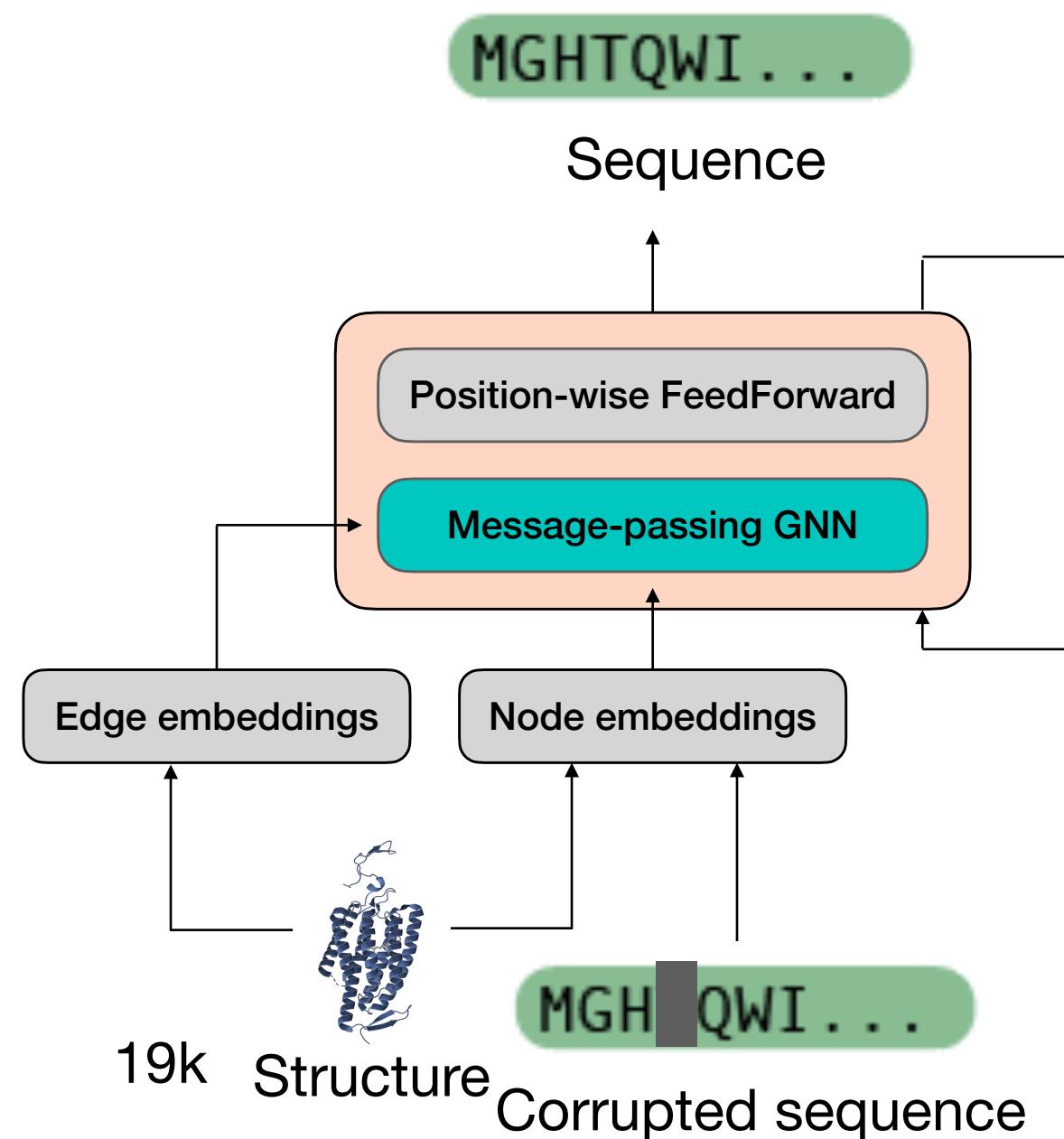
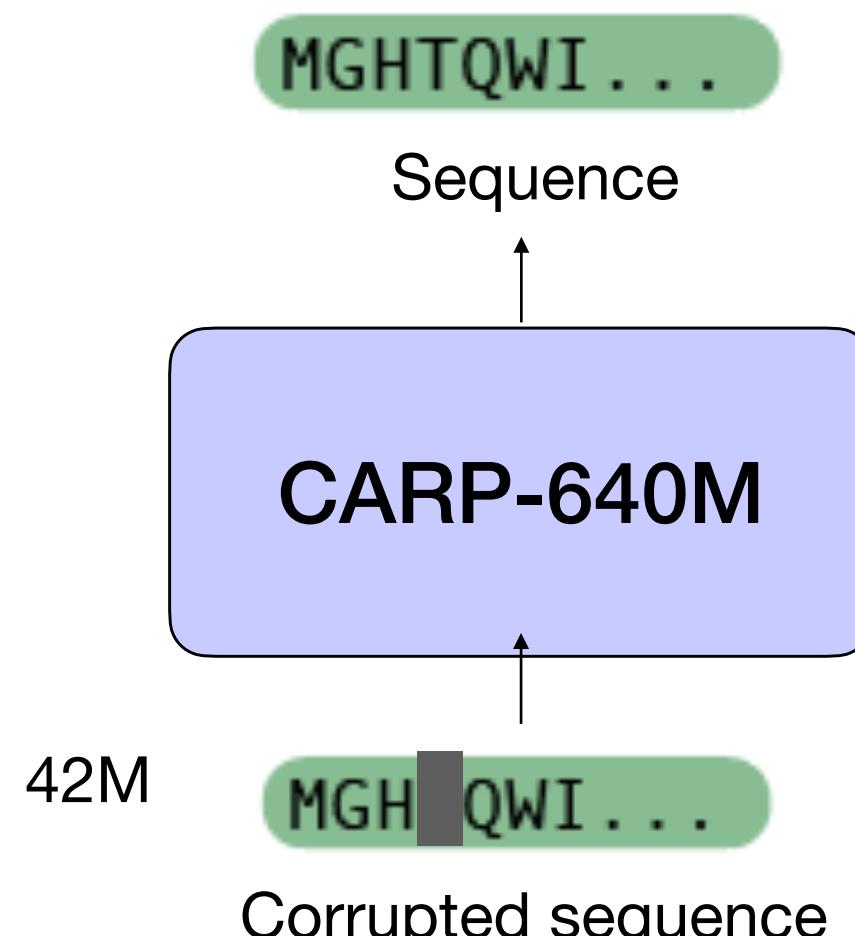
# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$



Masked language model



Model	Perplexity
-------	------------

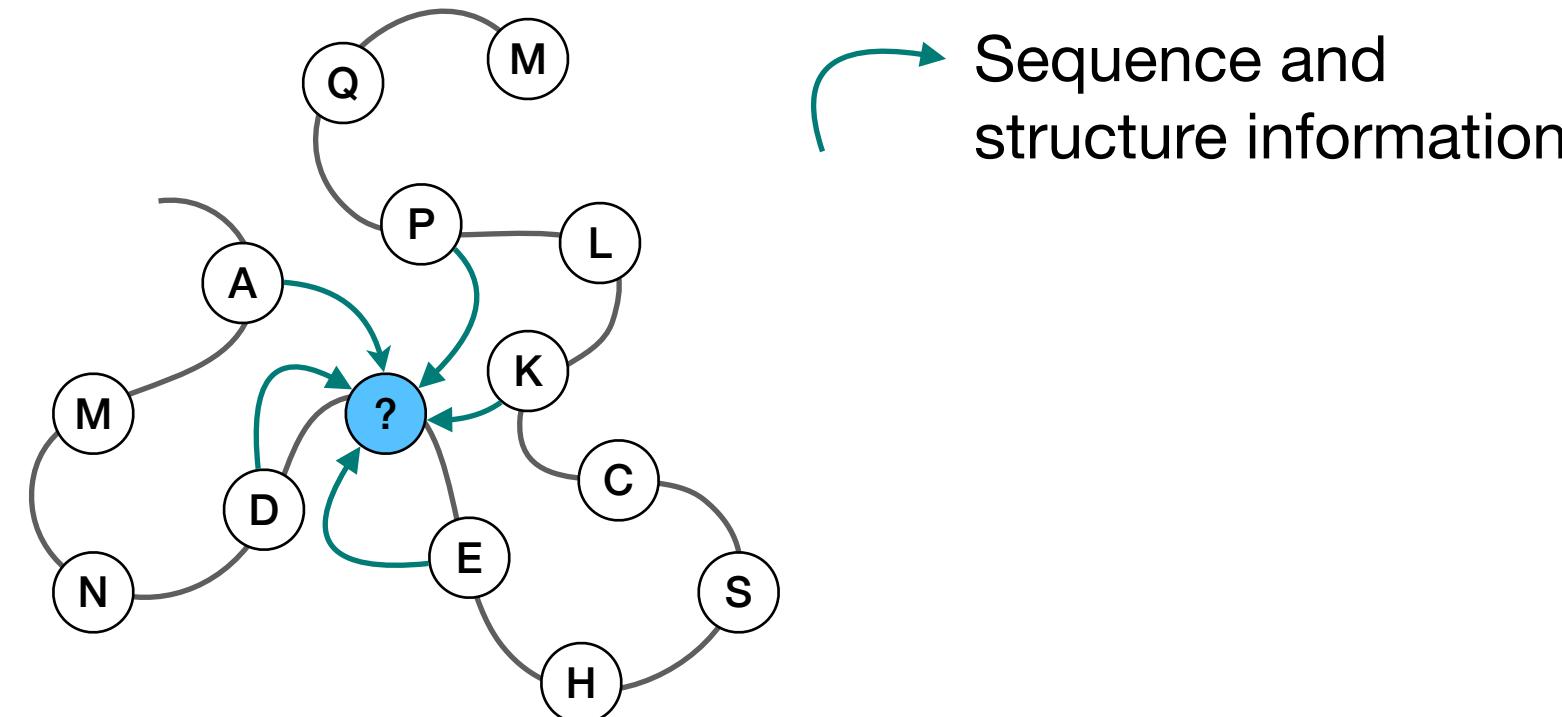
GNN	4.96
-----	------

Sequence only	7.06
---------------	------

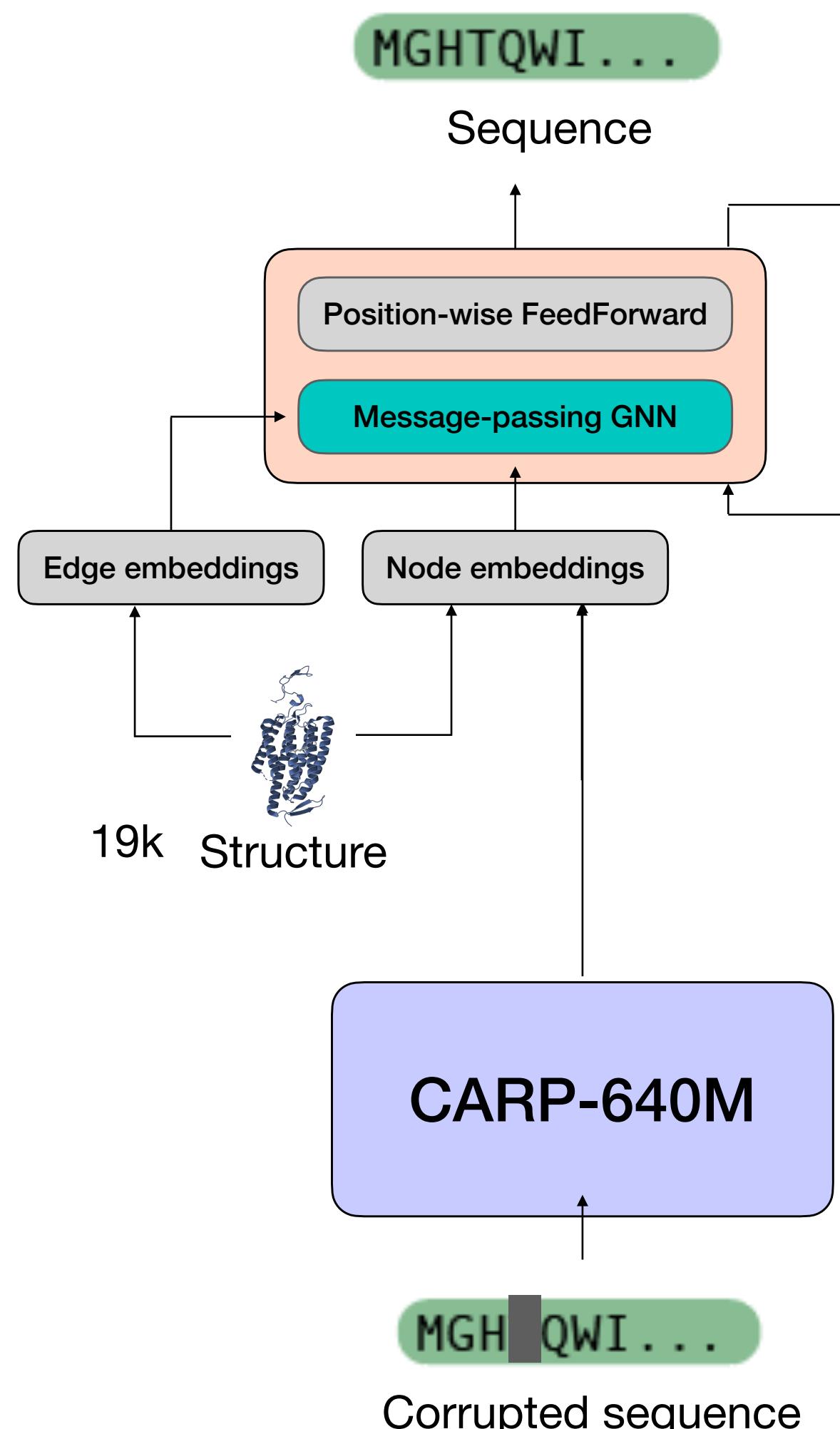
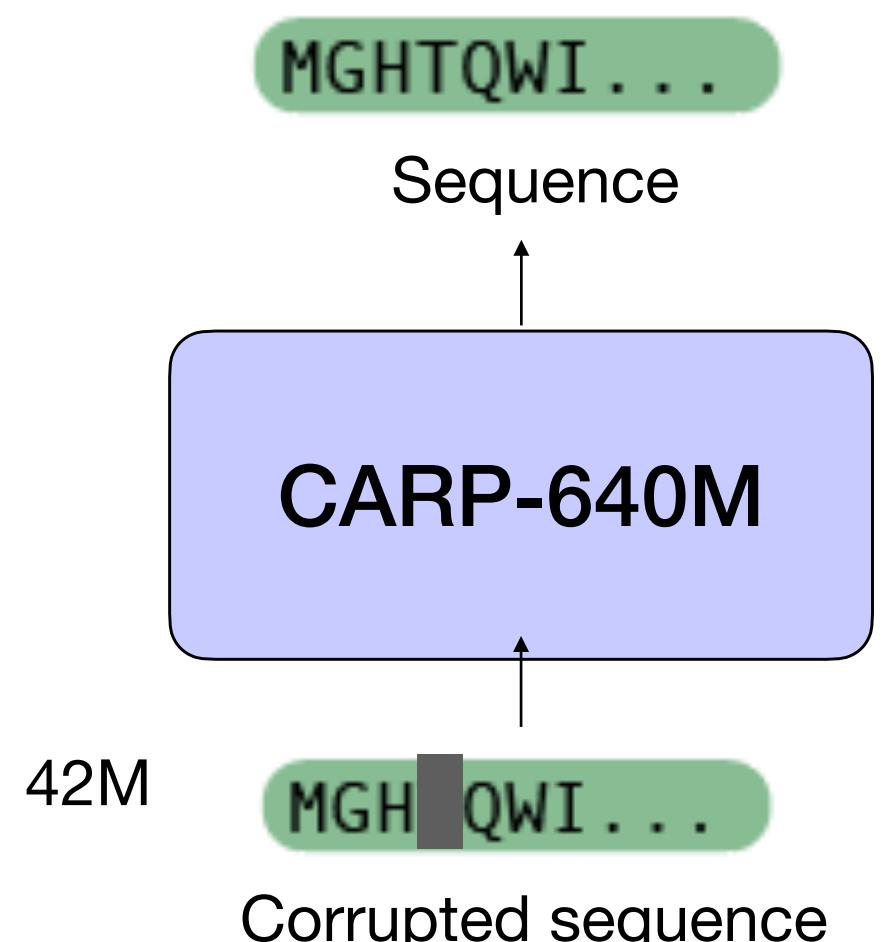
# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$



Masked language model

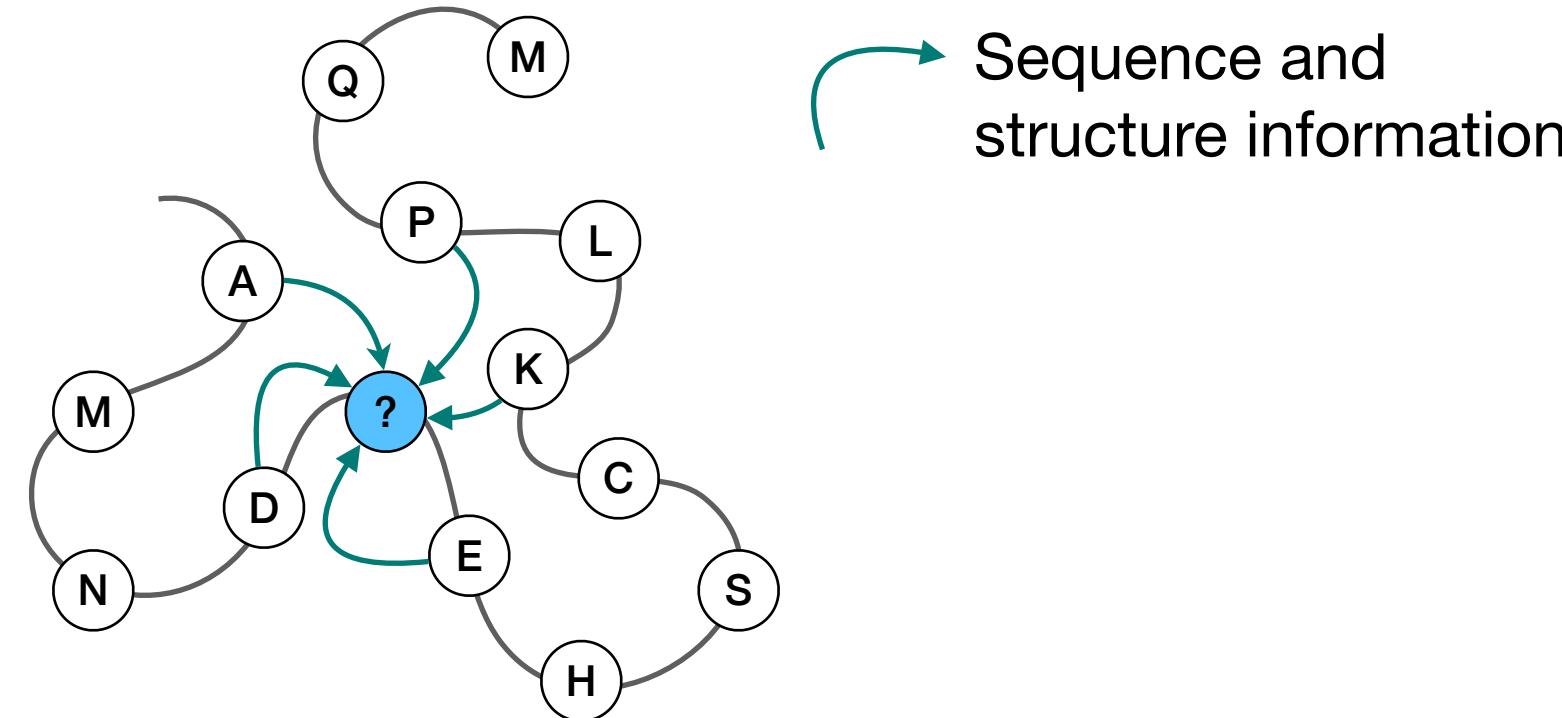


Model	Perplexity
GNN	4.96
Sequence only	7.06

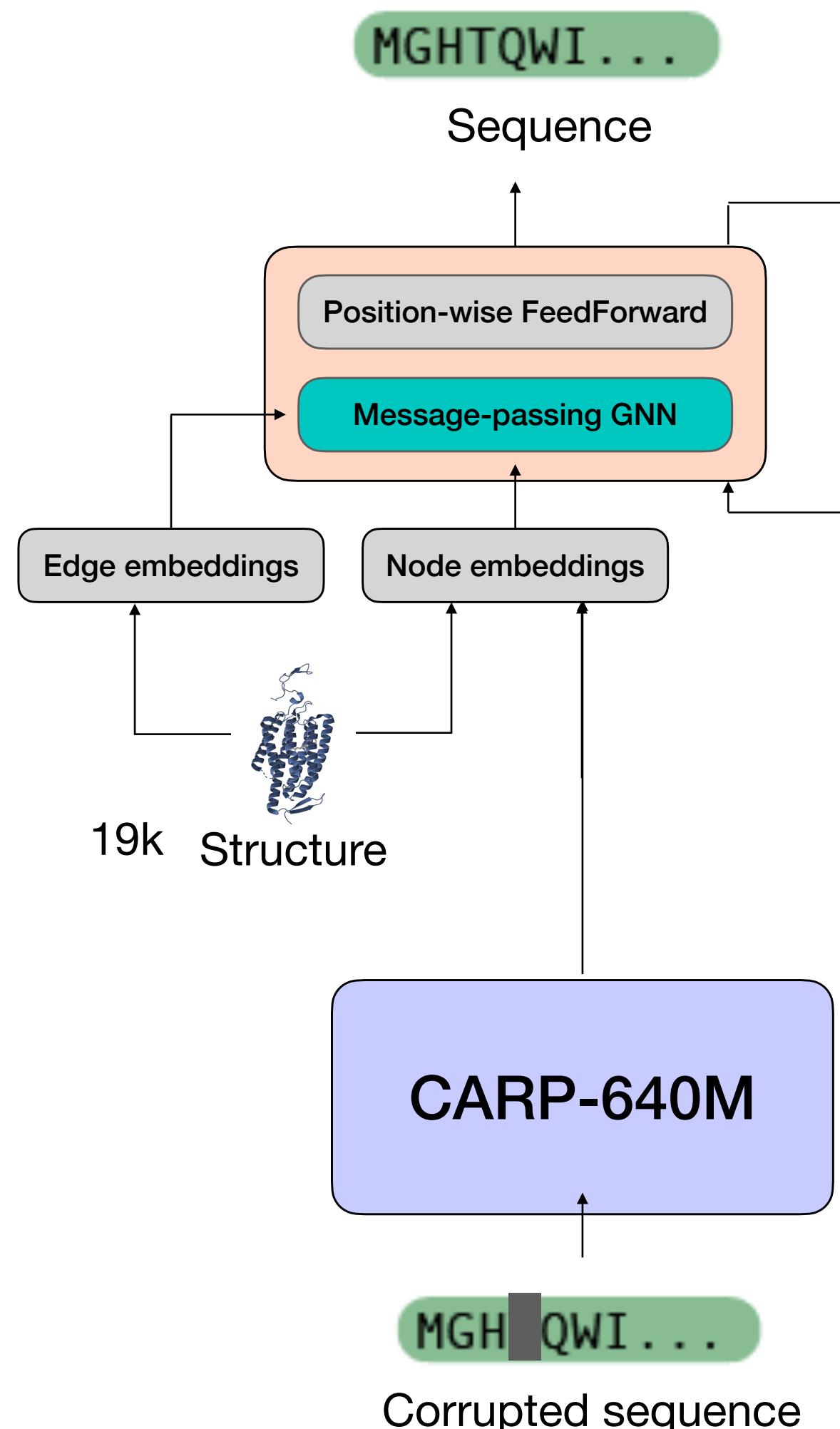
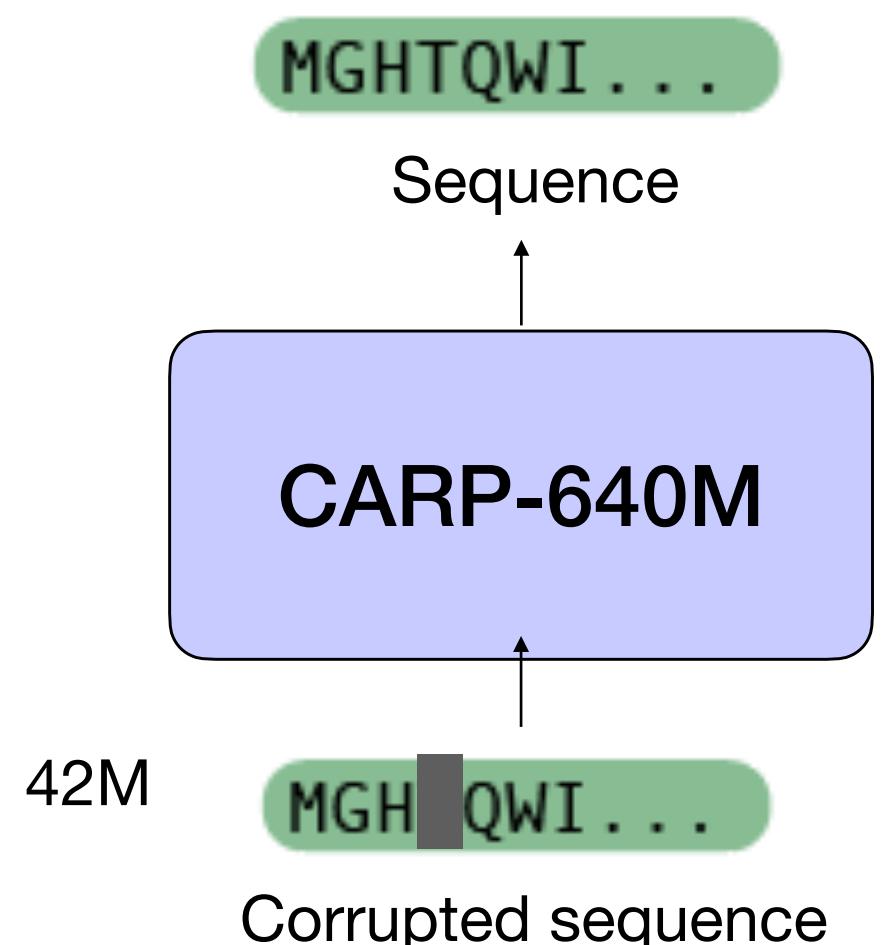
# Conditioning on structure improves sequence reconstruction

Structure-conditioned masked language model

Goal: Predict  $P(s_{\text{masked}} | \text{structure}, s_{\text{unmasked}})$

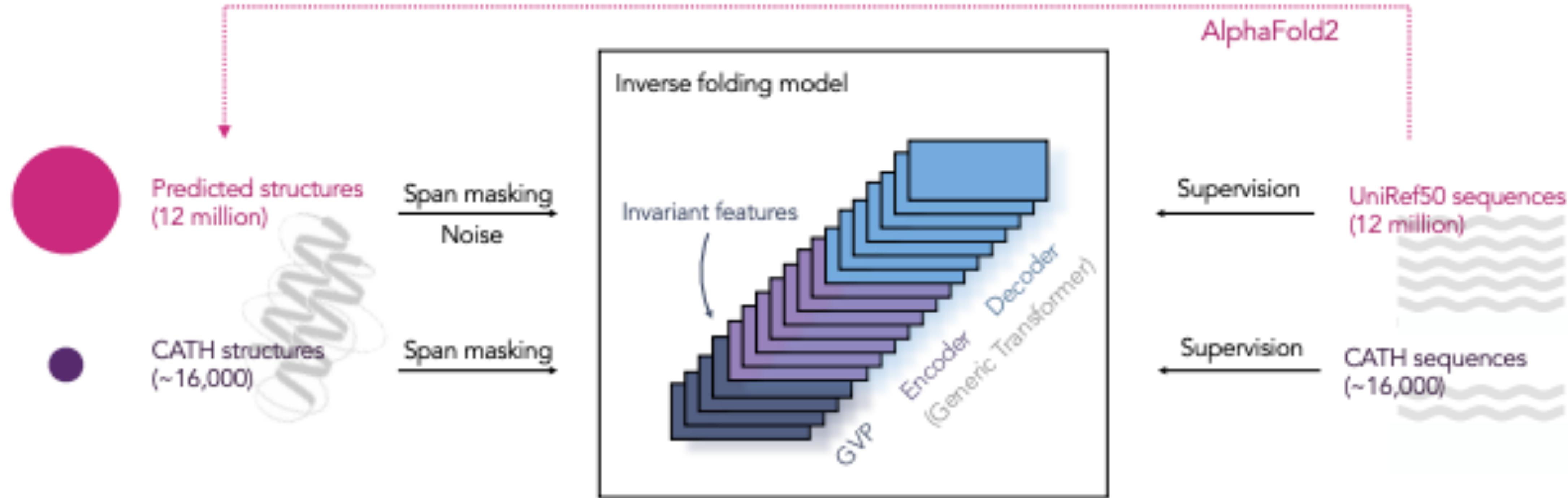


Masked language model

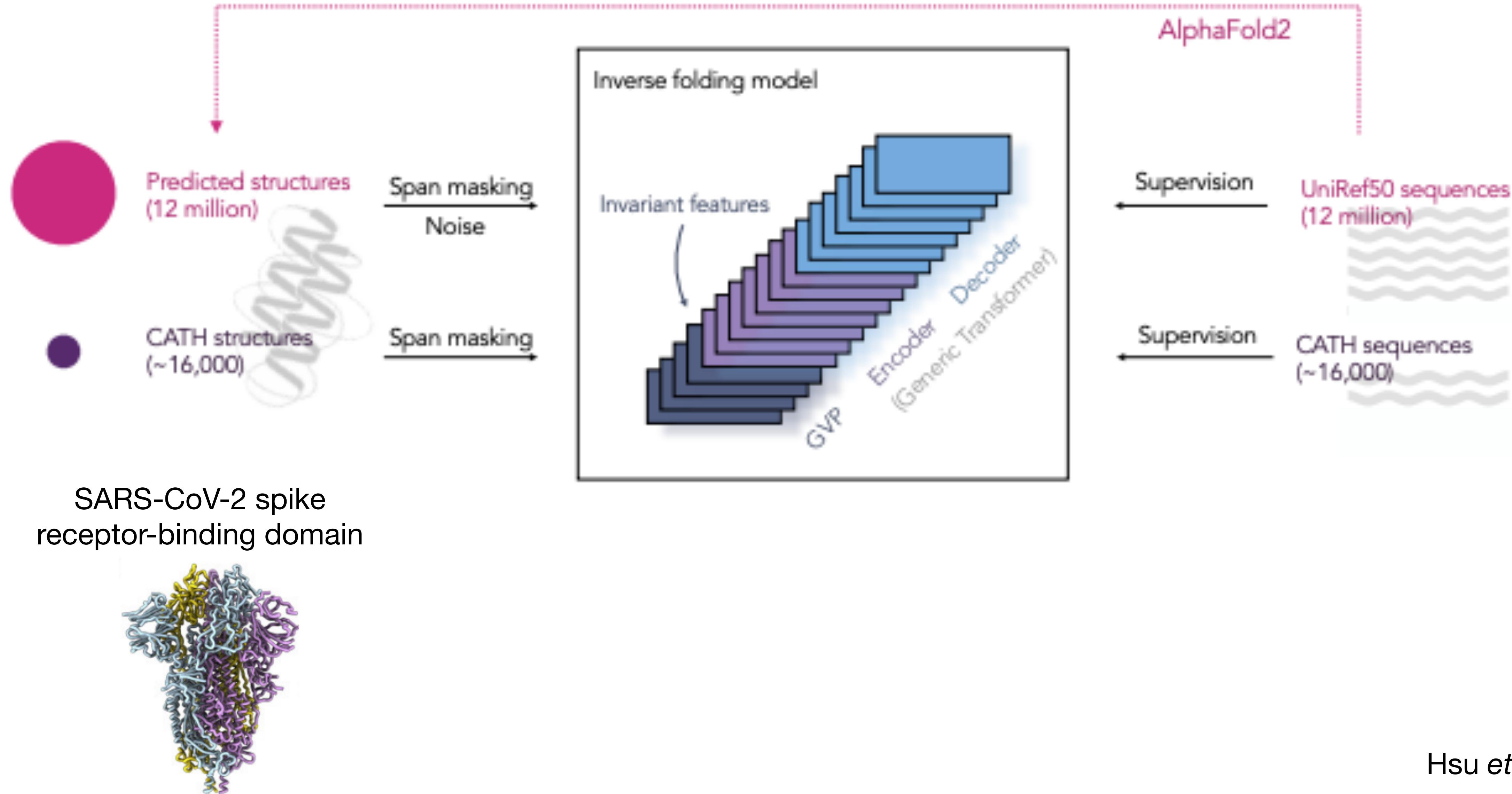


Model	Perplexity
GNN	4.96
Sequence only	7.06
Pretrained logits -> GNN	4.08

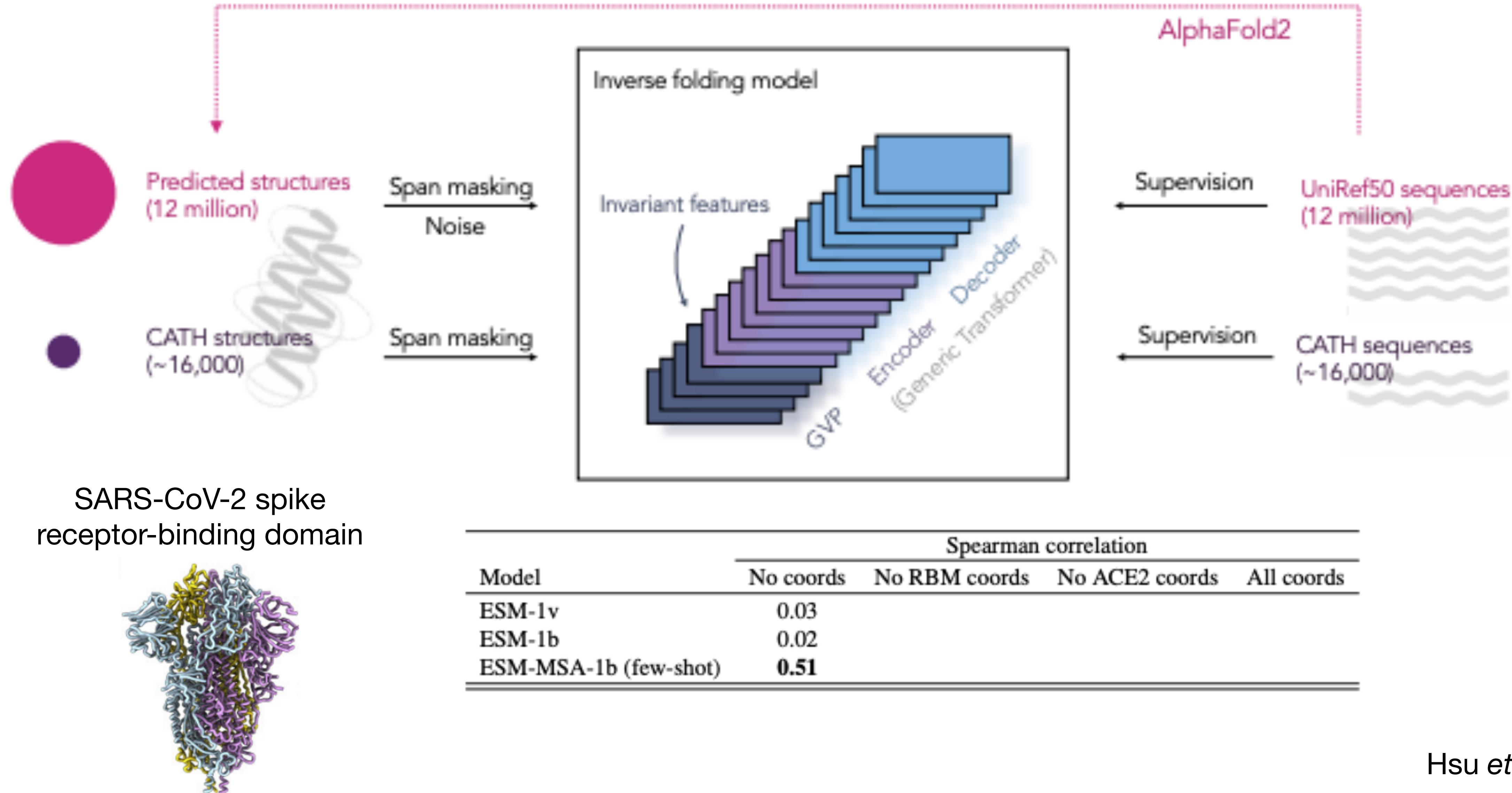
# Adding predicted structures improves sequence reconstruction



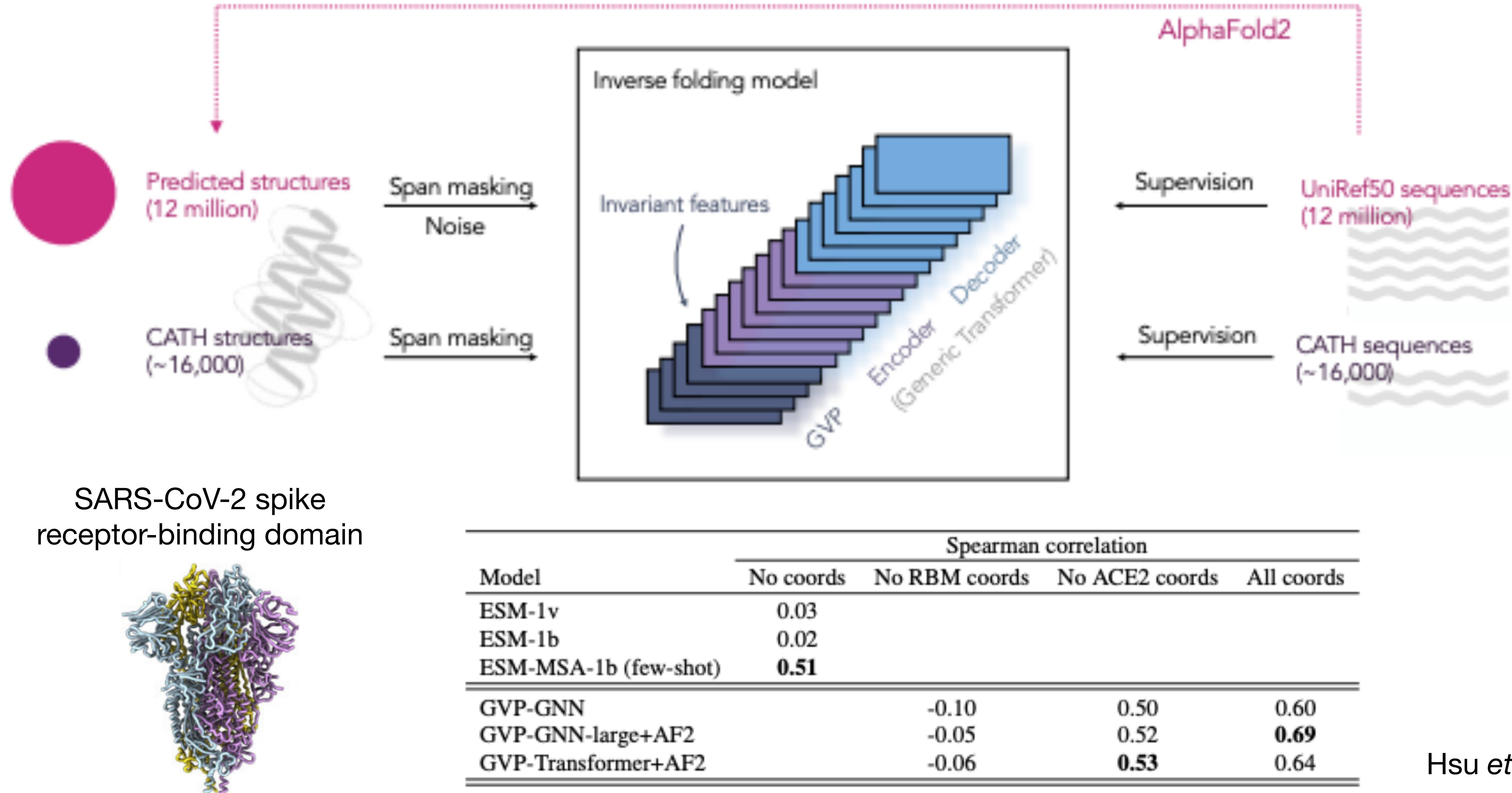
# Adding predicted structures improves sequence reconstruction



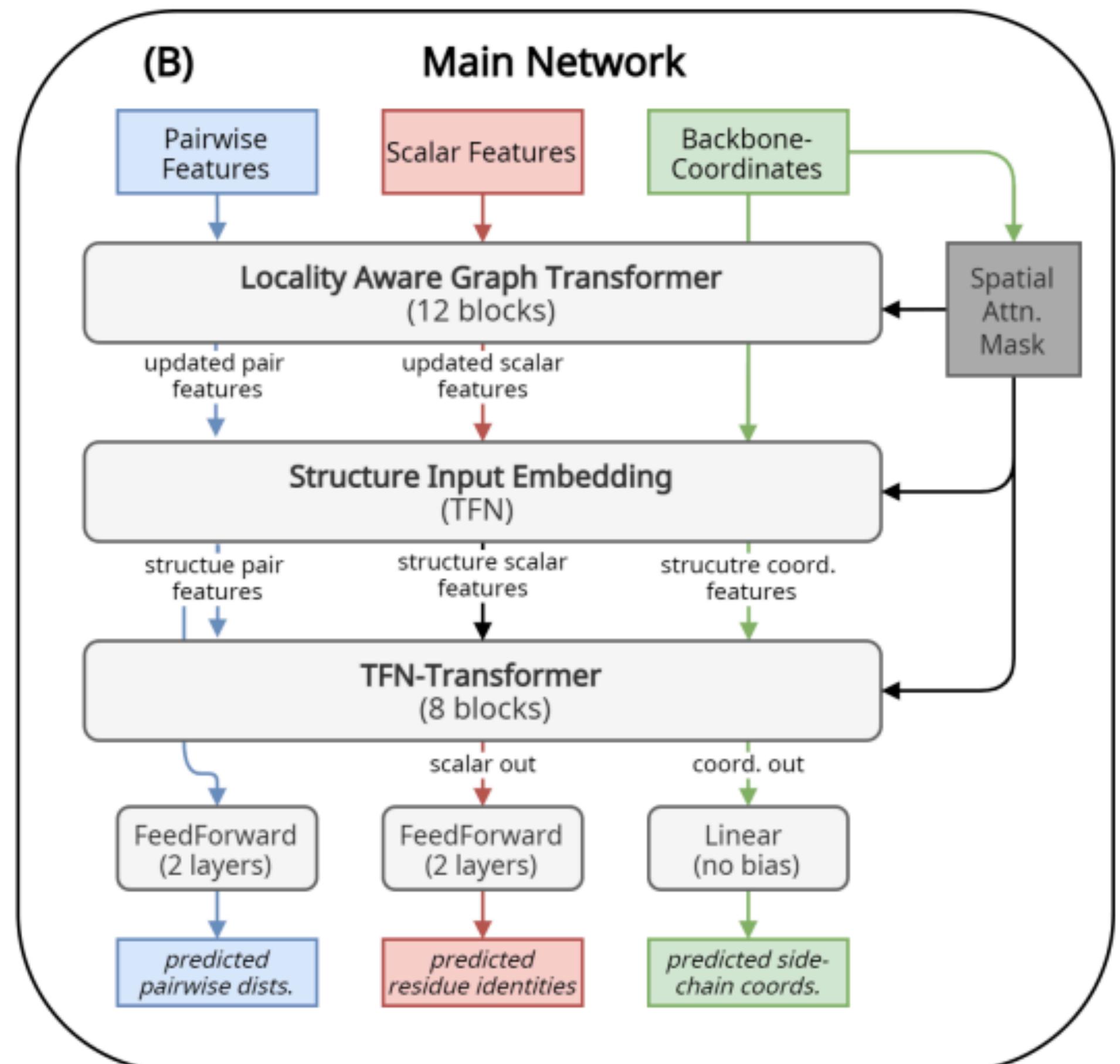
# Adding predicted structures improves sequence reconstruction



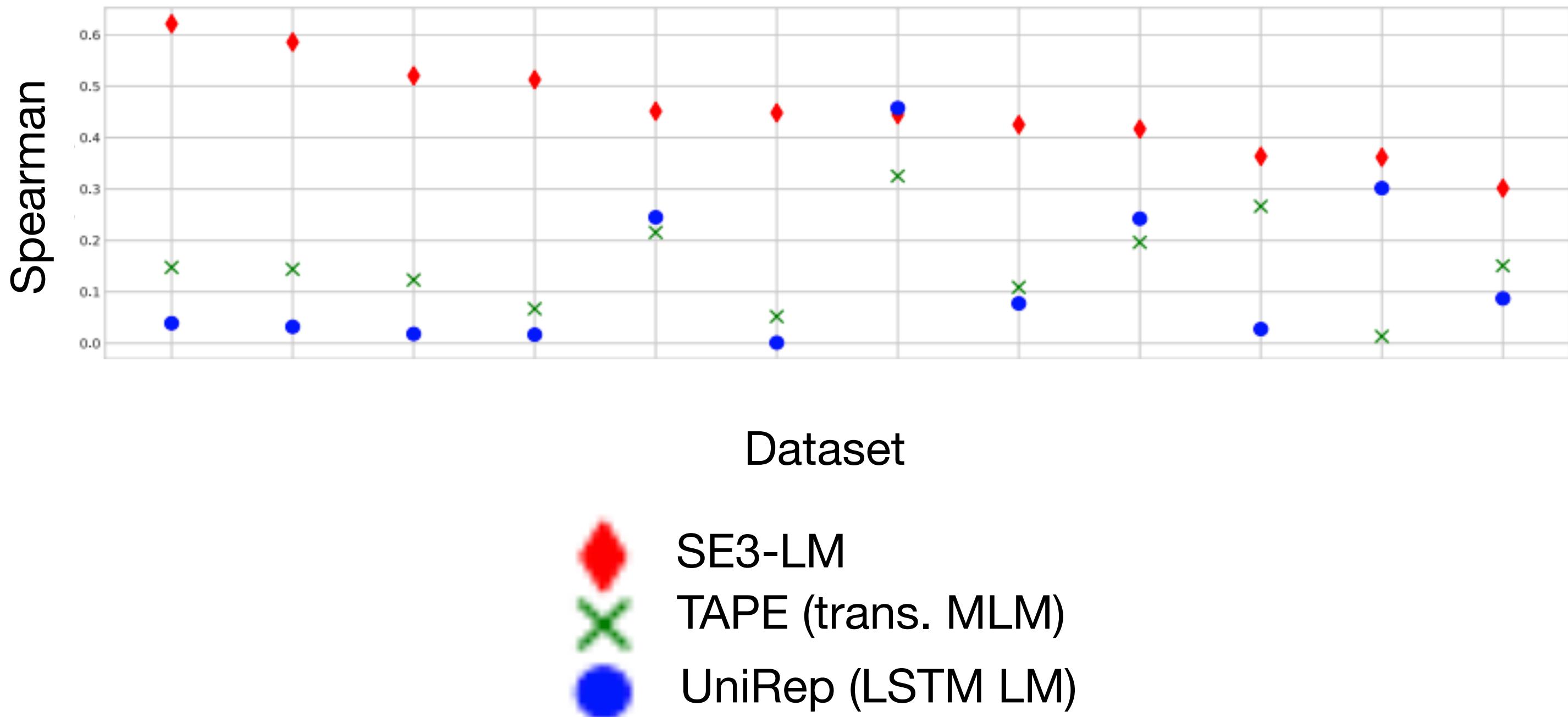
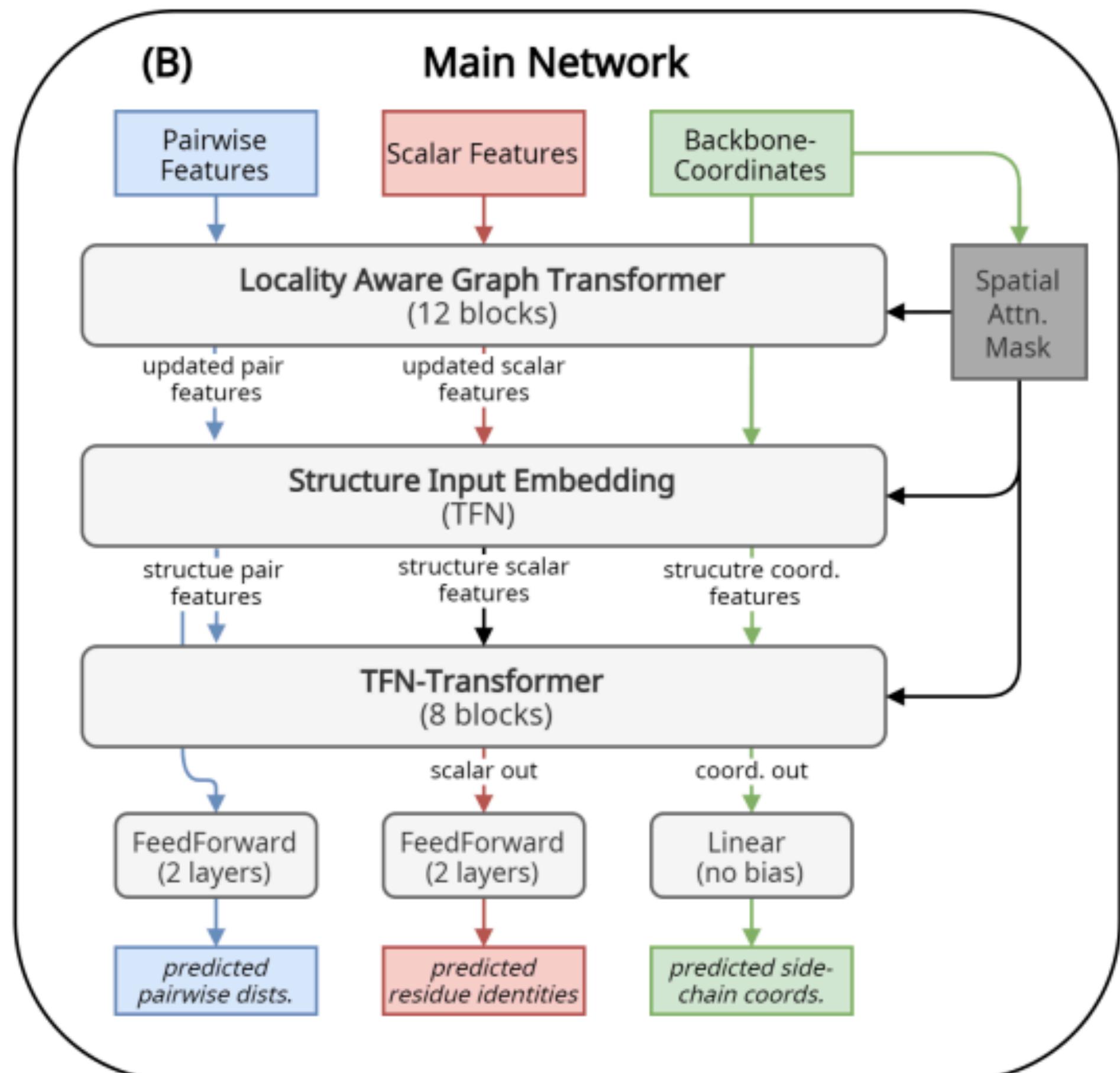
# Adding predicted structures improves sequence reconstruction



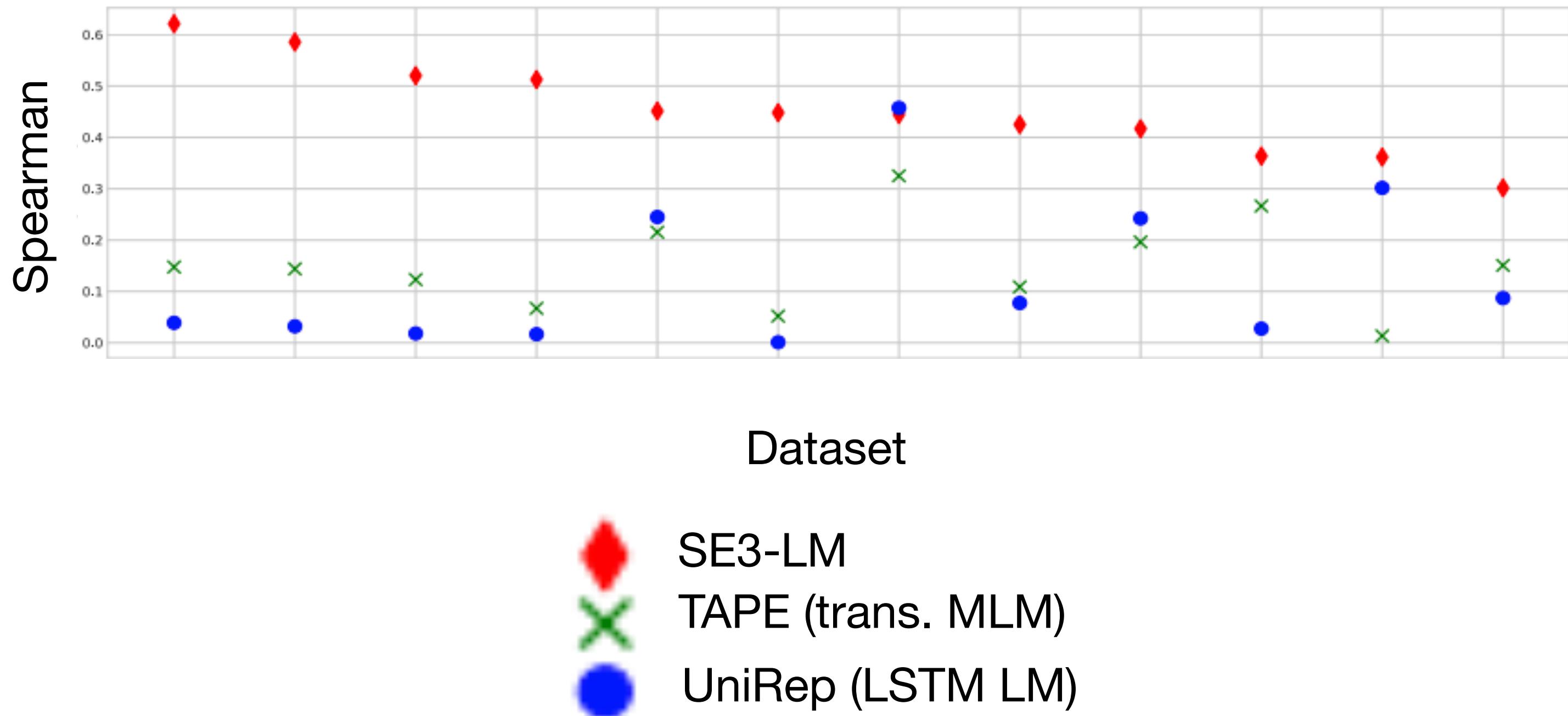
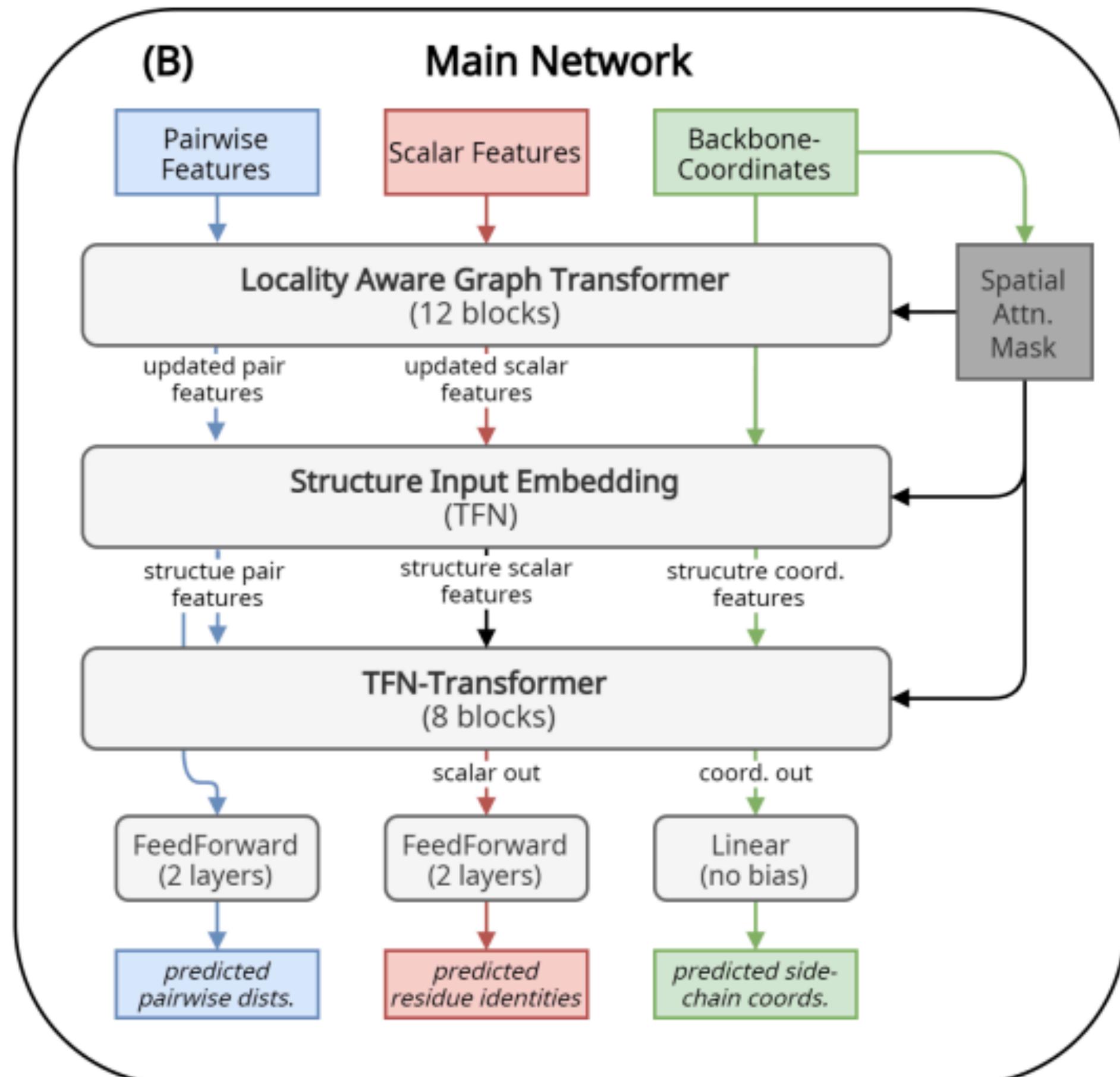
# Adding predicted structures improves sequence reconstruction



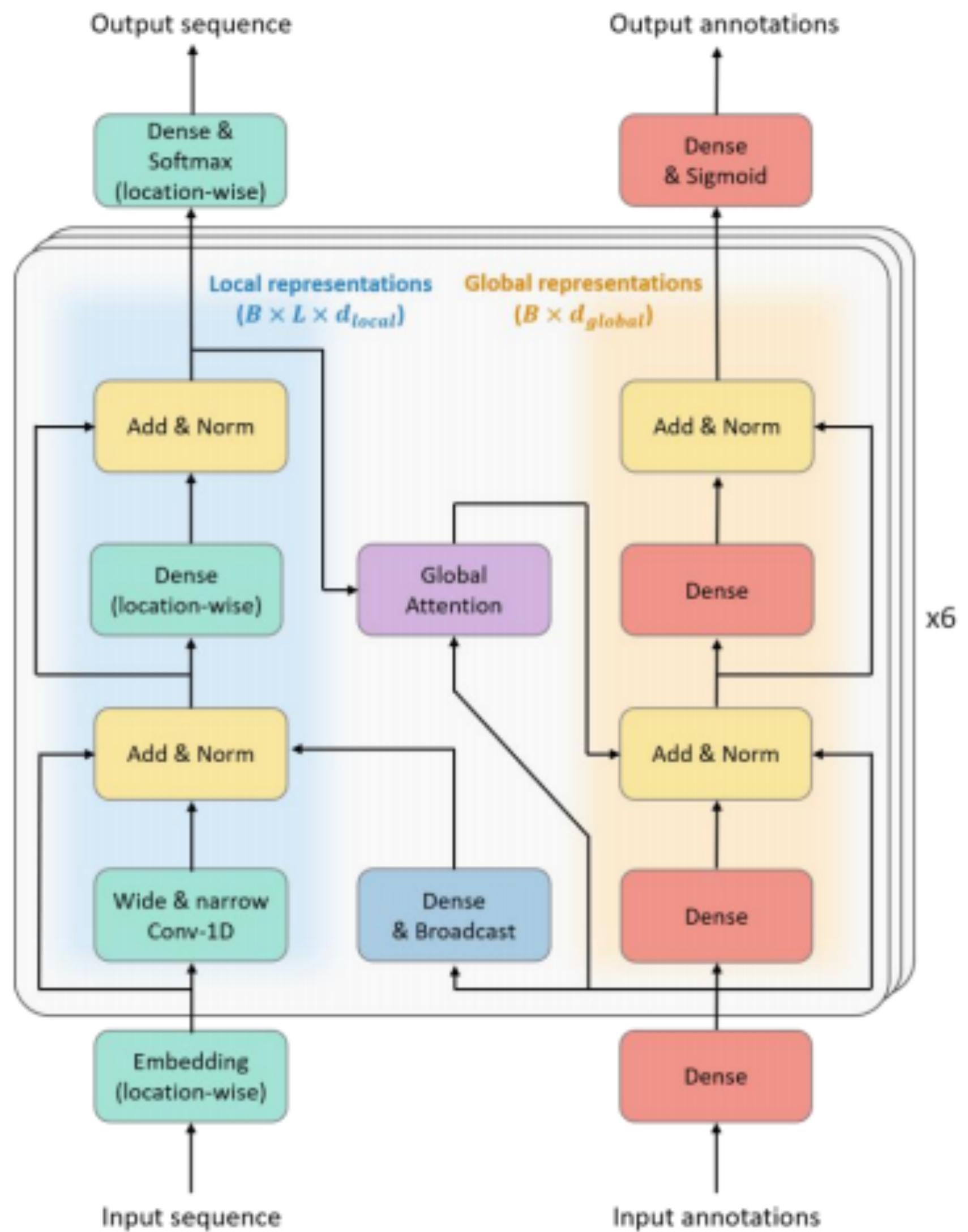
# Adding predicted structures improves sequence reconstruction



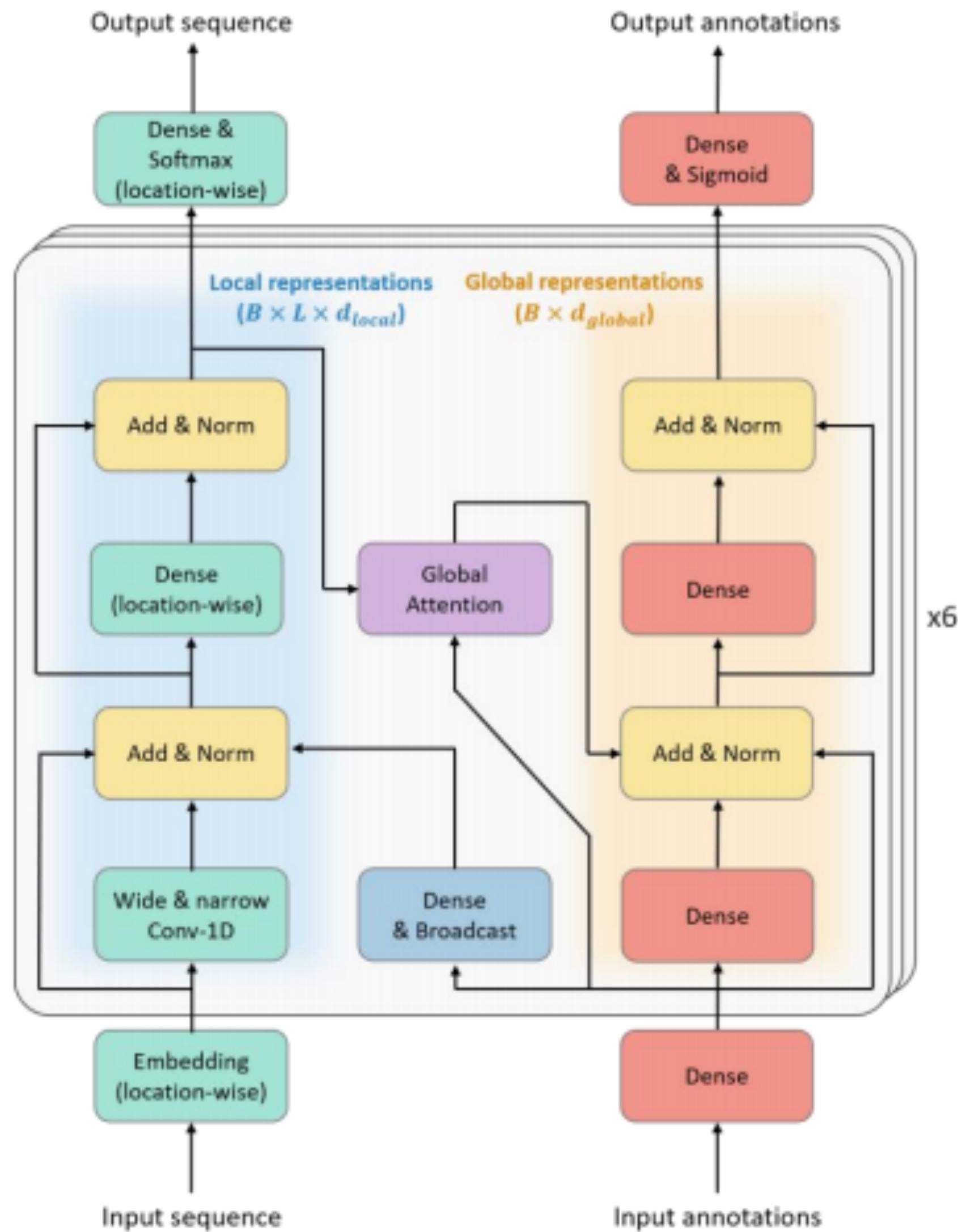
# Adding predicted structures improves sequence reconstruction



# Combine sequence and annotations

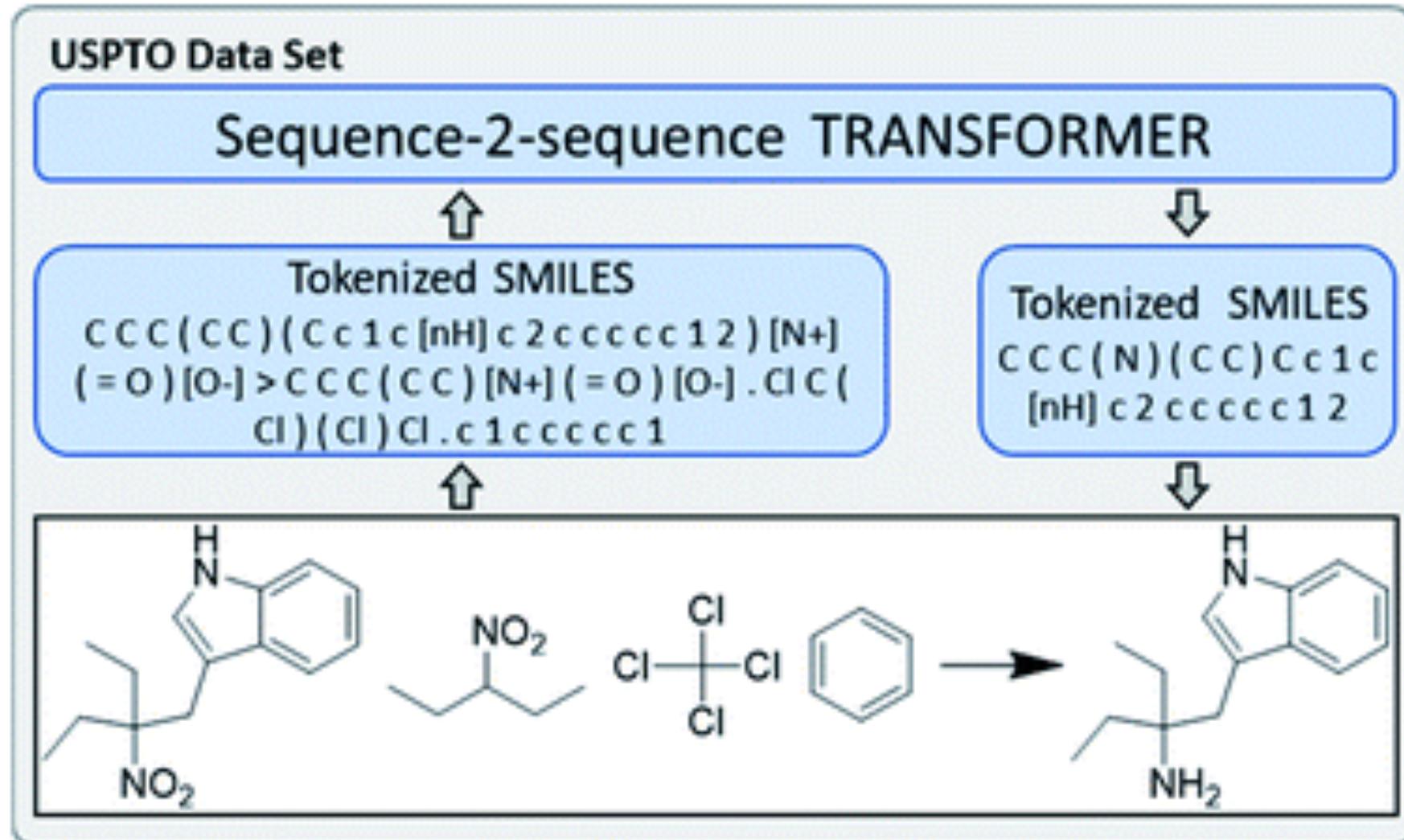


# Combine sequence and annotations

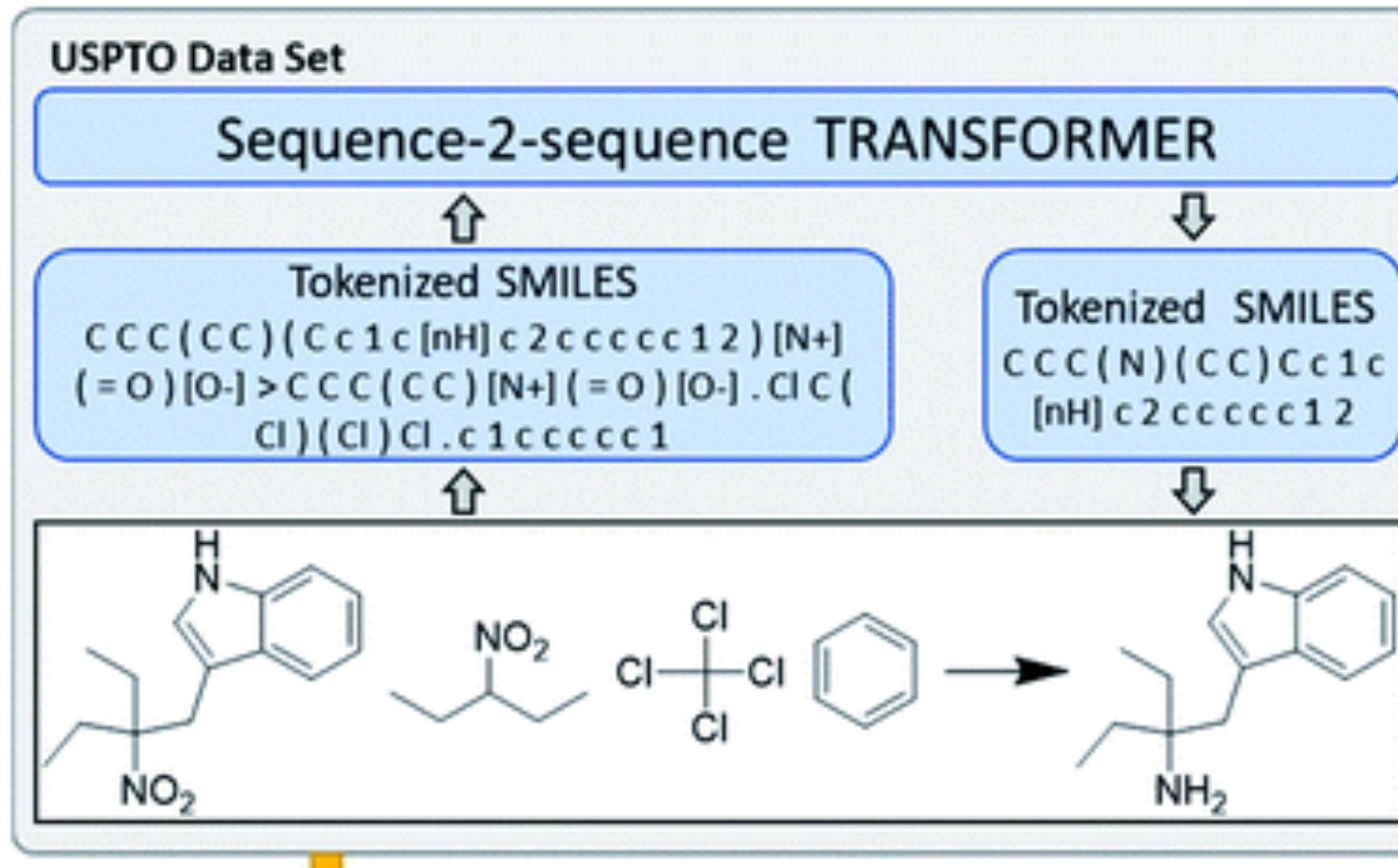


Method	Structure	Evolutionary	Engineering		
	Secondary structure	Remote homology	Fluorescence	Stability	
Without Pretraining	<b>TAPE Transformer</b>	0.70	0.09	0.22	-0.06
	<b>LSTM</b>	0.71	0.12	0.21	0.28
With Pretraining	<b>ProteinBERT</b>	0.70	0.06	0.65	0.63
	<b>TAPE Transformer</b>	0.73	0.21	0.68	0.73
	<b>LSTM</b>	0.75	0.26	0.67	0.69
	<b>UniRep mLSTM</b>	0.73	0.23	0.67	0.73
	<b>ProteinBERT</b>	0.74	0.22	0.66	0.76

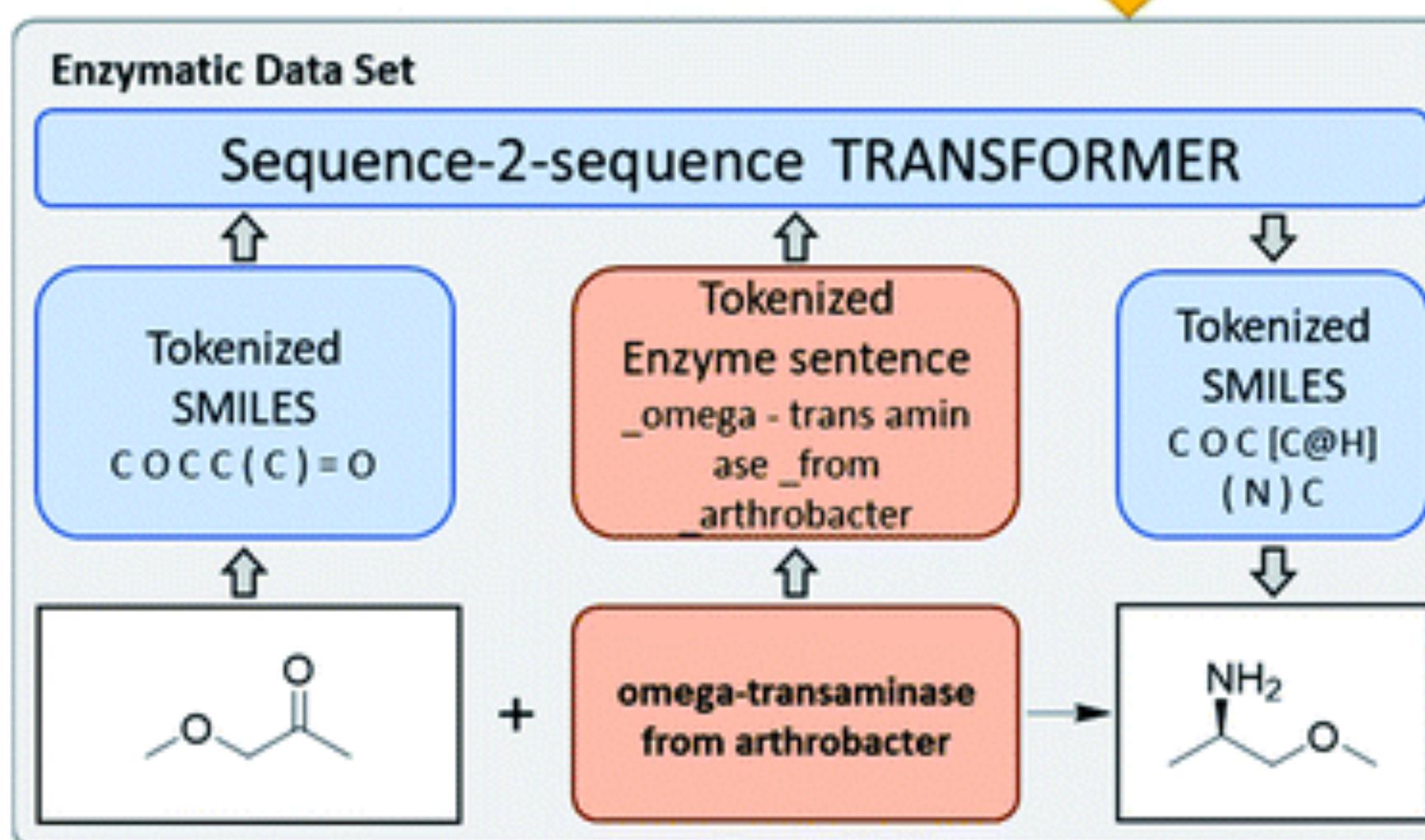
# Predict product from substrates and text



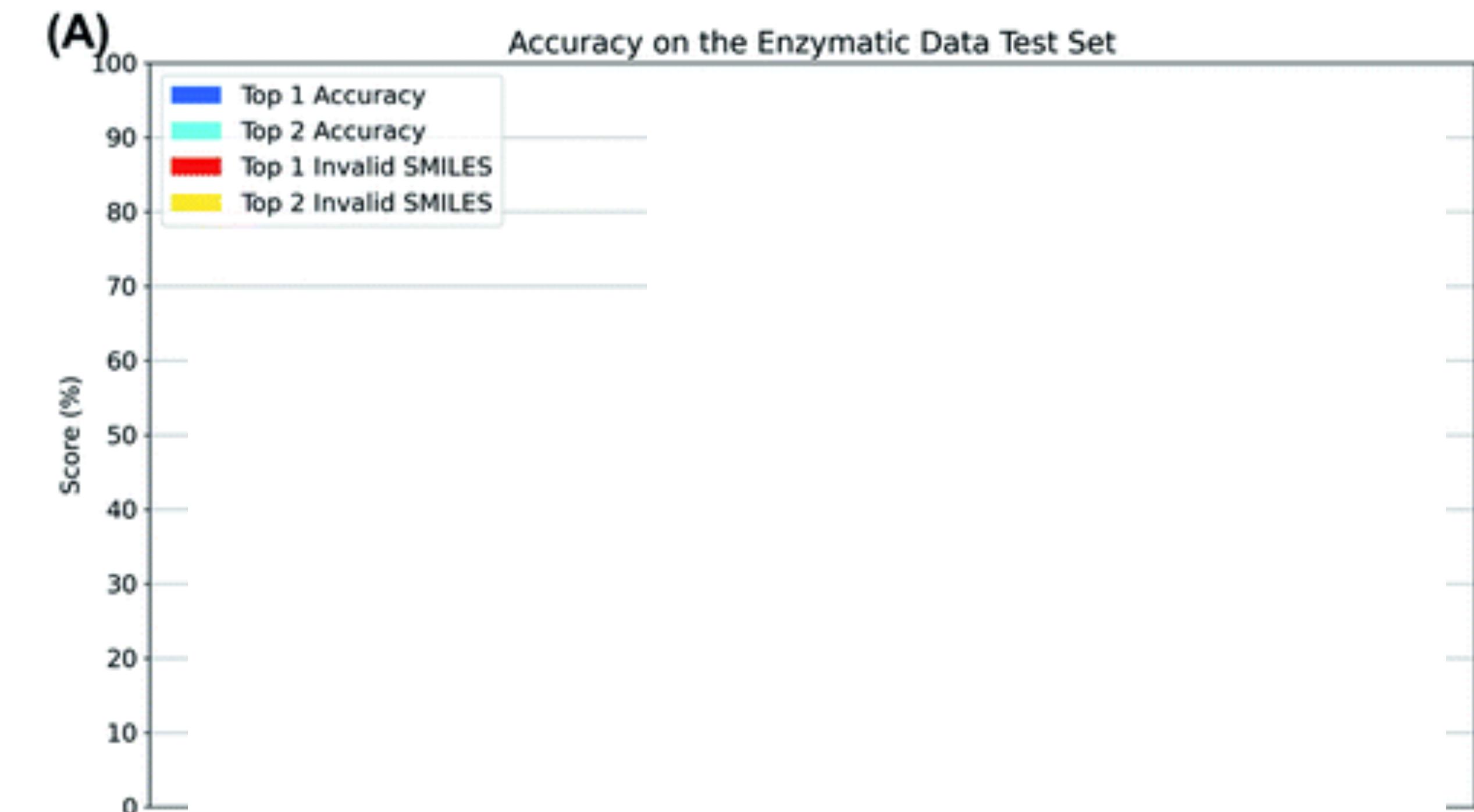
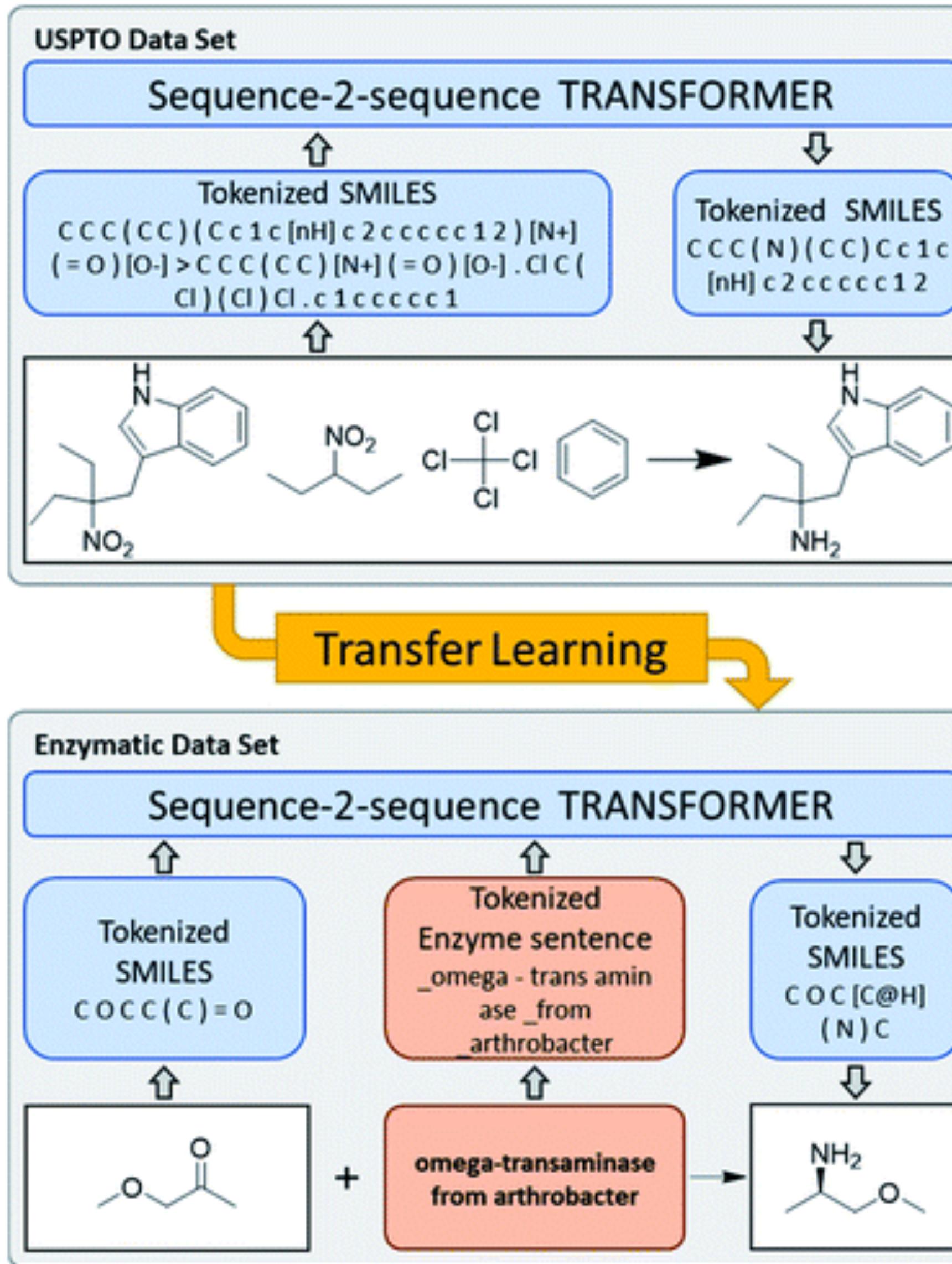
# Predict product from substrates and text



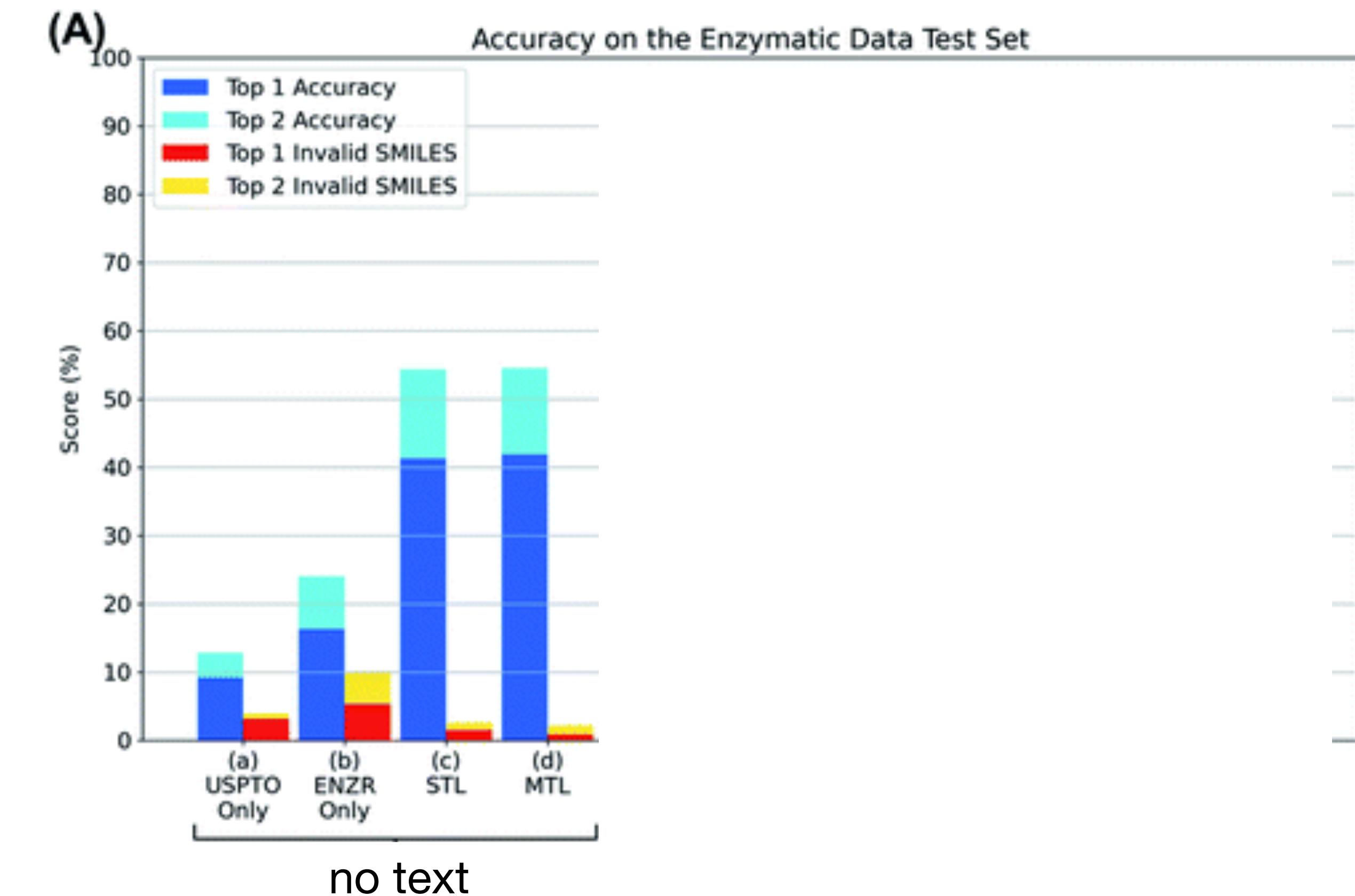
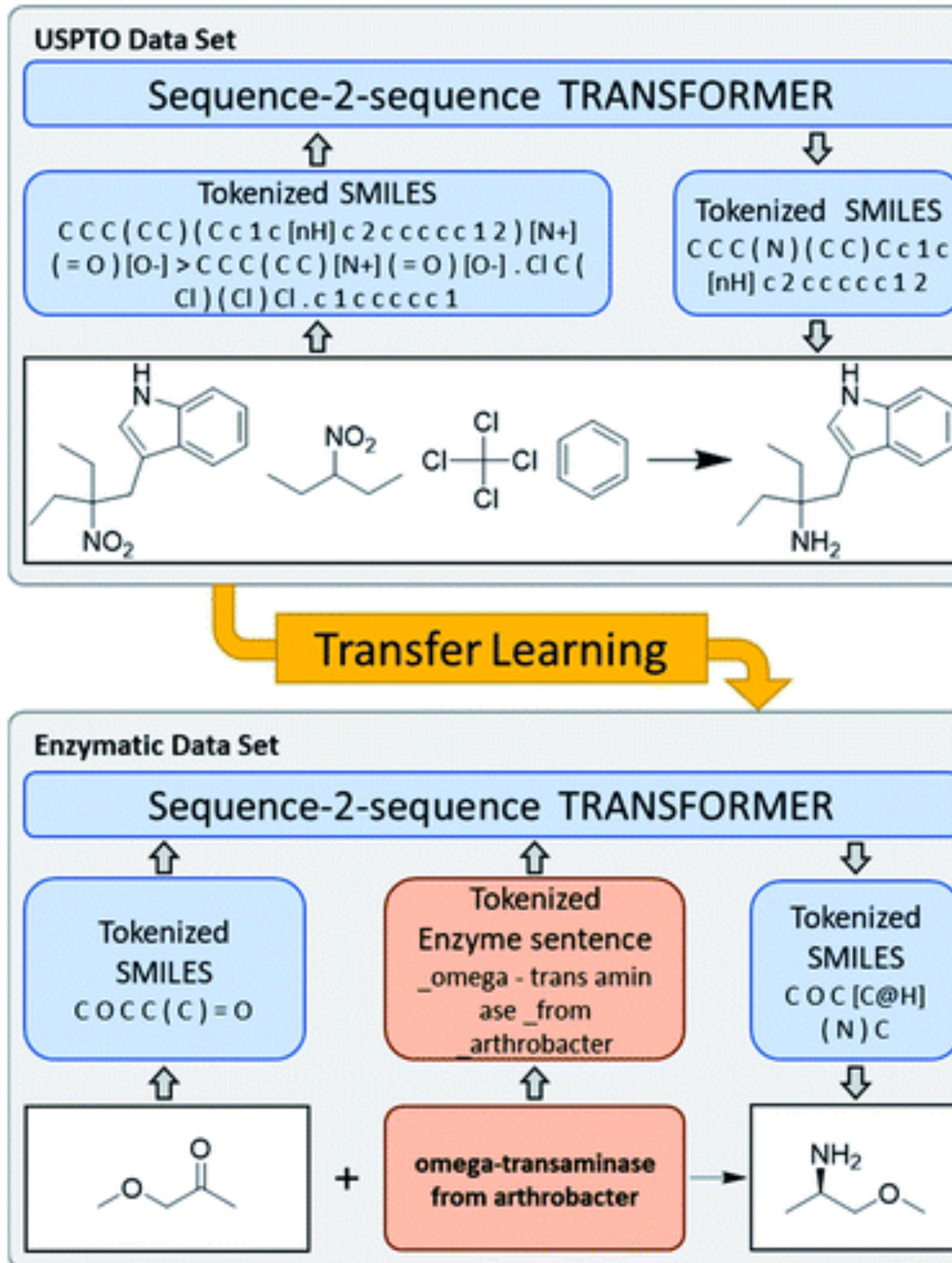
Transfer Learning



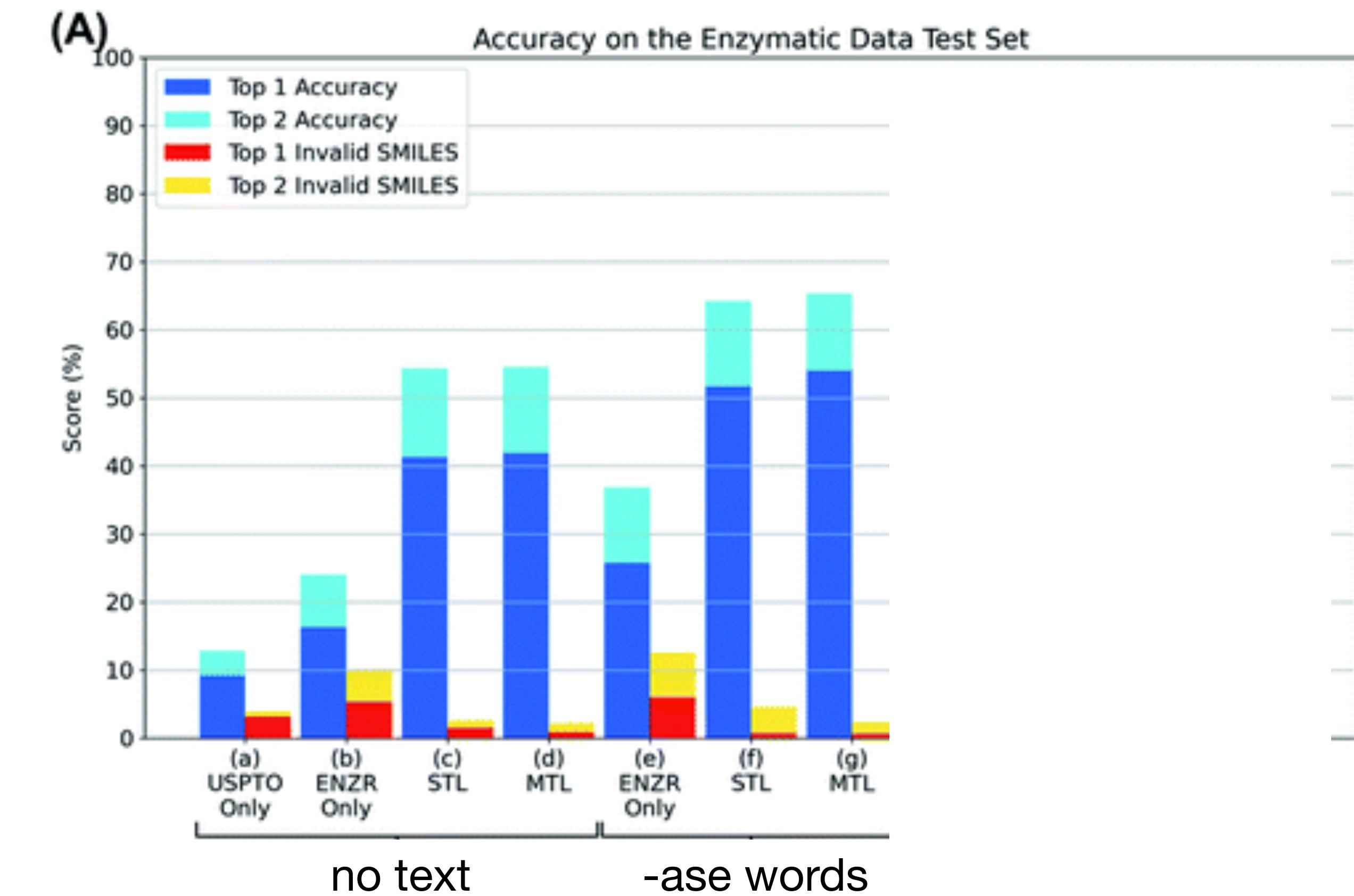
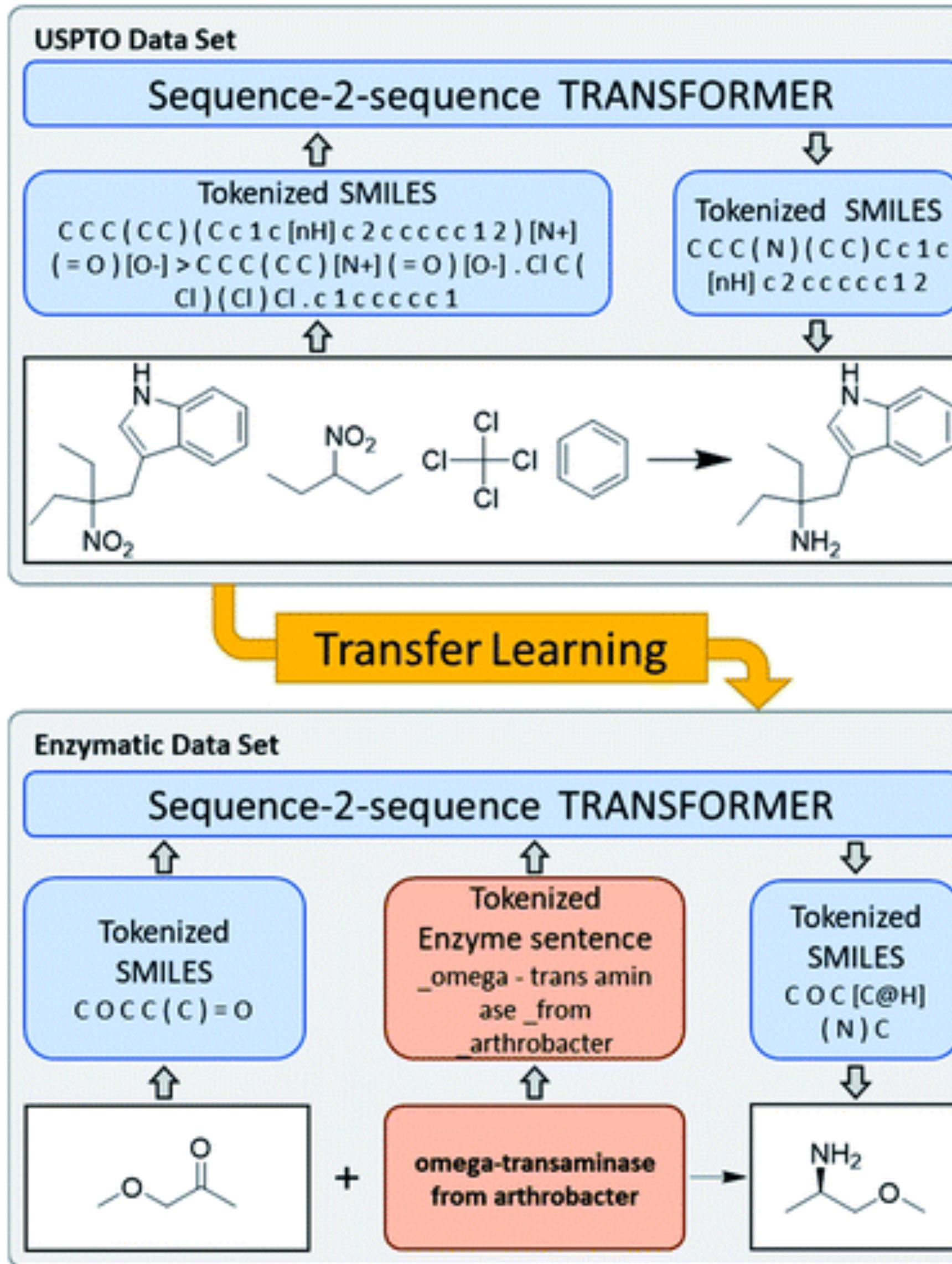
# Predict product from substrates and text



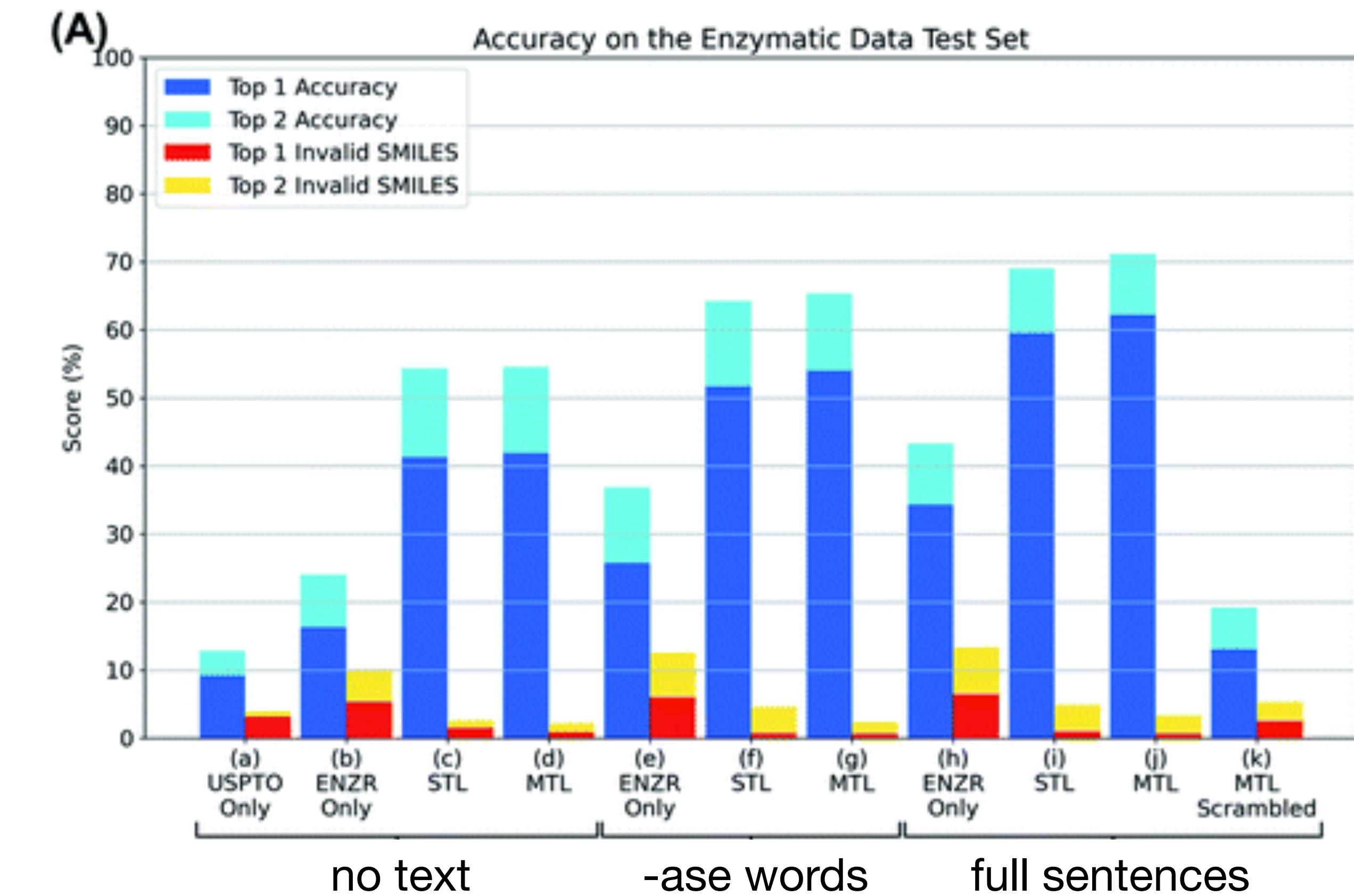
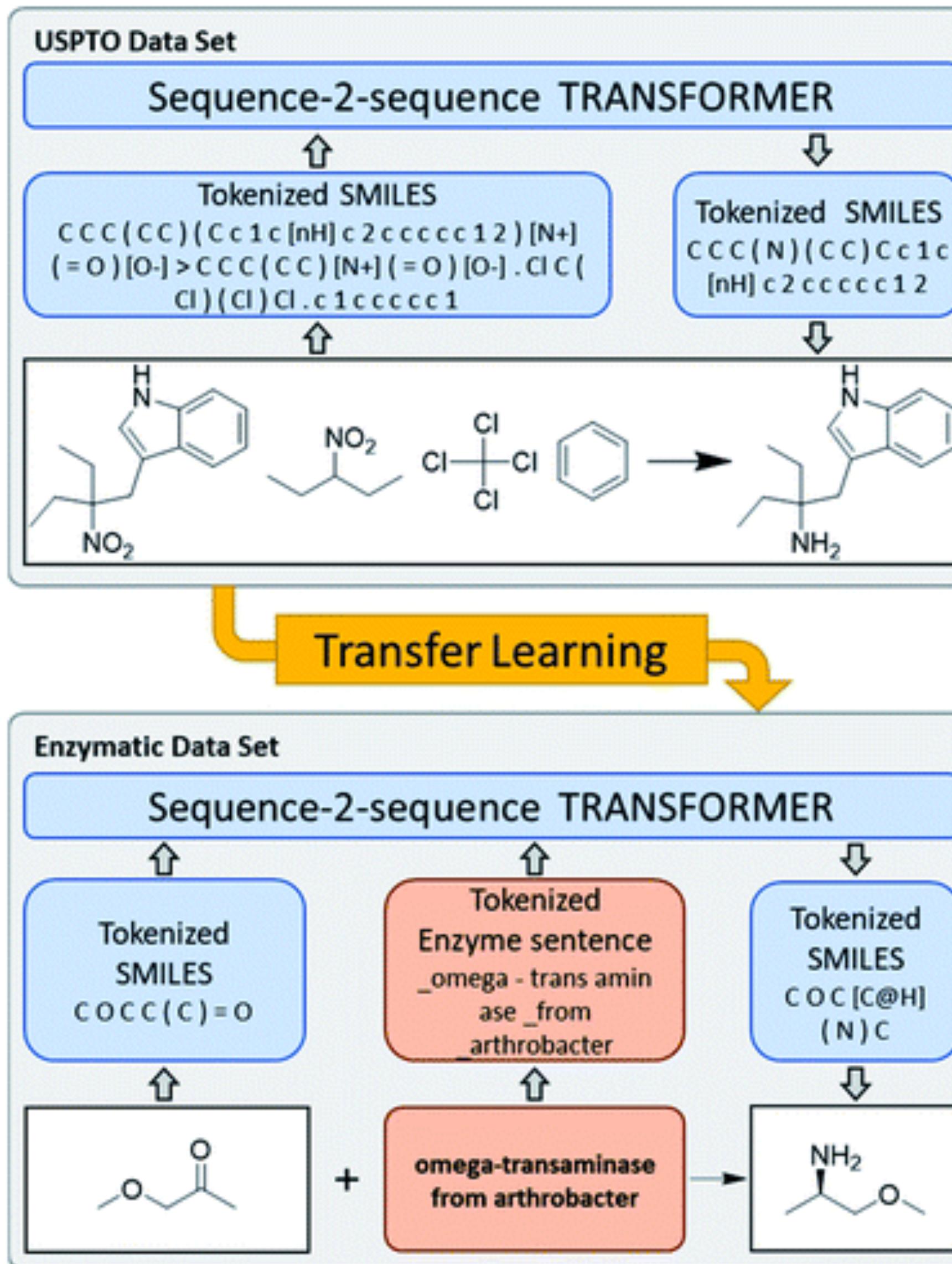
# Predict product from substrates and text



# Predict product from substrates and text



# Predict product from substrates and text



# Questions



@KevinKaichuang

[yang.kevin@microsoft.com](mailto:yang.kevin@microsoft.com)

<https://www.ml4proteinengineering.com/>