

**Genentech**  
*A Member of the Roche Group*

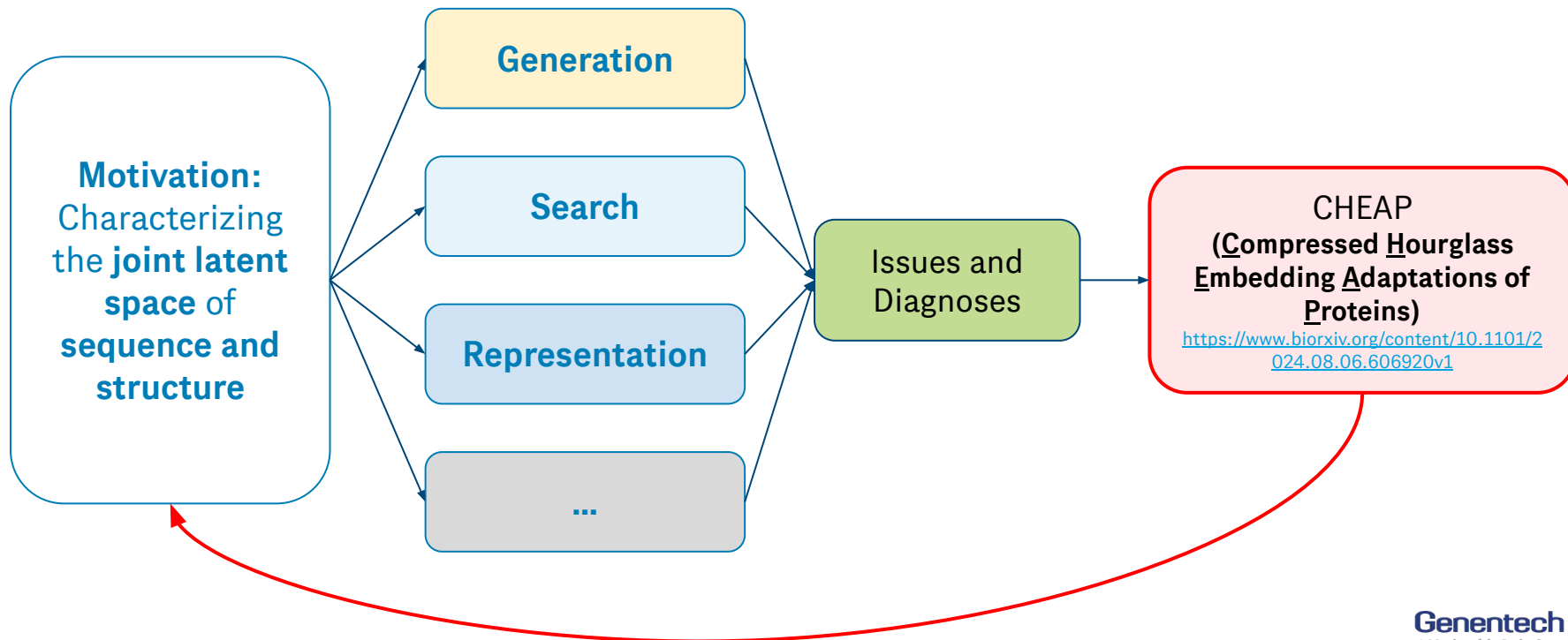
# *Tokenized and Continuous Embedding Compressions of Protein Sequence and Structure*

*October 22, 2024  
Stanford AI + Biomedicine Seminar*

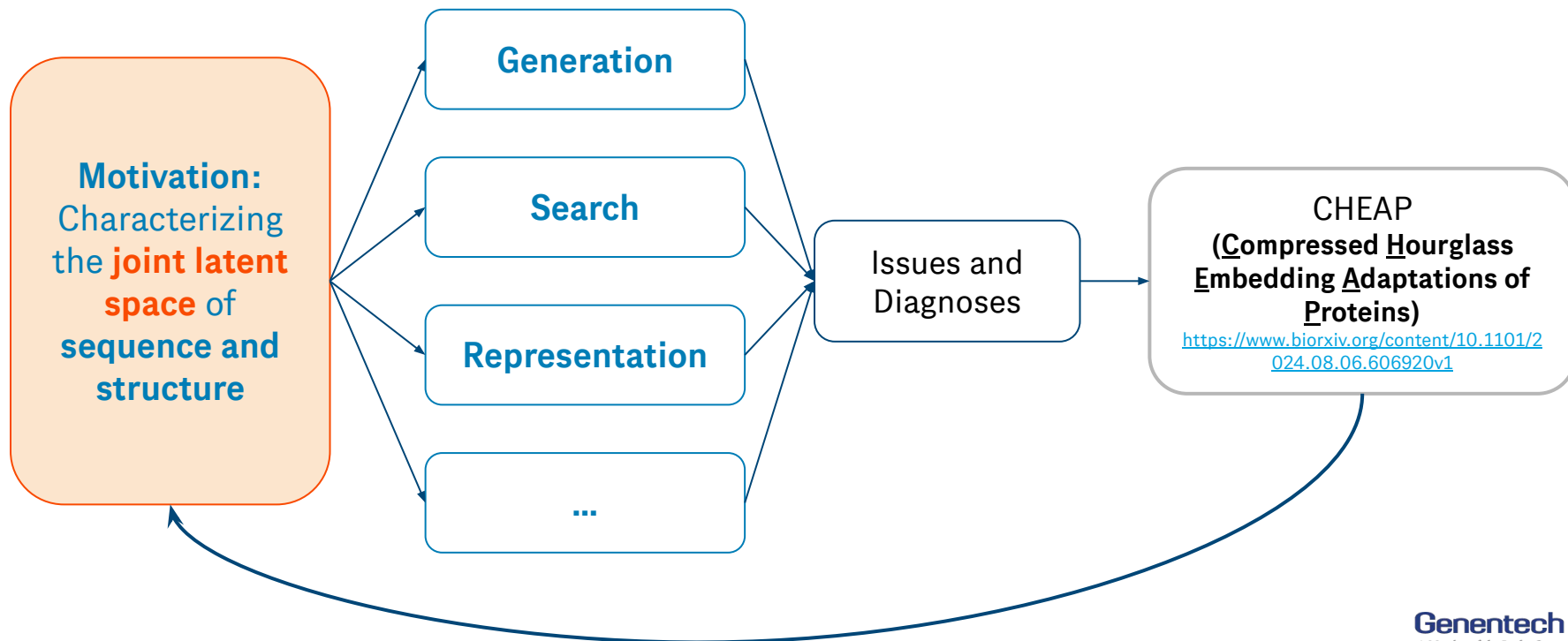
*Amy X. Lu  
UC Berkeley / BAIR  
Prescient Design / Genentech*

Paper: [bit.ly/cheap-proteins](https://bit.ly/cheap-proteins)  
GitHub: [github.com/amyxlu/cheap-proteins](https://github.com/amyxlu/cheap-proteins)

# Agenda

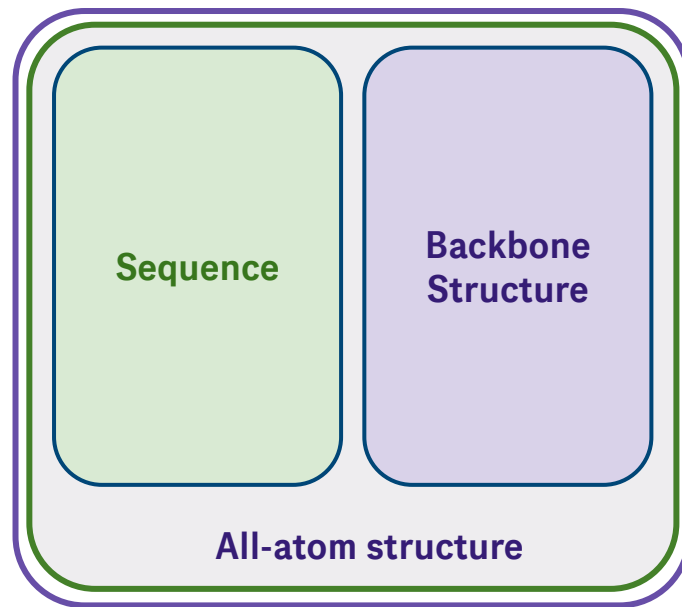


# Agenda



## Motivation: Obtaining a joint embedding of structure & sequence from sequence alone

- Existing protein representation models often capture either  $p(\text{sequence})$  or  $p(\text{structure})$ , limiting flexibility
- Desiderata:
  - Capture the joint embedding of **sequence** and **structure**
  - Can be explicitly decoded back to **structure** and **sequence**
  - Can be captured from **sequence** alone

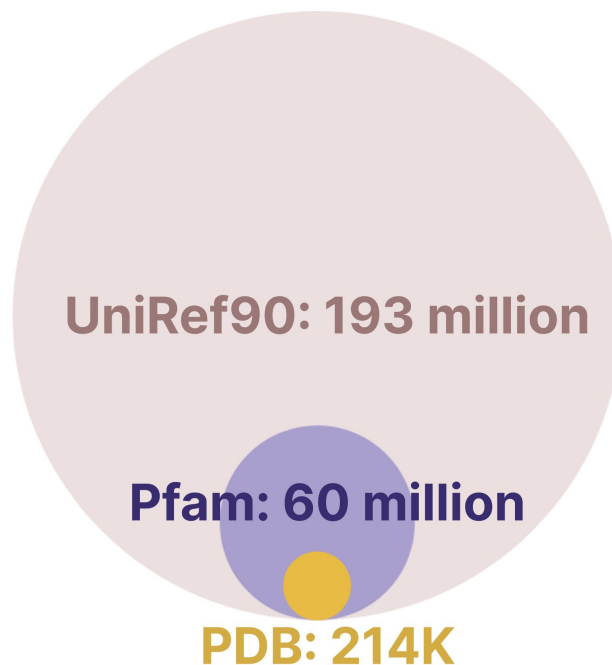


All-atom structure is a superset of sequence information!

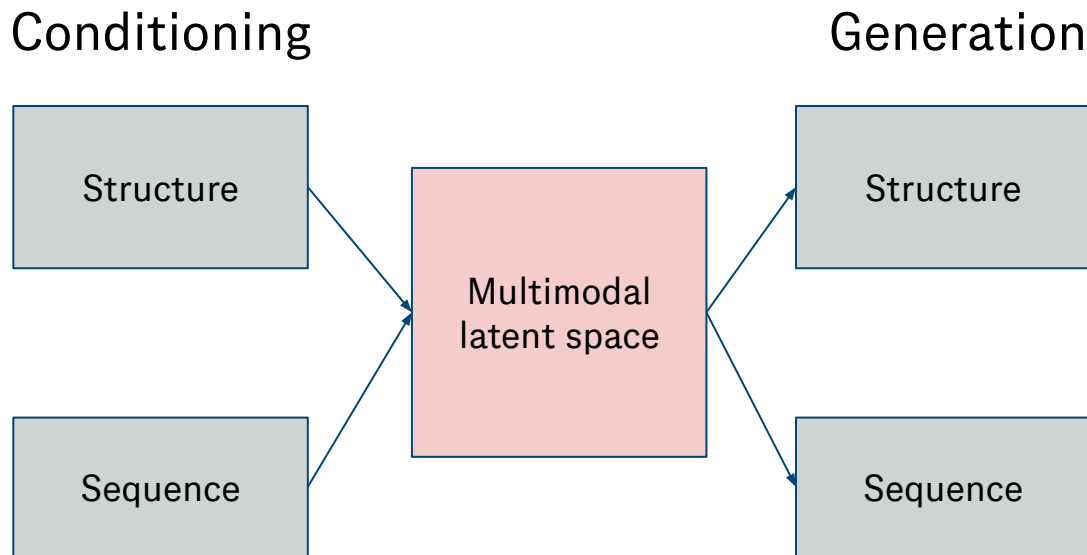


## Motivation: Sequence databases offer better data distribution coverage and function label abundance

- Structure databases have strong priors which may not always be useful:
  - biased towards crystallizable proteins
  - sequence database sizes approaches internet-scale data, while structure databases are much smaller



## Motivation: **Directly capturing the joint distribution is flexible**



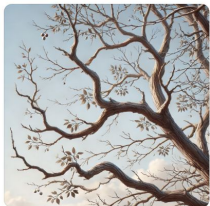
Being able to characterize a joint latent space allows flexibly conditioning by and generating either modality.

# Motivation: Direct sampling from the joint distribution is natural

## Structure generation + inverse folding

 $p(\text{branches})$ 

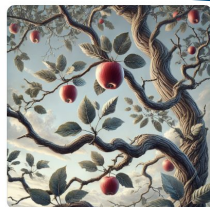
Generate a picture of tree branches.



Here is the image of tree branches you requested. If you need any modifications or additional details, feel free to ask!

 $p(\text{apples}|\text{branches})$ 

Add apples to this tree branch.

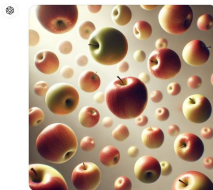


Here's the updated image with apples added to the tree branches. If you need any further adjustments, just let me know!

## Sequence generation + folding

 $p(\text{apples})$ 

Generate a picture of apples in midair.



Here is the image of apples suspended in midair. I hope it captures the surreal effect you were looking for!

 $p(\text{branches}|\text{apples})$ 

Add branches to these apples.



Here is the updated image with branches holding the apples in place. The branch touch while maintaining the surreal effect. I hope this meets your expectations!

## Co-generation

 $p(\text{apples, branches})$ 

Generate a picture of an apple tree.

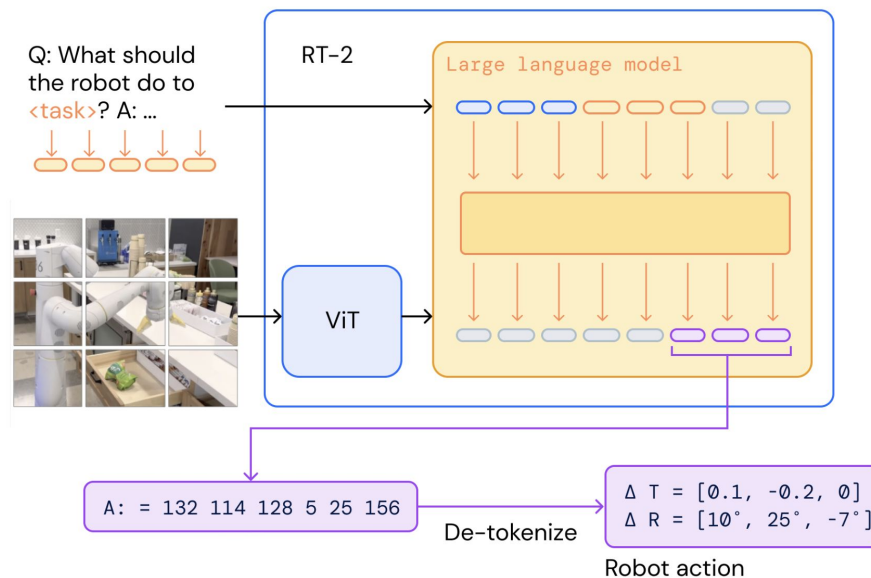


Here is the image of an apple tree you requested! Let me know if you need any adjustments or further assistance.

## Motivation: Large pretrained models capture useful priors for decision making

- Multimodal pretrained models offer useful priors
  - e.g. VLMs in robotics

→ can we use information captured by AlphaFold2, etc. as a “foundation model” for decision making in protein engineering?



[RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control](#)

---

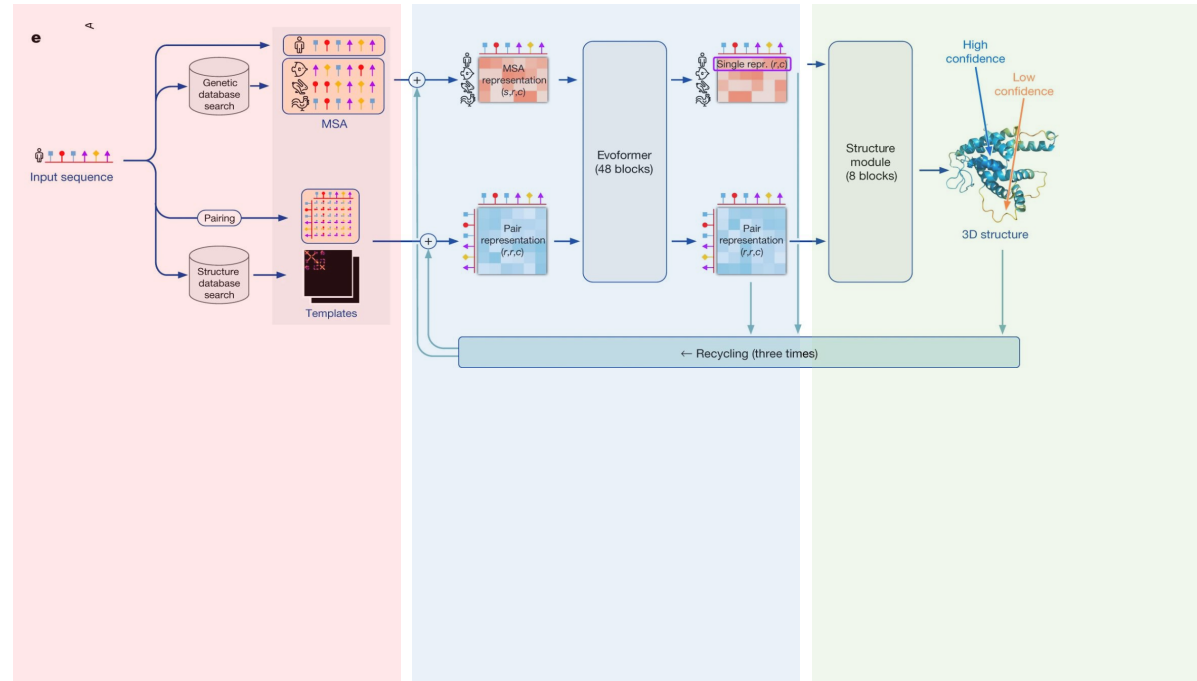
How can we repurpose the joint representation of  $p(\text{sequence, structure})$  in protein folding models for downstream tasks?

# Refresher: ESMFold for sequence-to-structure prediction



## AlphaFold2:

Uses an explicit retrieval step



harness additional  
sequence-based priors

learn structural features  
from sequence latents

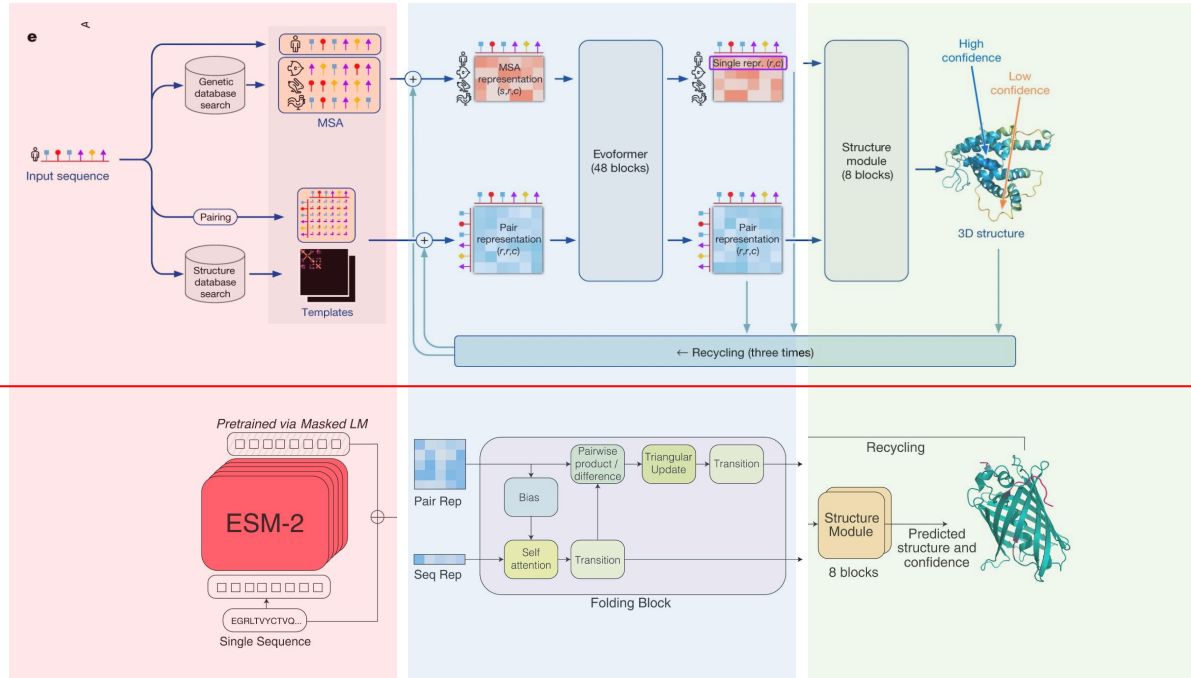
generate structures

# Refresher: ESMFold for sequence-to-structure prediction



## AlphaFold2:

Uses an explicit retrieval step



## ESMFold:

Replaces retrieval step with a **language model**

harness additional  
sequence-based priors

learn structural features  
from sequence latents

generate structures

```
esm / esm / esmfold / v1 / esmfold.py
Code Blame 364 lines (305 loc) · 13.6 KB
152     def forward(
185         # === ESM ===
186         esmaa = self._af2_idx_to_esm_idx(aa, mask)
187
188         if masking_pattern is not None:
189             esmaa = self._mask_inputs_to_esm(esmaa, masking_pattern)
190
191         esm_s, esm_z = self._compute_language_model_representations(esmaa)
192
193         # Convert esm_s to the precision used by the trunk and
194         # the structure module. These tensors may be a lower precision if, for example,
195         # we're running the language model in fp16 precision.
196         esm_s = esm_s.to(self.esm_s_combine.dtype)
197         esm_s = esm_s.detach()
198
199         # === preprocessing ===
200         esm_s = (self.esm_s_combine.softmax(0).unsqueeze(0) @ esm_s).squeeze(2)
201
202         s_s_0 = self.esm_s_mlp(esm_s)
203         if self.cfg.use_esm_attn_map:
204             esm_z = esm_z.to(self.esm_s_combine.dtype)
205             esm_z = esm_z.detach()
206             s_z_0 = self.esm_z_mlp(esm_z)
207         else:
208             s_z_0 = s_s_0.new_zeros(B, L, L, self.cfg.trunk.pairwise_state_dim)
209
210         s_s_0 += self.embedding(aa)
211
212         structure: dict = self.trunk(
213             s_s_0, s_z_0, aa, residx, mask, no_recycles=num_recycles
214         )
```



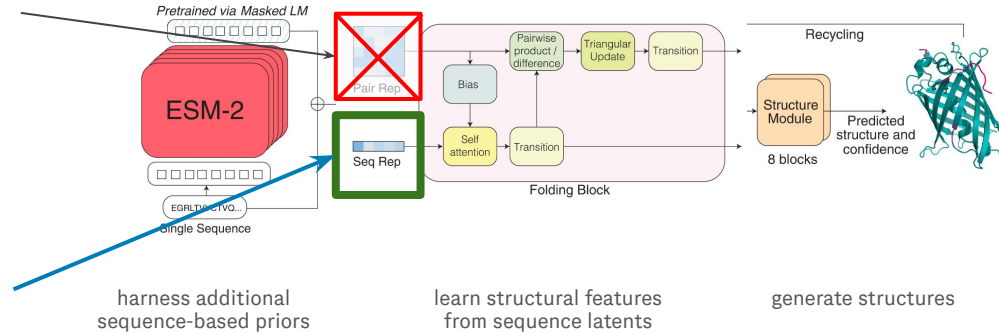
Observation: at inference time, the pairwise input is initialized as zeros...



Observation: at inference time, the pairwise input is initialized as zeros...

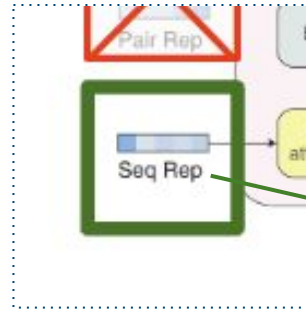


→ LM embedding captures sufficient inductive biases for structure, but requires **only sequence data** during training!

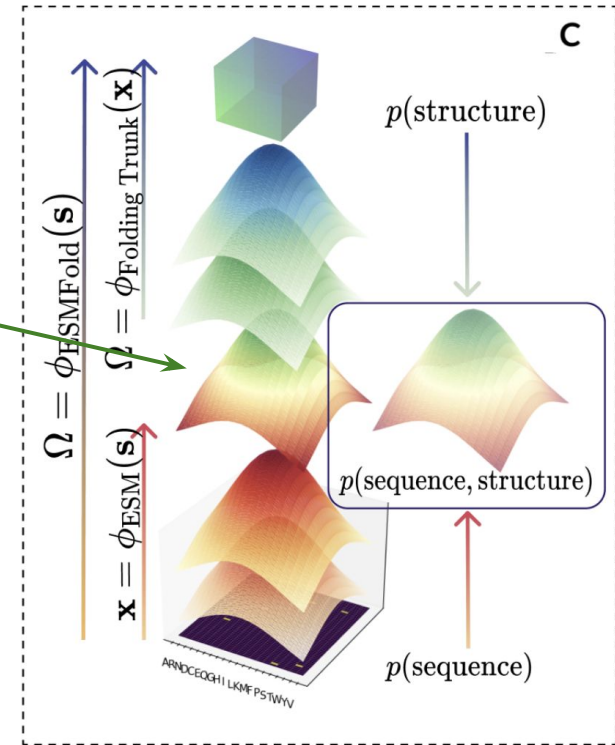


Observation: at inference time, the pairwise input is initialized as zeros...

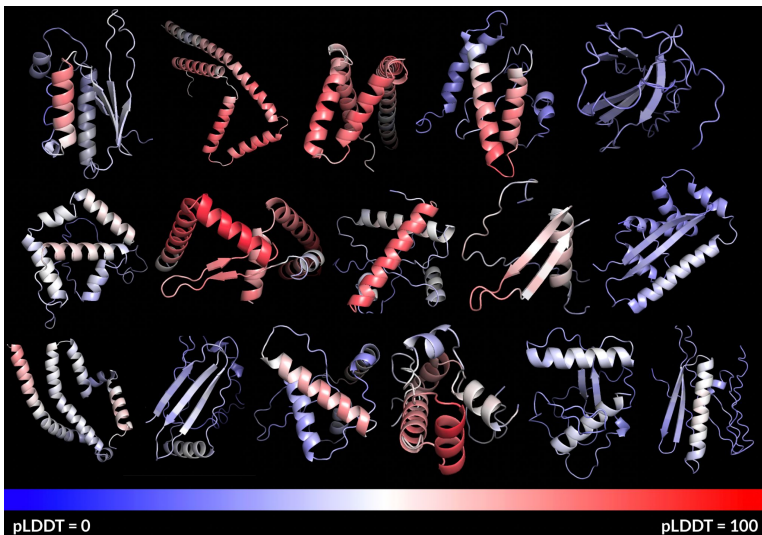
→ LM embedding captures sufficient inductive biases for structure, but requires **only sequence data** during training!



Consider this latent space as a **joint representation of protein sequence and structure that can be obtained from sequence only.**



## an early attempt at diffusing in this latent space...



We are able to learn structural folds, despite using only sequence inputs!

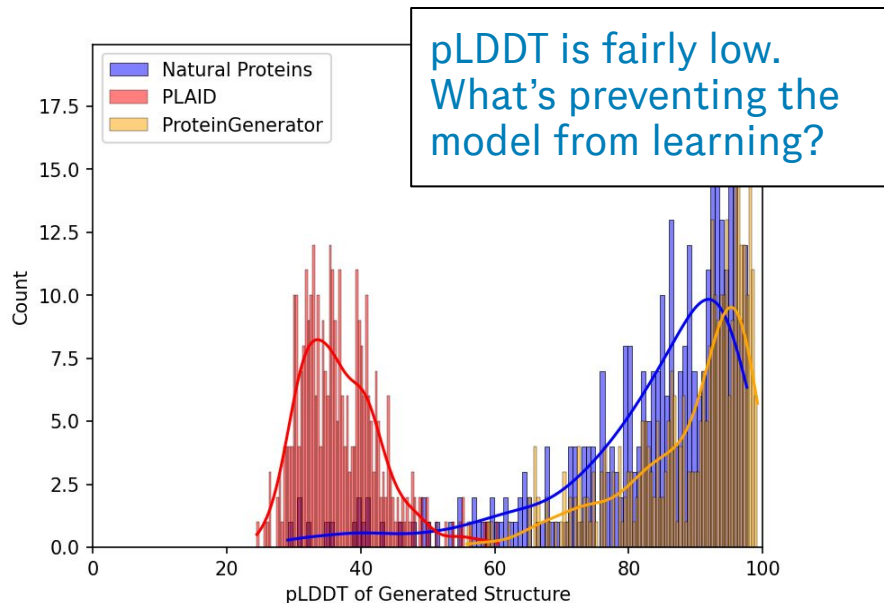
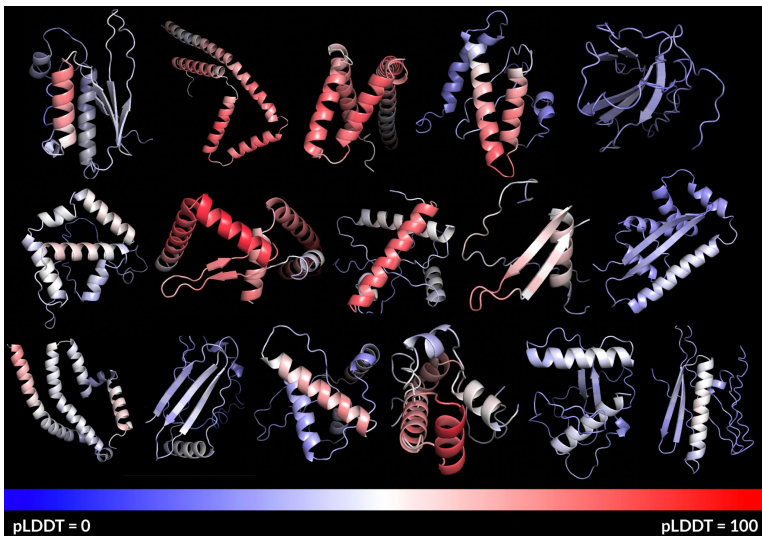
Empirically considering this latent space as a joint distribution is a go 

**PLAIID v0.5: Generating Protein Sequence and Structure Without Structural Training Data**

Amy X. Lu, Kevin K. Yang, Pieter Abbeel

*ICML 2024 Workshop on Machine Learning for Life and Material Sciences*

## an early attempt at diffusing in this latent space...



### PLAID v0.5: Generating Protein Sequence and Structure Without Structural Training Data

Amy X. Lu, Kevin K. Yang, Pieter Abbeel

## Issues and hypotheses -> CHEAP

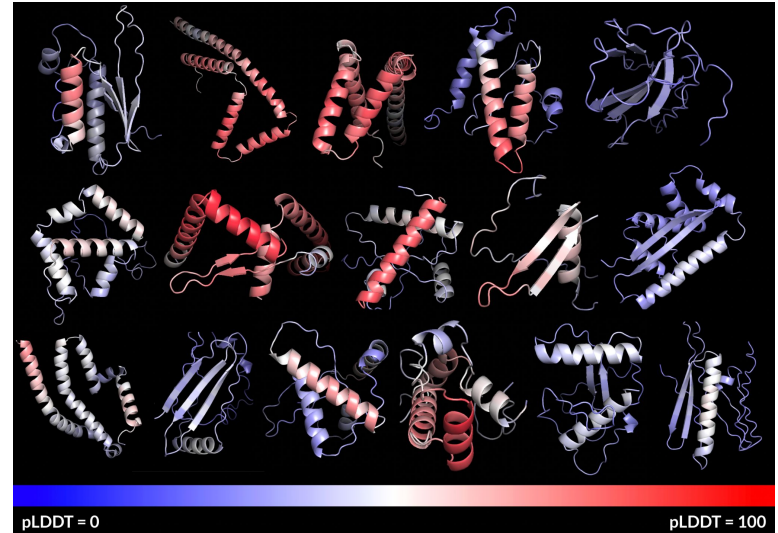
- Latent space requires regularization

In order to avoid arbitrarily high-variance latent spaces, we experiment with two different kinds of regularizations. The first variant, *KL-reg.*, imposes a slight KL-penalty towards a standard normal on the learned latent, similar to a VAE [46, 69], whereas *VQ-reg.* uses a vector quantization layer [96] within the decoder. This model can be interpreted as a VQGAN [23] but with the quantization layer absorbed by the decoder.

[High-Resolution Image Synthesis with Latent Diffusion Models](#)

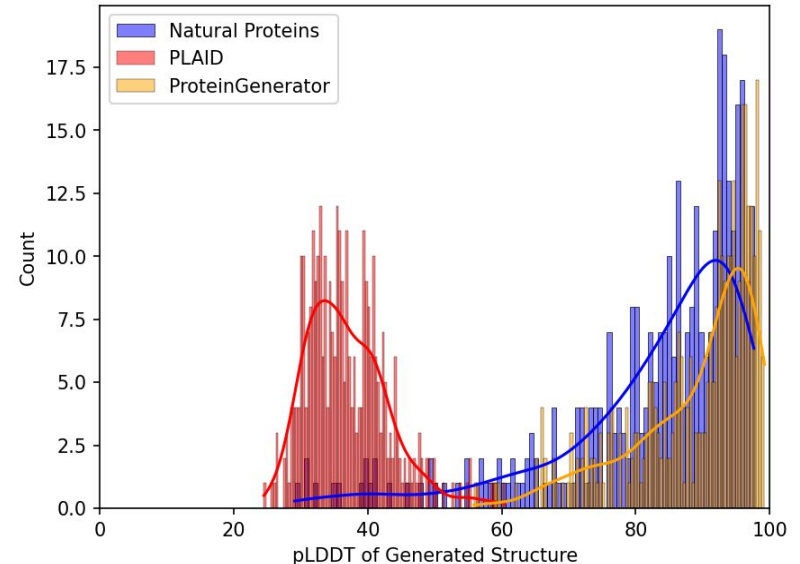
## Issues and hypotheses -> CHEAP

- Latent space requires regularization
- Training data only allows for length of 128 due to memory constraints
  - Some samples show the curvatures of a beta barrel, but sequence length limits seeing a full beta barrel



## Issues and hypotheses -> CHEAP

- Latent space requires regularization
- Training data only allows for length of 128 due to memory constraints
  - Some samples show the curvatures of a beta barrel, but sequence length limits seeing a full beta barrel
    - Need to shorten the protein?
- pLDDT is not designed to assess generation from evolutionary scale datasets
  - Biased towards generative models trained on the same data as AF2, i.e. PDB



## Issues and hypotheses -> CHEAP

- Latent space requires regularization
- Training data only allows for length of 128 due to memory constraints
  - Some samples show the curvatures of a beta barrel, but sequence length limits seeing a full beta barrel
    - Need to shorten the protein?
- pLDDT is not designed to assess generation from evolutionary scale datasets
  - Biased towards generative models trained on the same data as AF2, i.e. PDB
- Large latent space corresponds to high-resolution image generation
  - in LDMs, latent space is  $64 \times 4 \times 4$ , as opposed to ours, which is  $512 \times 1024$

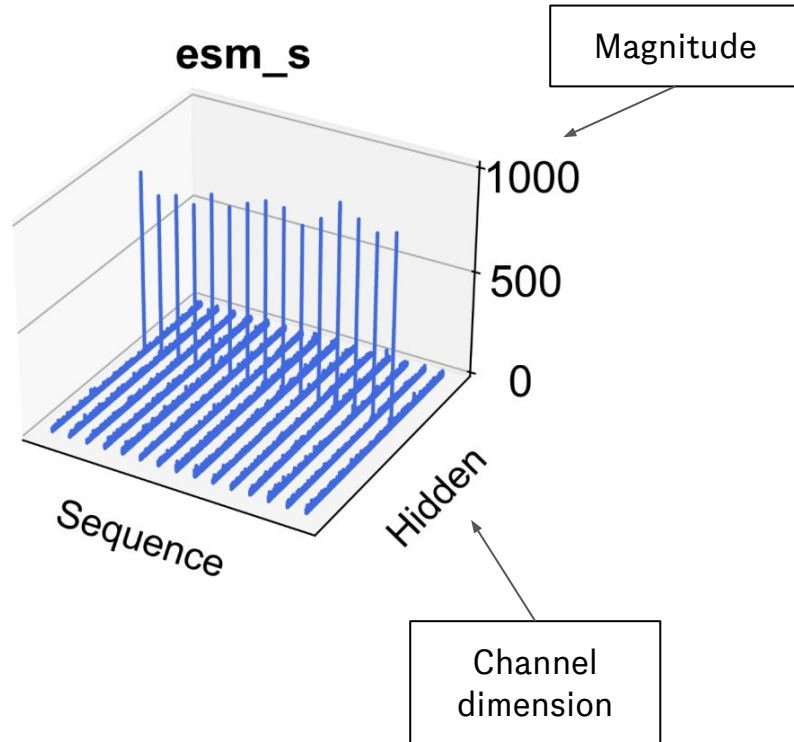
G. NCSN++ (Song et al., 2021) FFHQ-1024<sup>2</sup> Reference Samples



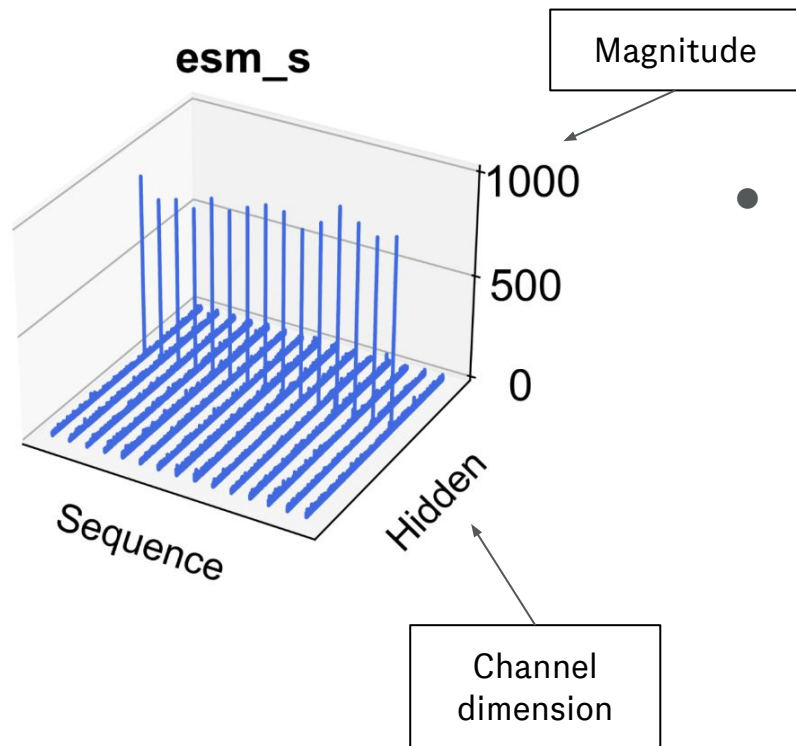
Diffusion models in their naive formulation often fail for  $1024 \times 1024$  resolution generation.



## A closer look at the latent space of ESMFold...

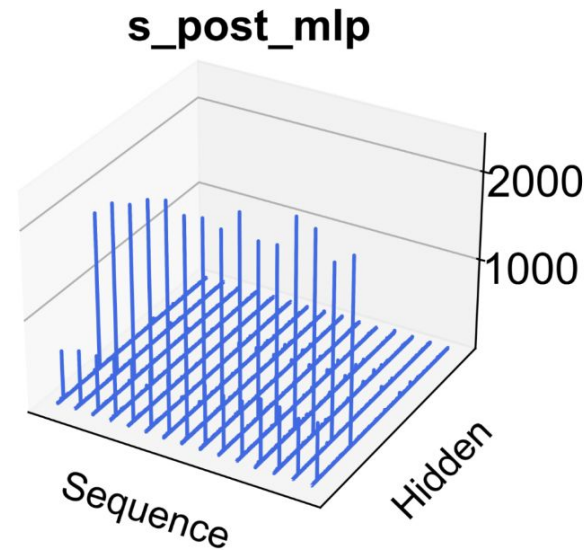
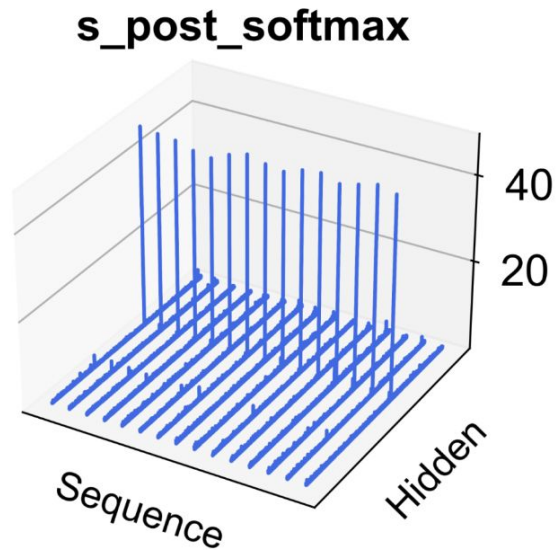
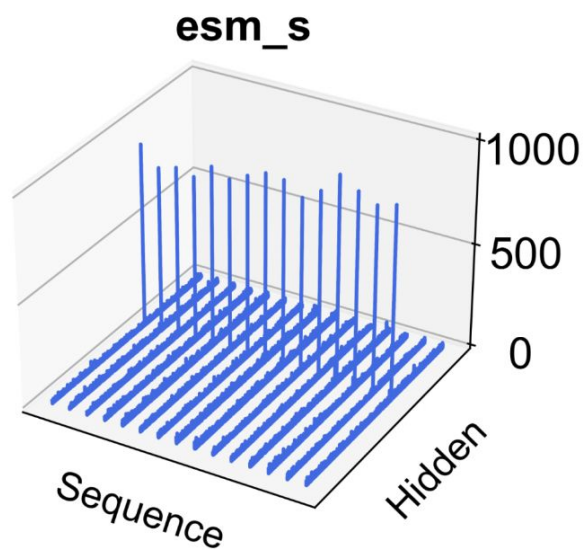


## ...ESMFold latent space exhibits pathologically large values



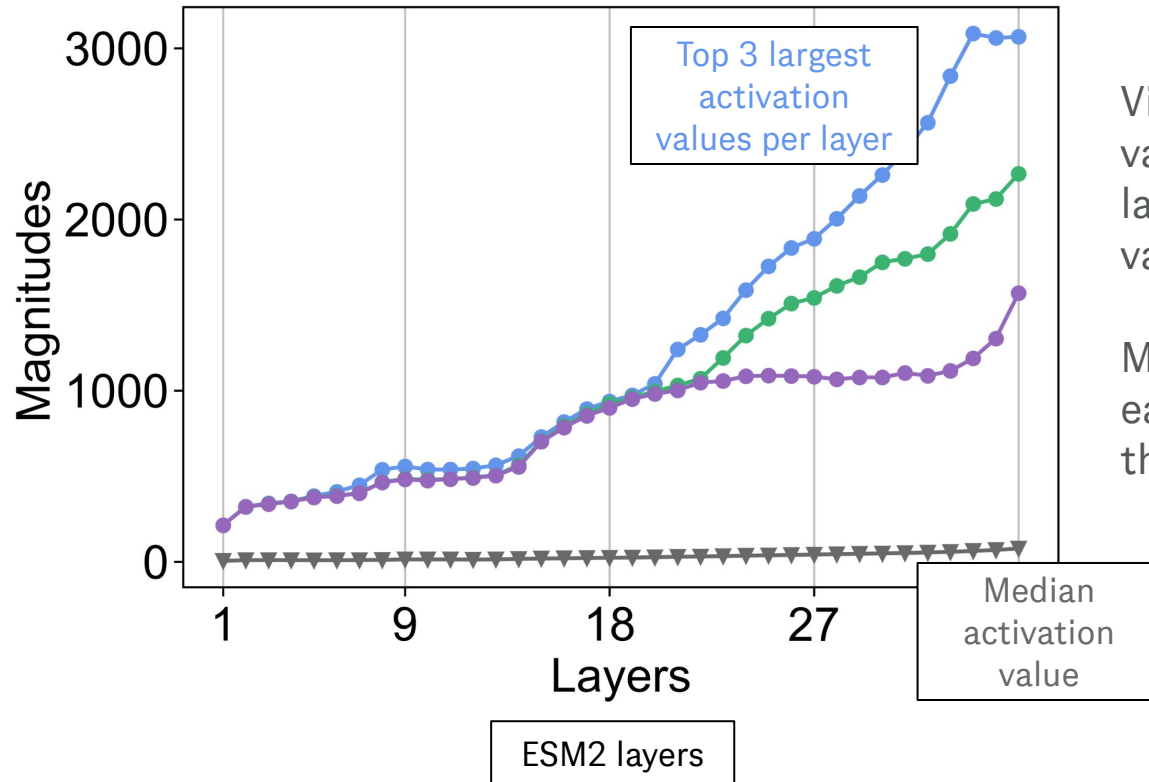
- Some channels exhibit very high mean values, regardless of the input.
  - Implications for generation: data distribution is no longer Gaussian distributed

## ...ESMFold latent space exhibits pathologically large values



Not just an issue for this particular layer...

## ~~ESMFold~~ ESM2 latent space exhibits pathologically large values



Visualizing the top 3 highest values in intermediate ESM2 layers, against the median value.

Massive activations begin in early layers, and accumulate throughout the model.

# ~~ESMFold~~ Large transformer model latent spaces exhibits pathologically large values

A pervasive issue across large transformer models!

*[Submitted on 27 Feb 2024 (v1), last revised 14 Aug 2024 (this version, v2)]*

## **Massive Activations in Large Language Models**

[Mingjie Sun](#), [Xinlei Chen](#), [J. Zico Kolter](#), [Zhuang Liu](#)

We observe an empirical phenomenon in Large Language Models (LLMs) -- very few activations exhibit significantly larger values than others (e.g., 100,000 times larger). We call them massive activations. First, we demonstrate the widespread existence of

# ESMFold Large transformer model latent spaces exhibits pathologically large values

A pervasive issue across large transformer models!

[Submitted on 27 Feb 2024 (v1), last revised 14 Aug 2024 (this version, v2)]

## Massive Activations in Large Language Models

Mingjie Sun, Xinlei Chen, J. Zico Kolter, Zhuang Liu

We observe an empirical phenomenon in Large Language Models (LLMs) -- very few activations exhibit significantly larger values than others (e.g., 100,000 times larger). We call them massive activations. First, we demonstrate the widespread existence of

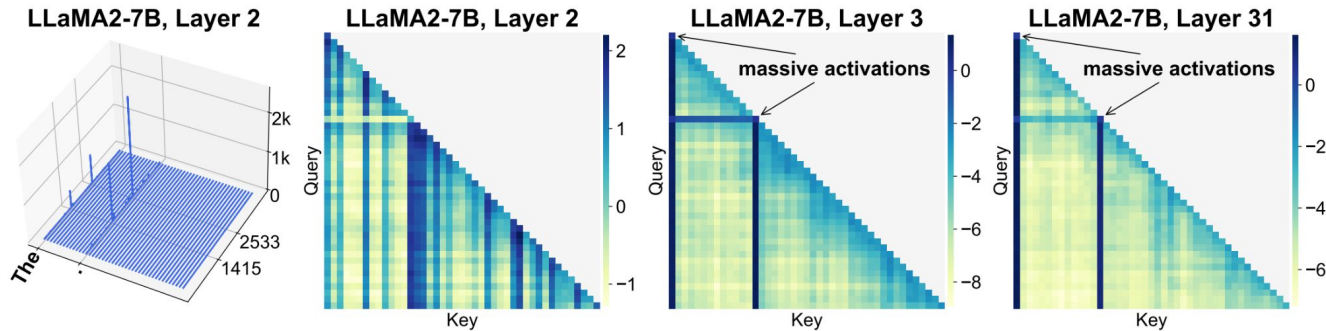
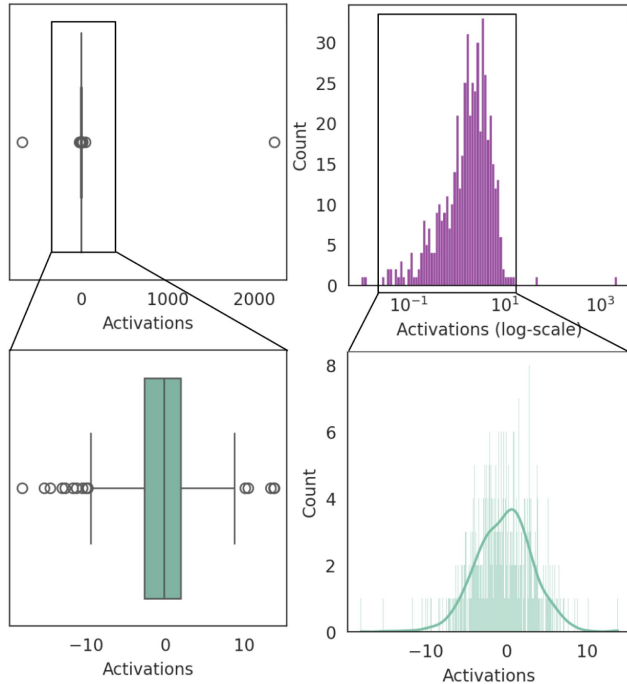
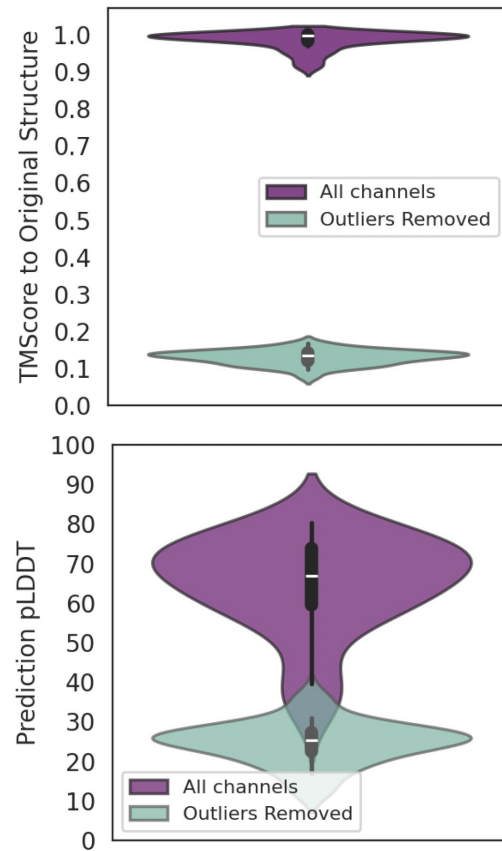
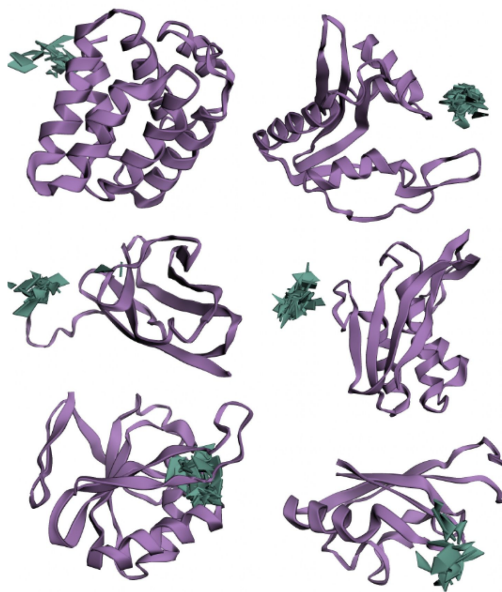
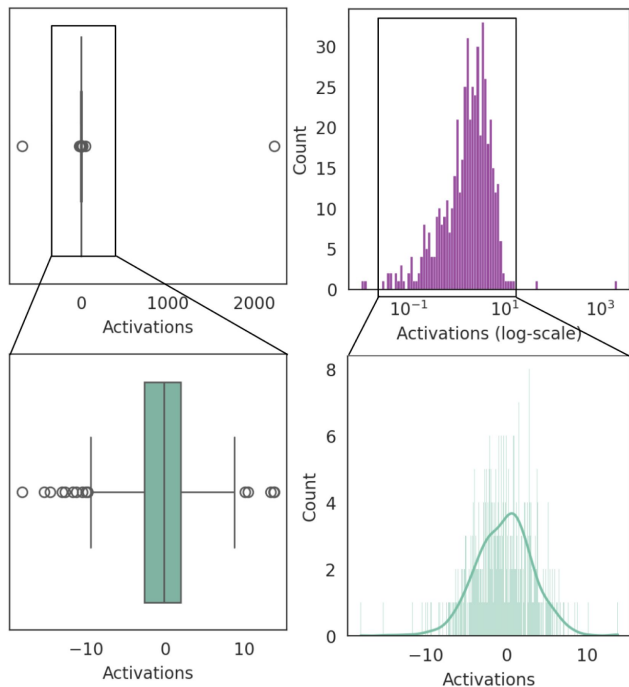


Figure 5: Attention patterns *before* and *after* massive activations appear in LLaMA2-7B. For each layer, we visualize average attention logits (unnormalized scores before softmax) over all heads, for an input sequence.

## What if we just remove these wacky channels?

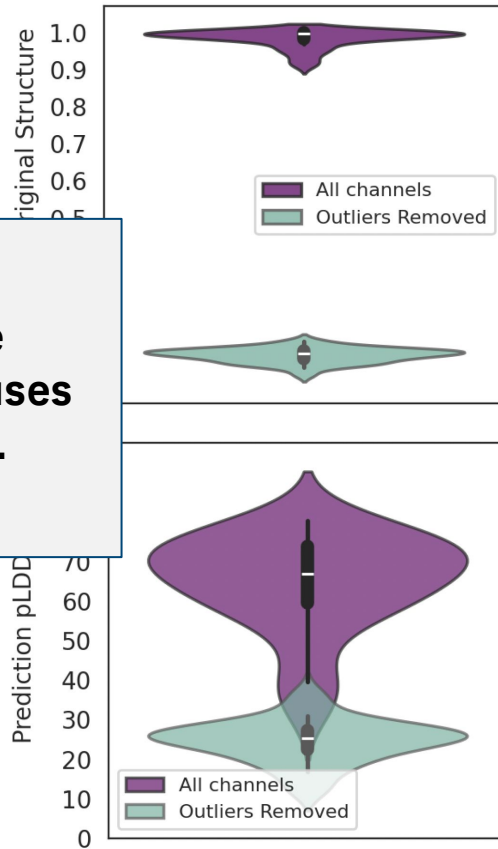
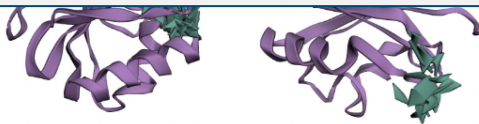
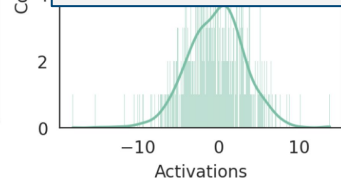
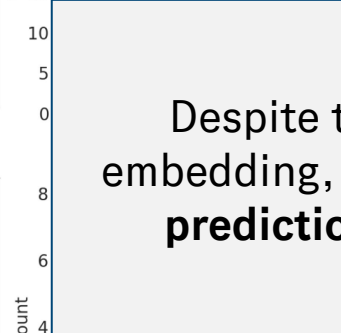
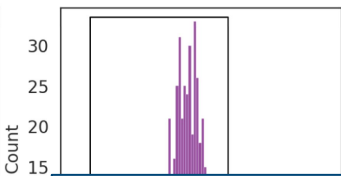
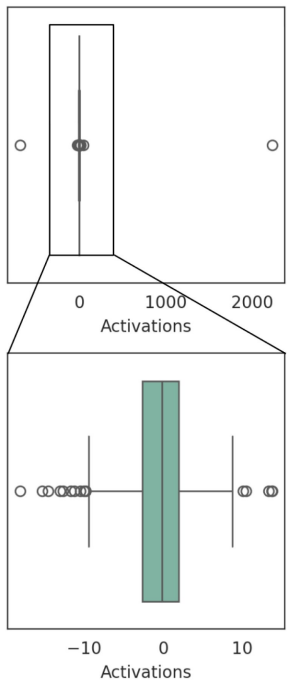


## What if we just remove these wacky channels?



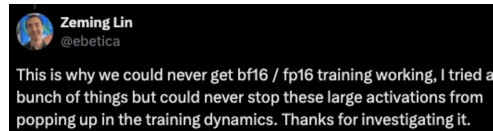


## What if we just remove these wacky channels?



## Why should we care about these massive activations?

- Training stability
- Model compression and 8-bit quantization
- Model interpretability
- ...



---

LLM.int8(): 8-bit Matrix Multiplication  
for Transformers at Scale

---

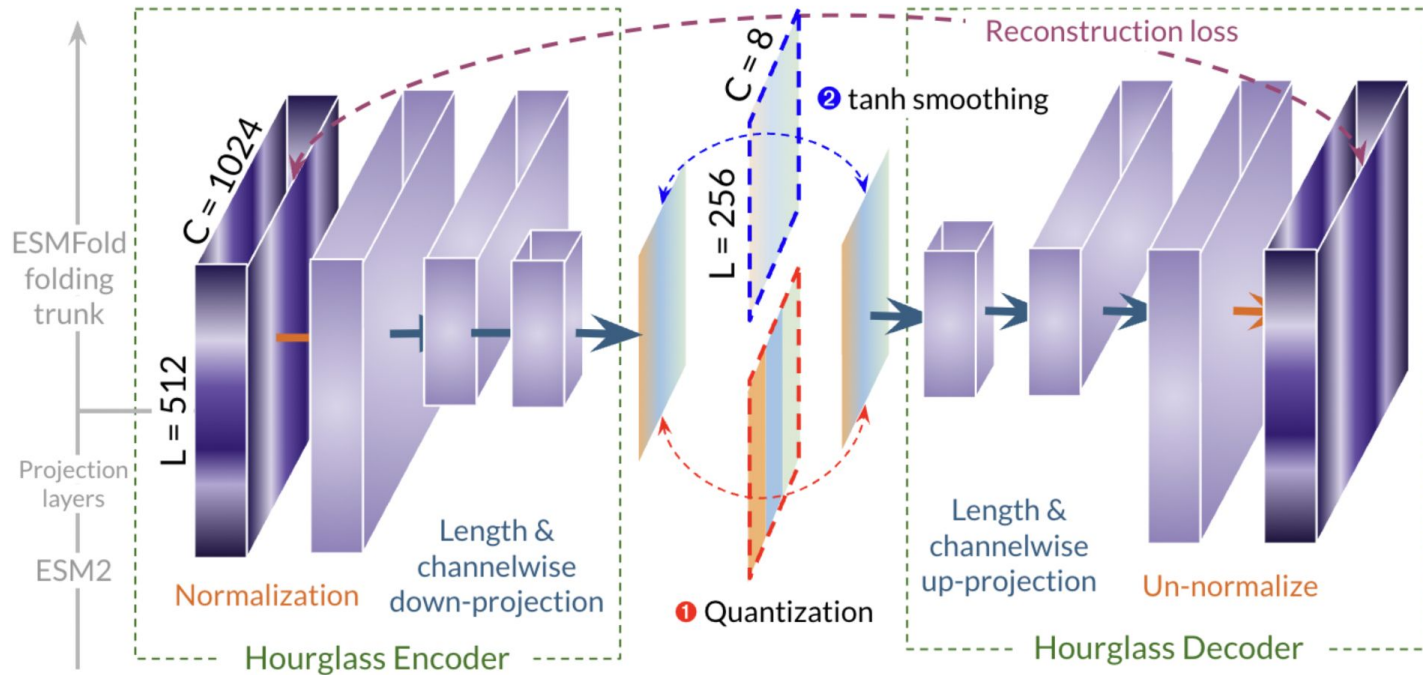
If removing 3 channels can remove performance, is the information evenly distributed through all the channels?

If not, can we **compress these channels?**

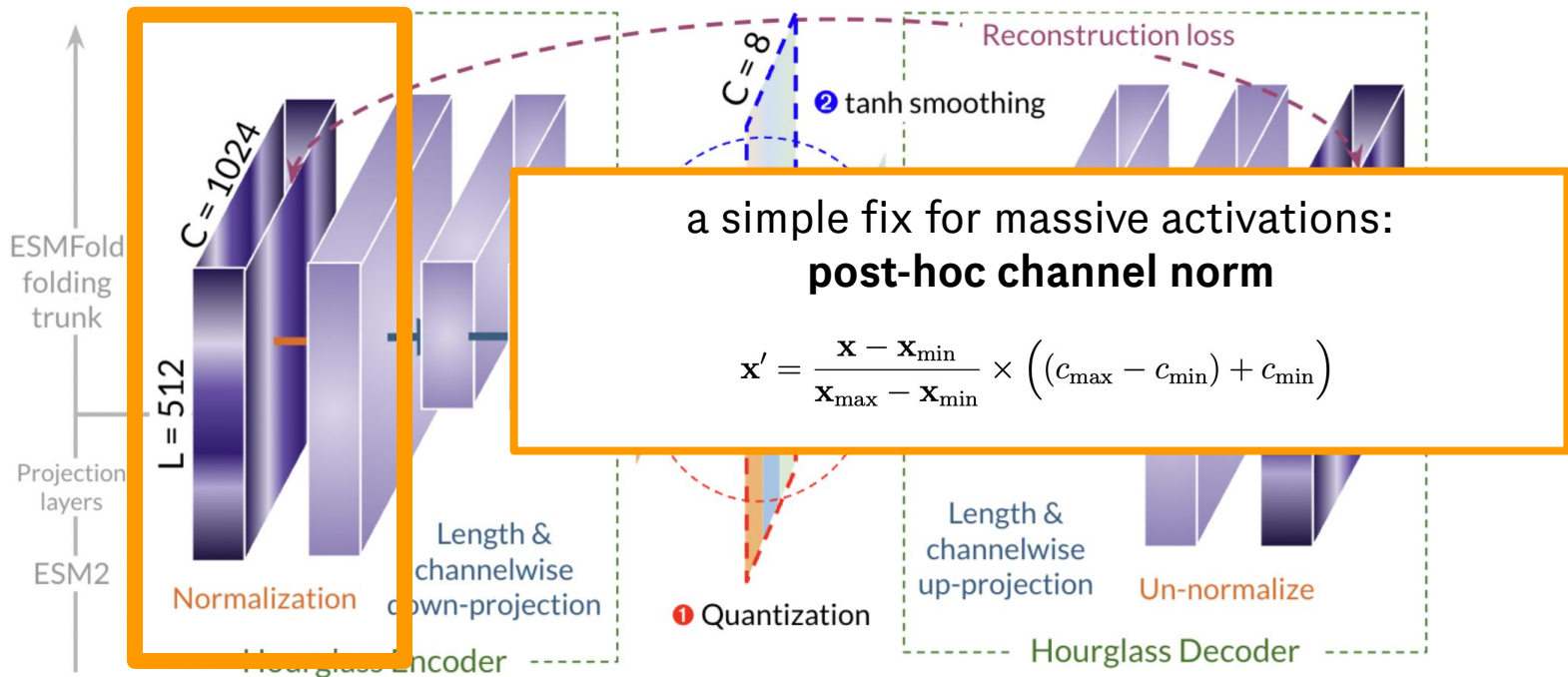
## Why compress?

- More portable representation
- Better understanding of protein folding internals
- Compressed data distributions are easier to learn during generative modeling

# An autoencoder for protein embedding compression



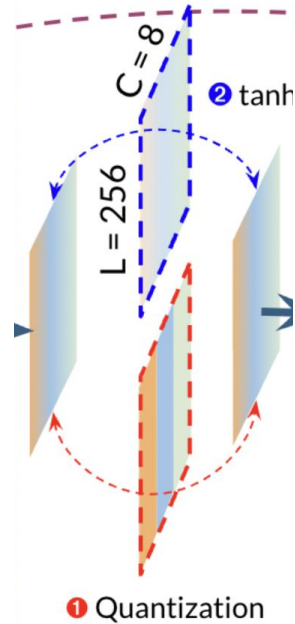
# An autoencoder for protein embedding compression



# Obtaining CHEAP embeddings

## 1. Tokenized

- Discretize embeddings using FSQ
  - 'snaps' continuous encoder values to discrete bins

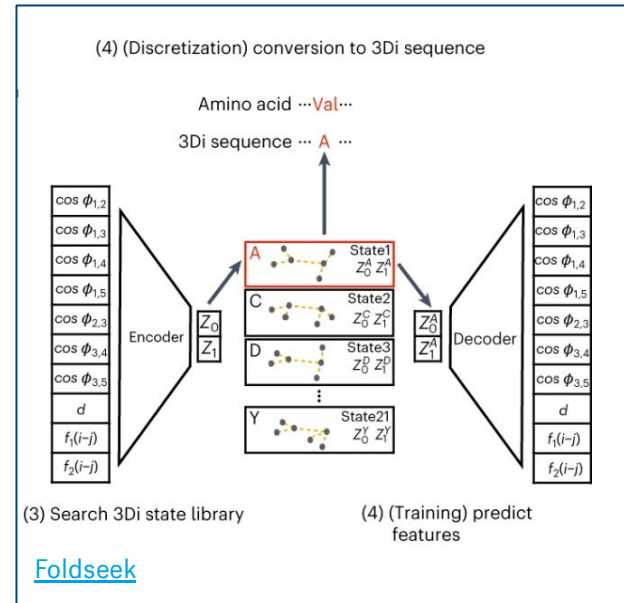
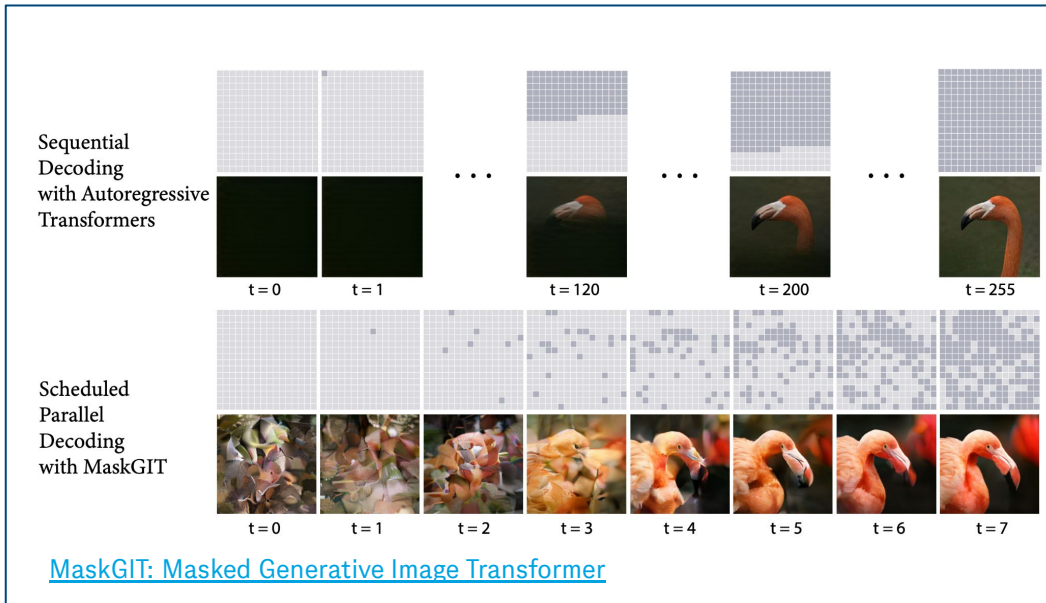


## 2. Continuous

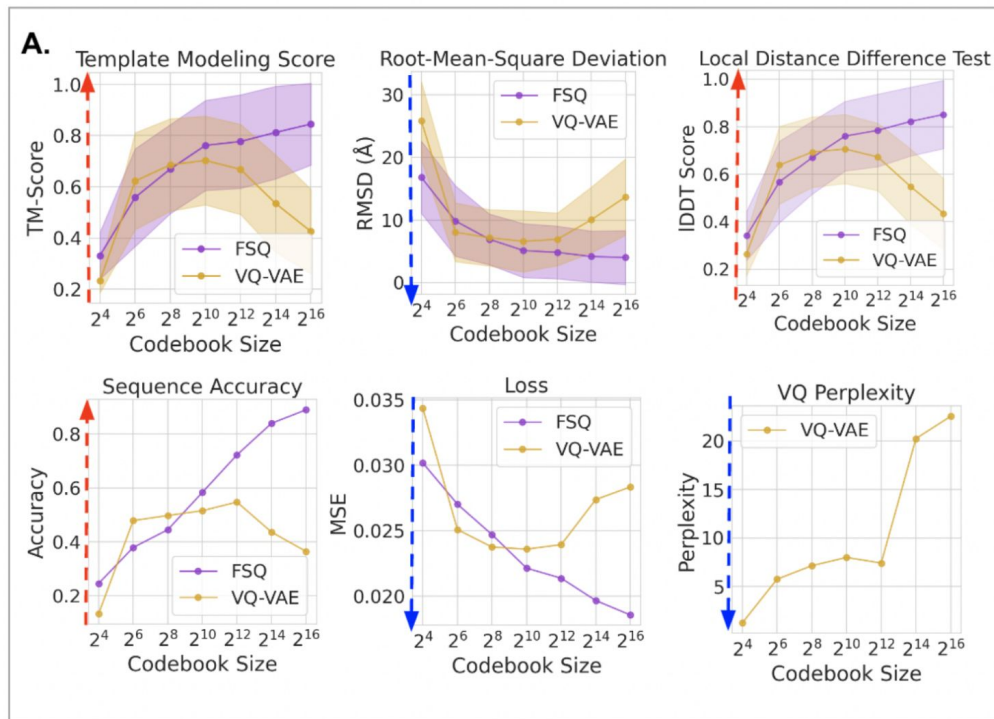
- Take the output of the downprojecting autoencoder
  - apply tanh to bound values between  $[-1, 1]$ , to bound values during diffusion

## Side note: why tokenized representations?

Tokenized representations can be helpful for our downstream aims of generation and search:

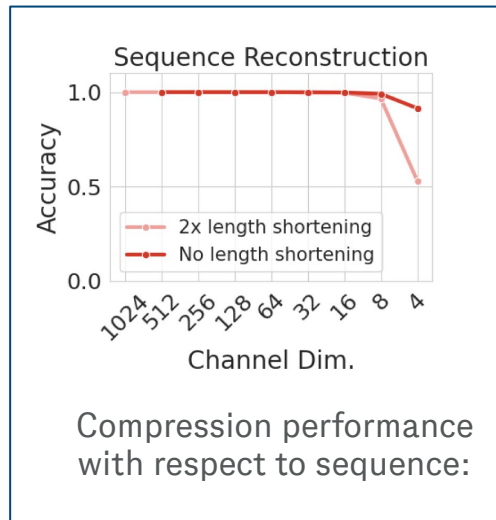


# All-atom structural tokenizer, obtained from sequence alone



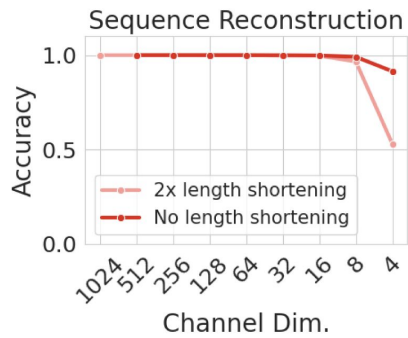


## ..yes, we *can* compress the embeddings:

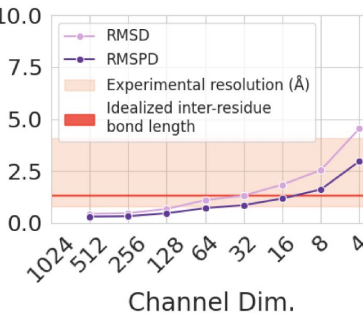
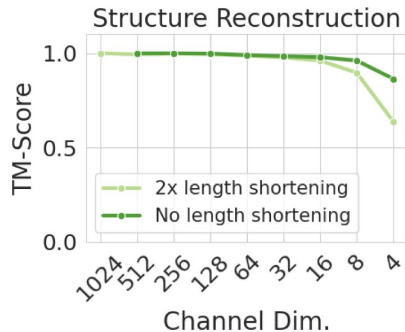


We can compress up to 8x, and sacrifice very little performance.

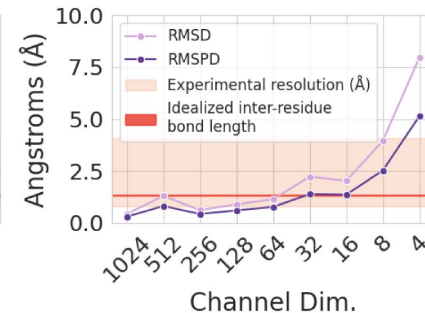
## ..yes, we *can* compress the embeddings:



Compression performance with respect to sequence:

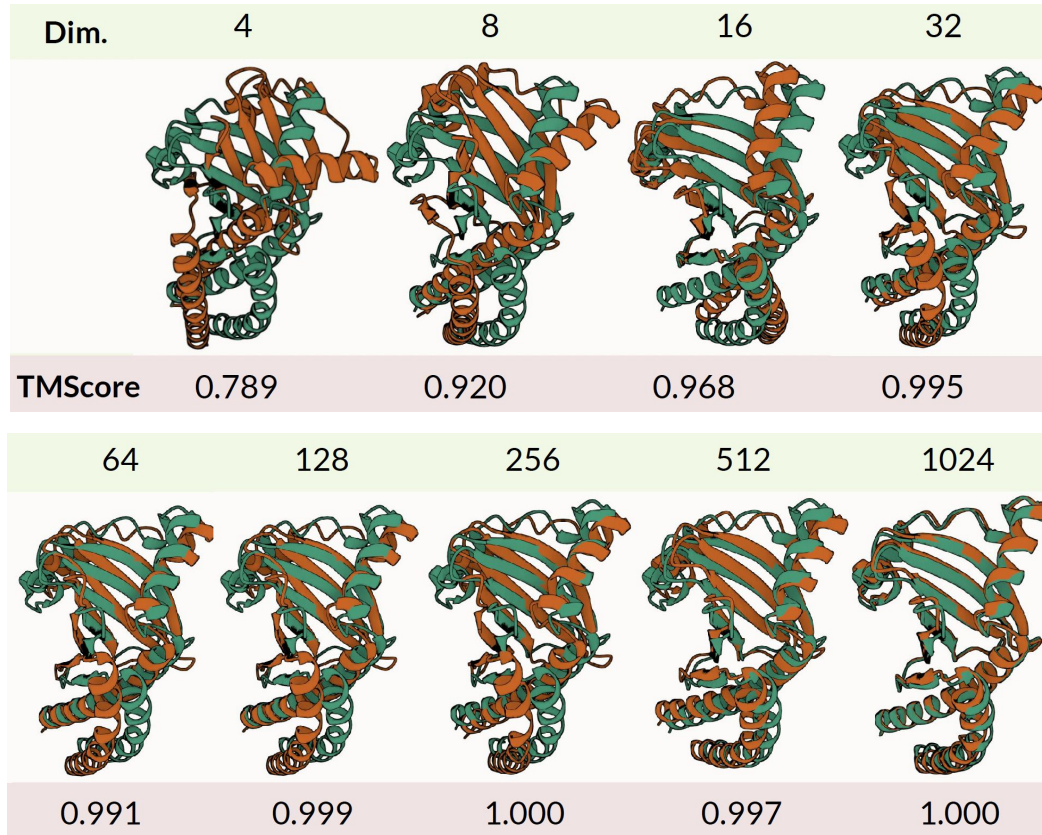


Compression performance with respect to structure:

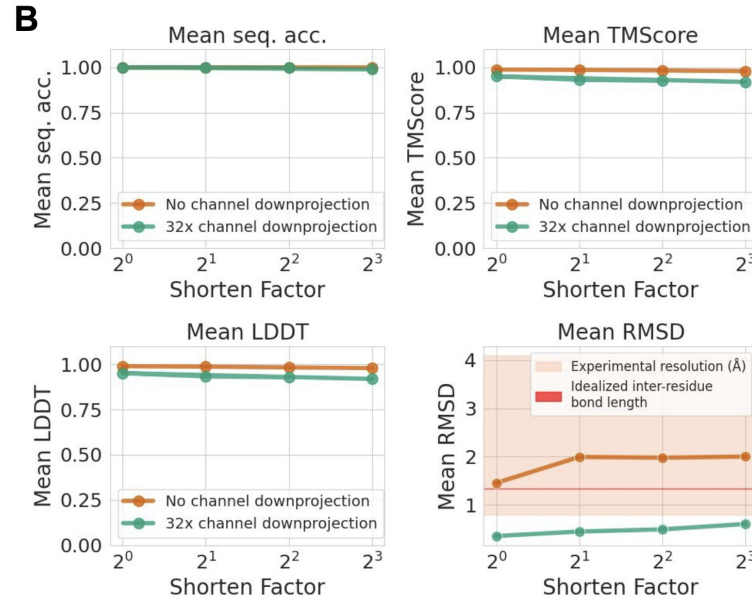


Sequence information is easier to retain than structure.

## ..yes, we *can* compress the embeddings:

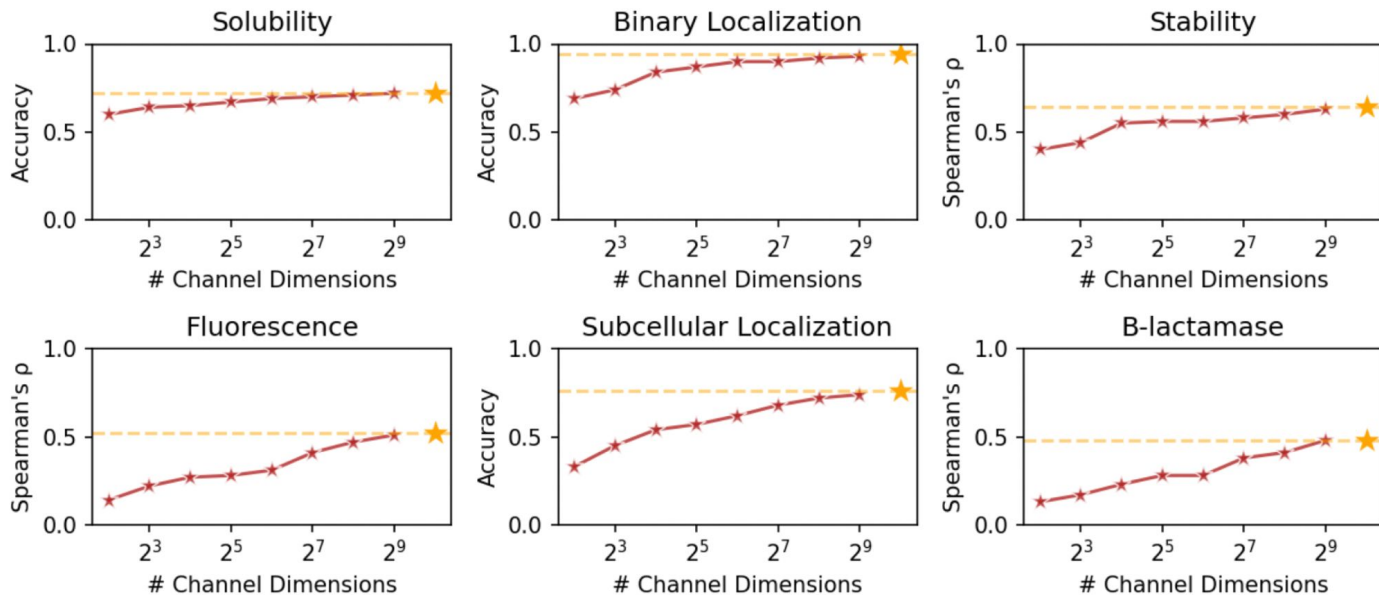


# We can compress lengthwise and channelwise:



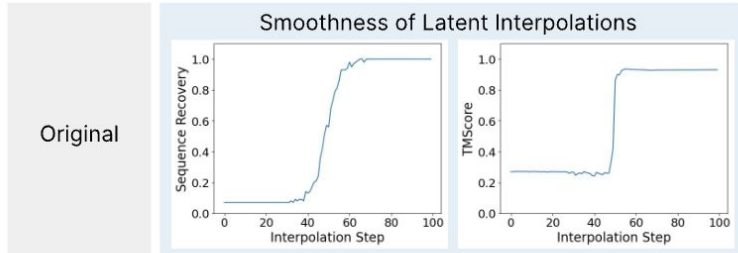
What does this mean for how structural information is shared across residue positions?

## What about function information?



Performance degradation with compression is much more gradual.  
What does this imply about the information content captured in pLMs  
with respect to downstream tasks?

# Does the autoencoding scheme “fix” the irregular latent space?

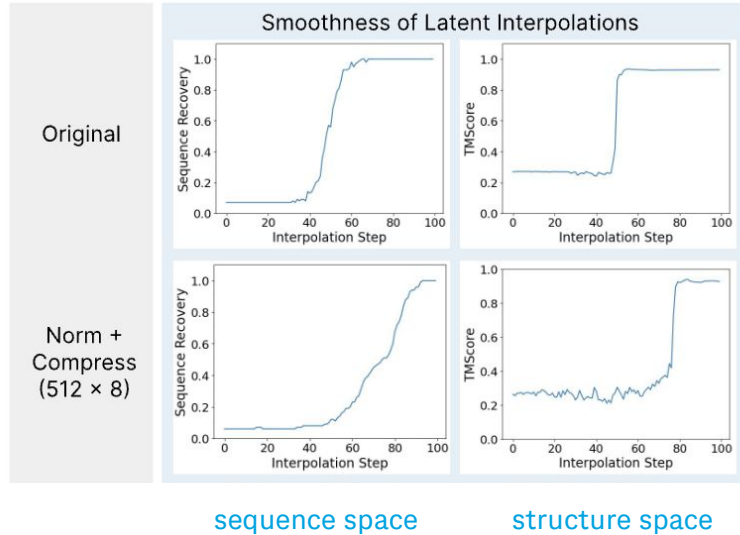


- Despite linearly interpolating in the latent space, the decoded sequence and structure changes very abruptly.

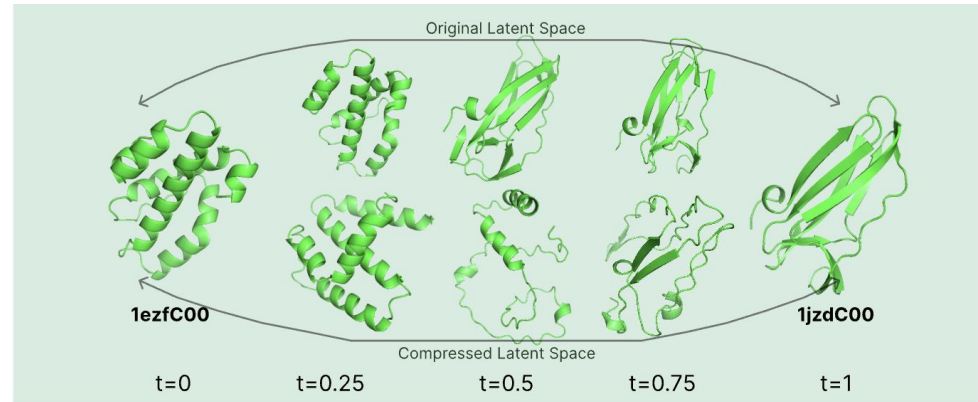
sequence space

structure space

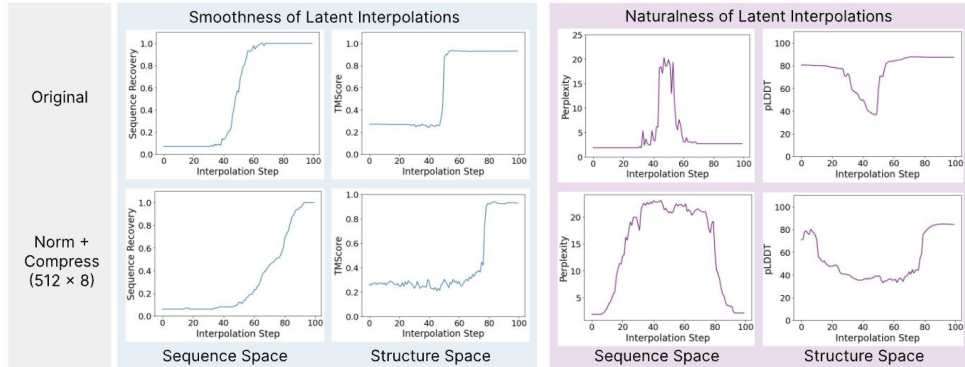
# Does the autoencoding scheme “fix” the irregular latent space?



- Despite linearly interpolating in the latent space, the decoded sequence and structure changes very abruptly.
- After CHEAP regularization, the change is more gradual



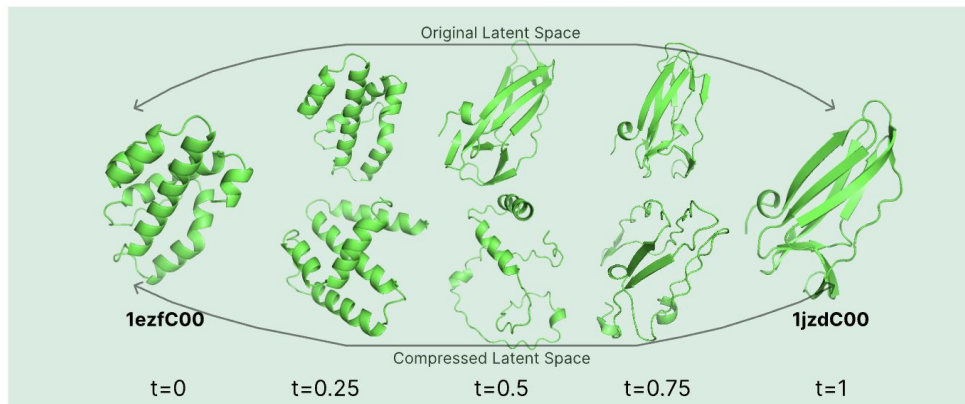
# PLM latent manifolds might be less “rugged” than true protein fitness landscapes



What makes for a good latent space?

Should we want more of the latent space to map back to a “valid protein” for sampling purposes, or properly model the rugged protein landscape?

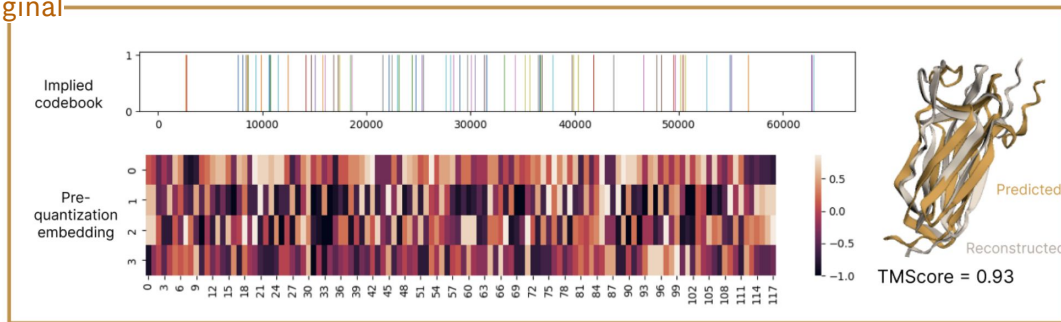
Do current PLM embeddings actually recapitulate protein fitness landscapes?



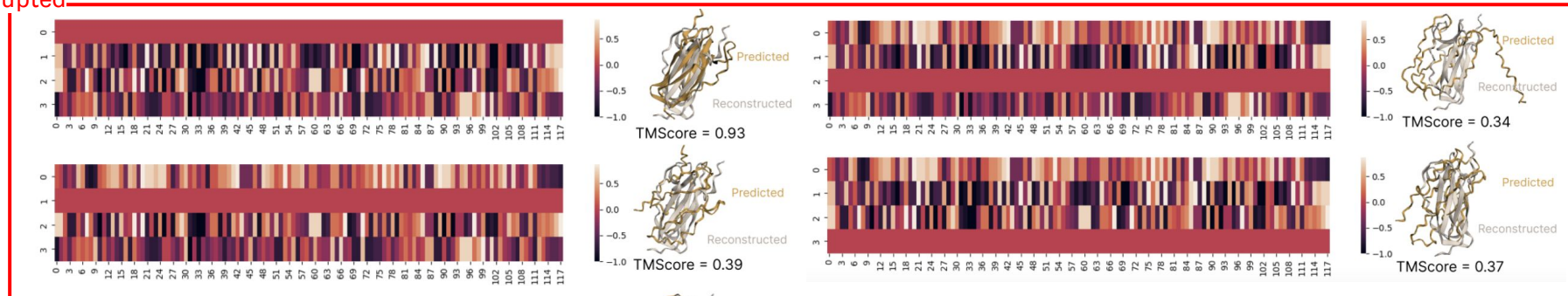


# “Disrupting” and reconstructing in the token space

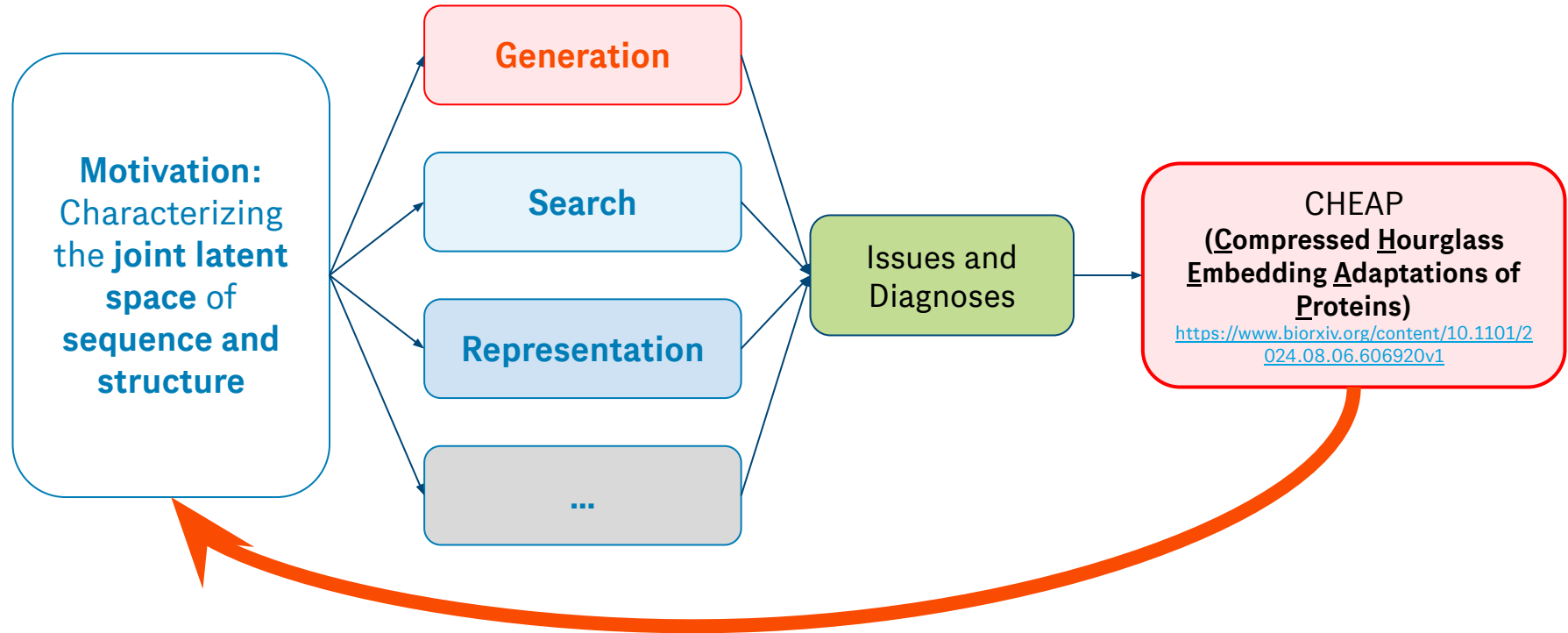
original



corrupted



# Agenda



---

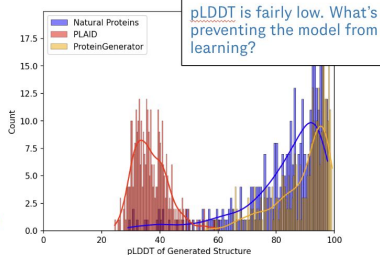
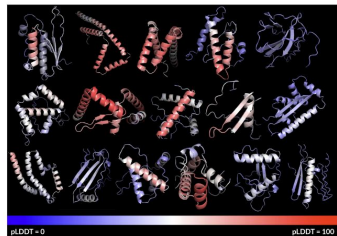
# PLAID (Protein LAtent Induced Diffusion)

*ongoing work!*

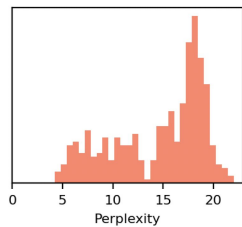
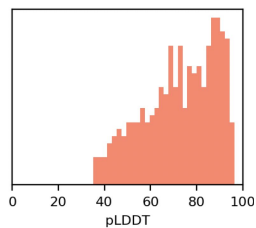
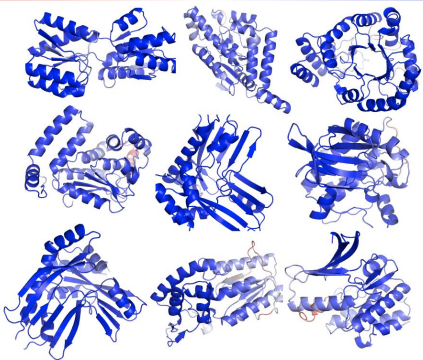
*tl;dr – now that we have a regularized & compressed embedding of  $p(\text{sequence, structure})$ , can we train a latent diffusion model for co-generation?*

# PLAID, again

an early attempt at diffusing in this latent space...



pLDDT < 50    50 < pLDDT < 70    70 < pLDDT < 90    pLDDT > 90



- Learn diffusion model in regularized and compressed latent space
  - mirrors the regularized autoencoder in LDM
- Can learn on longer sequences due to CHEAP shortening
- Use DiT instead of U-triangular self attention
  - allows for scaling up to higher parameter counts
- Scale up to 2B parameters with BS=2048

# PLAID, again

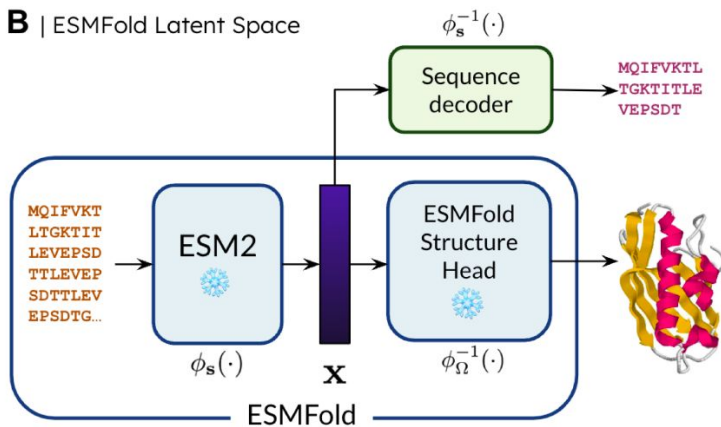
## A | Database Comparison

UniRef90: 193 million

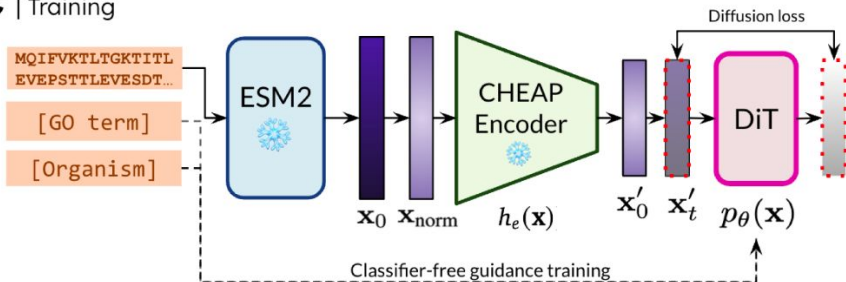
Pfam: 60 million

PDB: 214K

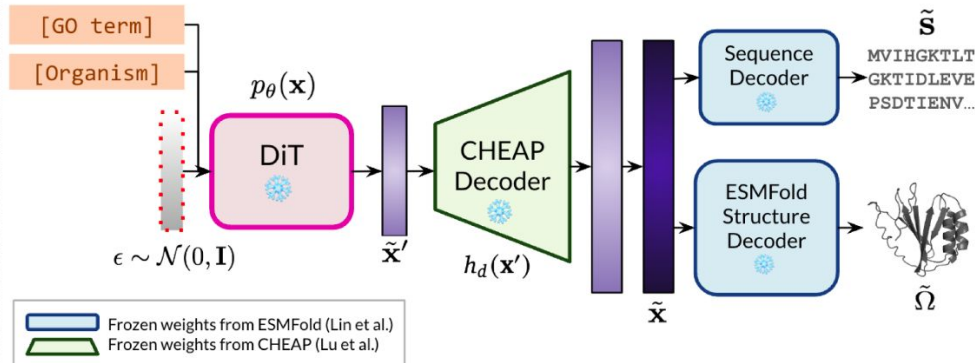
## B | ESMFold Latent Space



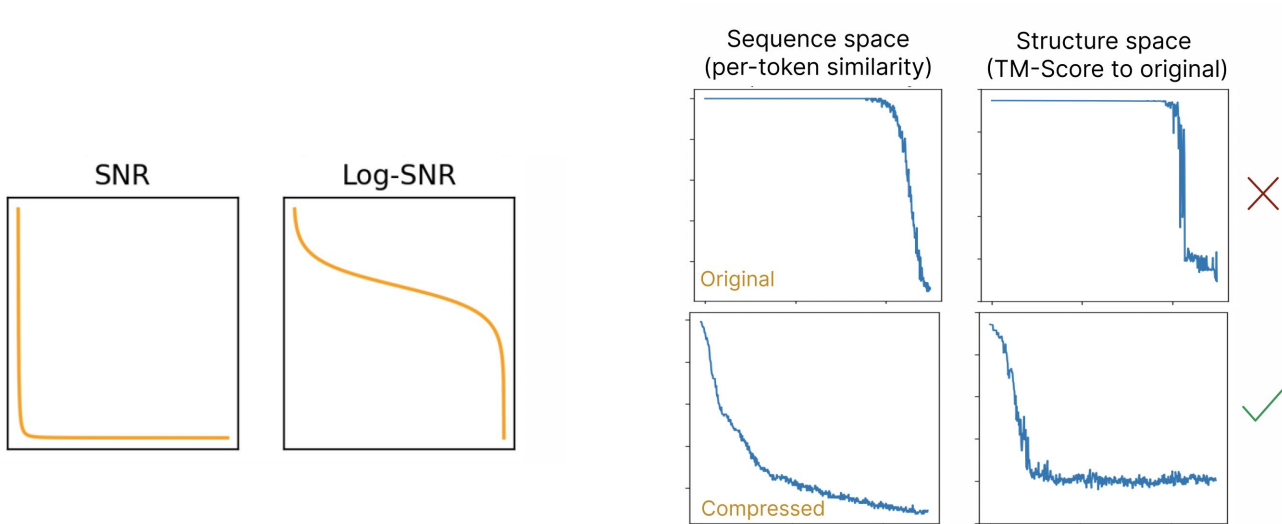
## C | Training



## D | Inference



## Comparing noise schedules in original and compressed latent space:

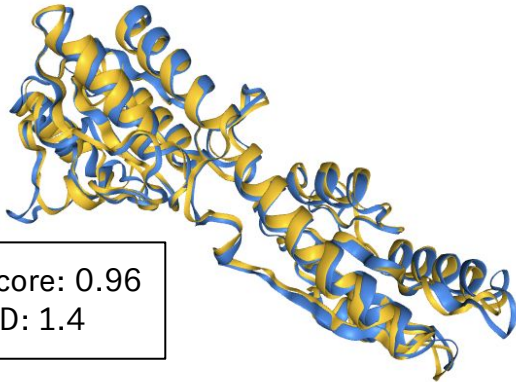


Noising in the CHEAP compressed space maps to noise in the sequence and structure space that is closer to the true signal-to-noise ratio.

# Samples demonstrate sequence and structural conservation

prompt: “yeast” AND “6-phosphofructokinase activity”

Search against the **structure database (PDB100)** to see if our samples are sensible...



- closest match: 3o8o [**Structure of phosphofructokinase**]
- organism: **Saccharomyces cerevisiae** (i.e. yeast)
- Sequence identity: 47.9%

Search against the **sequence database (UniRef90)** to see if our samples are sensible...

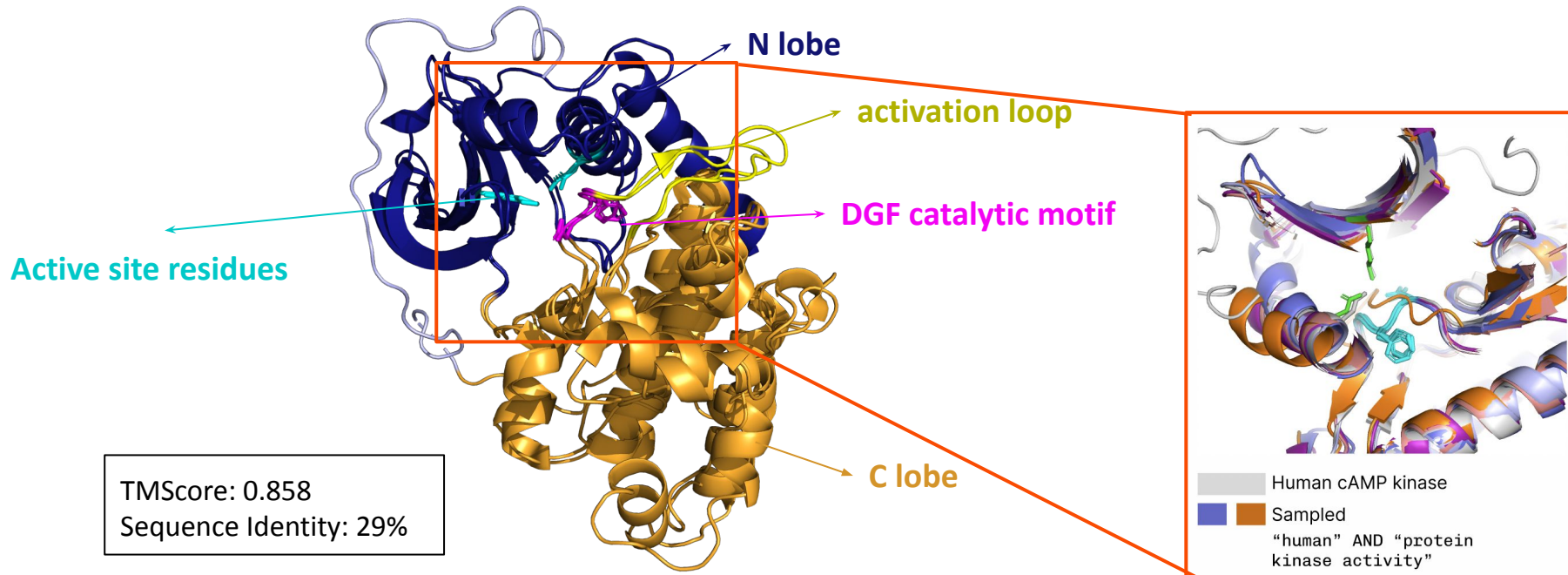
Score	Expect	Method	Identities	Positives	Gaps
327 bits(838)	3e-102	Compositional matrix adjust.	151/298(51%)	219/298(73%)	4/298(1%)
Query 2	MAIVNVGAPASGLNSAVRSLVRHCLSQGHTALAVINGFNLCK--NDSAMKII--ECKWEQ				58
Sbjct 409	+AI++VGAPA G+NSA R+ V +CL++GHT +A+ NGF GLC+ +D + + E KW IAIIHVGPAGGMNSATRAAVAYCLTRGHTPIAIHNGFPGLCRHDDKPLGAVREKWL				468
Query 59	VNLWFAKGGSQFGTARTIFNSNDLELIFDKFEANEINGLLIIGGFNSYNALTVLRHNRQQ				118
Sbjct 469	V W +KGG+ GT R++ S D+E FE + +GL +IGGF ++ A+ LR R+ VEGWISKGGSEIGTNRSL--PSEMEQTAKCFEQYKFDGLFVIGGFVAVGELRKARKD				527
Query 119	YPEFKIPMIIIPATISNNVPGTAYSLGSESSLNALCTCVDKIKQTASAKRRVVFVETLG				178
Sbjct 528	YP F IP++I+PATISNNVPGT YS+GS++ LNAL + D IKQ+ASA++RRVVFVET G YPAFNIPVILPATISNNVPGTEYSIGSDTCLNALVSYCDAIKQSASATRRRVVFVETQG				587
Query 179	GTSGYIATMAGVCCGARSYIYPEQIGIDLHKLDKDCDFLKQAFDKDPSYKSGRIIKNEA				238
Sbjct 588	G SGYIAT+AG+ GA ++Y PE+GID+ L +D + L+++F D N++G++I++NE GRSGYIATIAGLSIGATAVYTPPEGIDIKMLSRDIEHLRESFANDKQNRAGKILLRNEH				647
Query 239	ASKVYSTNIIAQLIRDESNGKFDTRTSPGHQKGGTPTSLDRVYATKMGAKAMHF				296
Sbjct 648	ASK Y+T +IA +IR+ES G+F++R ++PGH O+GGTP+ +DRV A ++ K M FI+ ASKTYTTELIANMIREESKGRFESRLAVPGHVQGGTSPMDRVRVRLAVKCMQFIE				705

- closest match: PFK1 [**6-phosphofructokinase, alpha subunit**]
- organism: **Hypocomyce scalaris** (also in the fungus kingdom)
- sequence identity: 50.67%

## Examining active site conservation

prompt: "human" AND "protein kinase activity"

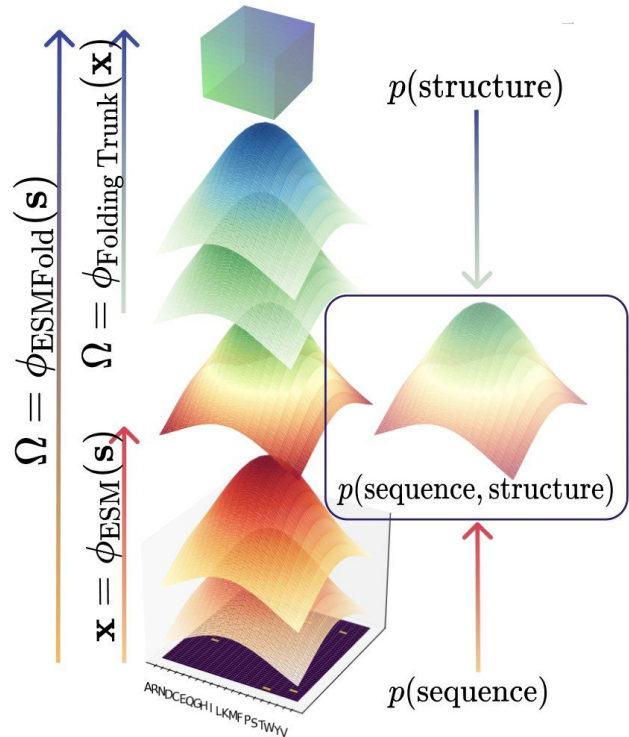
Closest Foldseek neighbor: 6cd6 (human calcium/calmodulin-dependent protein kinase kinase 1)



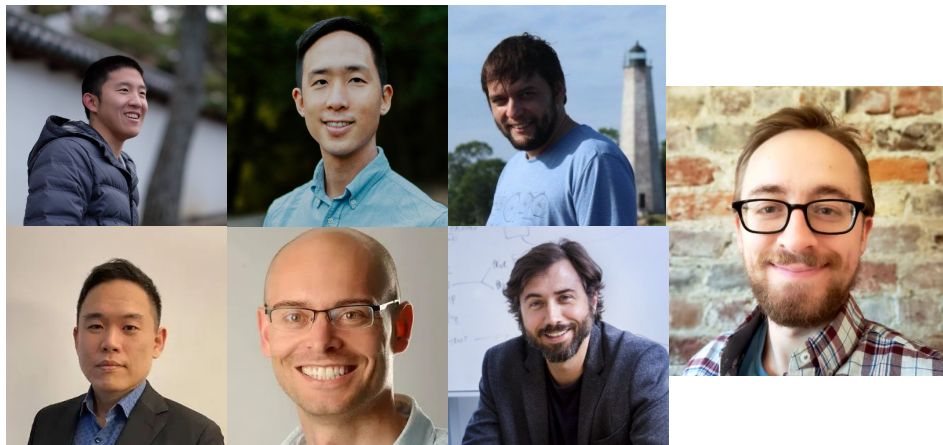


## Takeaways

- The latent space of ESMFold is disorganized with massive activations
- Compressing the latent space shows that *many* channels might be extraneous for structure prediction
- Information content relating to sequence, structure, and function is not symmetrical
- CHEAP regularization helps with latent diffusion model training, leading to **an all-atom co-generation model with sequence database scale coverage**



# Thanks!



## Berkeley

Amy X. Lu  
Wilson Yan  
Pieter Abbeel

## Prescient Design

Sai Pooja Mahajan  
Sarah Robinson  
Vladimir Gligorijevic  
Kyunghyun Cho  
Richard Bonneau  
Nathan C. Frey

## Microsoft Research

Kevin Yang

Paper: [bit.ly/cheap-proteins](https://bit.ly/cheap-proteins)

Code & weights: [github.com/amyxlu/cheap-proteins](https://github.com/amyxlu/cheap-proteins)



Paper



GitHub



@amyxlu



amyxlu.github.io



amyxlu@berkeley.edu

# prompt: “mouse” AND “6-phosphofructokinase activity”

Select target residues to highlight their structure.  
Click on highlighted sequences to dehighlight the corresponding chain.

CLEAR SELECTION ✕

→ 4xz2-assembly1\_C

```
Q 2 LAVHVGAP SAGINAAVRS AVRTGINNGYEVLFIQDGFQGLLKGE SHLHEVHWN SIA
  +AV+V G A P +A G +N A A V R S A V R G I +G + L I D G F G + K G ++ E + W ++
T 363 VAVINVGAP AAGMNAAVRS AVRVGTADGHRMLAIYDGFDFAGK--Q I K E I G W T D V C
Q 62 QTGGSDLHTARGRAMTEEQGLAEAAKAL EDHG INGLMVI GGFDNL SGVNM LRQARSK
  GGS L T R + L E A + H I N + I G G F + G + L A R K
T 421 GQGGSLGTGRVLP G--KYLEE I A T Q M R T H S I N A L L I I G G F E A Y L G L L E L S A A R E K
Q 122 LTNQIPLVAVPCTINNDVPGTMDLTGDSACNAIAEIVDRIKLSASATKSRVFIET
  + +P+V VP T++N+VPG+D+++G D+A N I + D R I K S A S + T K R V F + I E T
T 478 FC--VPMVMPATVSNNVPGSDFSIGADTALNTITDTCDRIKQASAGTKRRVFIET
Q 182 FCGYLATCAGIACGADACYVMEEEGKISVKNVPIQFEIMVTHLRGMRGLILHLER
  +CGYLA +G+A GADA Y++EE +++++ + E + ++ ++RGL+L E
T 536 YCYLANMGLAAGADAAYIFEFP--FDIRLQSNVEHLTEKMKTTIQRLGLVLRNES
Q 242 QYTTQFINKLFSEEGKGVFDIRINVLGYMQGGSPTPHNRNFGARCGMKCLLWL
  +YTT FI +L+SEEGKGVFD R NVLG+MQGG+P+P DRNFG + + + W+
T 594 NYTTDFIQLYSEEGKGVFDRCRNVLGHMQGGSPFPDRNFGTKISARAMEWI
```

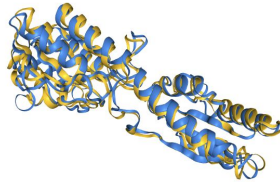
Select target residues to highlight their structure.  
Click on highlighted sequences to dehighlight the corresponding chain.

CLEAR SELECTION ✕

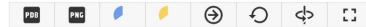
→ 3o8n-assemblyL\_A

```
Q 2 LAVHVGAP SAGINAAVRS AVRTGINNGYEVLFIQDGFQGLLKGE SHLHEVHWN SIA
  +AVM+V G A P +A G +N A A V R S V R G + G V L ++DGF+G K G ++ E + W ++
T 395 VAVMNVGAP AAGMNAAVRSTVRIGLIQGNRVLVHDGFEGPAGK--Q I E E A G W S V Y G
Q 62 QTGGSDLHTARGRAMTEEQGLAEAAKAL EDHG INGLMVI GGFDNL SGVNM LRQARSK
  GGS L + R + + + + + + + + I + G L + I G G F + + G L + R +
T 453 GQGGSKLGSKRT--LPK-KSFEQISANITKFNIIQGLVIIGFEAYTGGLELMGRKQ
Q 122 LTNQIPLVAVPCTINNDVPGTMDLTGDSACNAIAEIVDRIKLSASATKSRVFIET
  L I P + V + P T + + N + V P G + D + + + G D + A N I D R I K S A + T K R V F + I E T
T 510 LC--IPFVVIPTVSNNVPGSDFSVGADTALNTICTDRIKQSAAGTKRRVFIET
Q 182 FCGYLATCAGIACGADACYVMEEEGKISVKNVPIQFEIMVTHLRGMRGLILHLER
  +CGYLAT AG+A GADA Y++EE +++++ + E + V ++ + RGL+L E+
T 568 YCYLATMAGLAAGADAAYIFEFP--FTIRDLQANVEHLVQMKTTIKRGLVLRNFK
Q 242 QYTTQFINKLFSEEGKGVFDIRINVLGYMQGGSPTPHNRNFGARCGMKCLLWL
  +YTT FI L+SEEGK+FD R NVLG+MQGGSPTP DRNF+ + G K + W+
T 626 NYTTDFIFNLSEEGKGFDRKRVLGHMQGGSPFPDRNFATKMGAKAMNWM
```

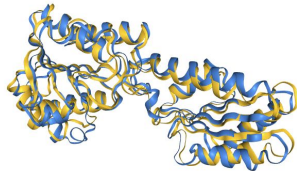
TM-Score: 0.93557  
RMSD: 1.78



4x2  
human  
e-value=48.2



TM-Score: 0.92514  
RMSD: 1.93



3o8n  
rabbit  
e-value=46.5



- species conditioning is biased database composition
  - e.g. performance on “HUMAN” and “ECOLI” is better, since they are better represented in the database

# Why GO terms and organism?

- generative protein design should propose designs that might be useful. What are some possible use cases?
  - being able to express in model organisms
  - humanization efforts
  - enzyme engineering

Organism: encourages generating samples that might express.

GO term: gives us finer control over monomer generation

