



# Repurposing Protein Folding Models for All-Atom Generation via Latent Space Compression

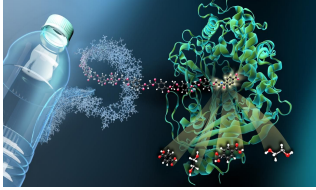
**Amy X. Lu**

PhD Student, UC Berkeley

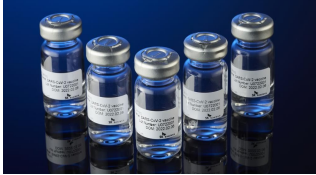
March 3rd, 2025

Baker Lab Journal Club

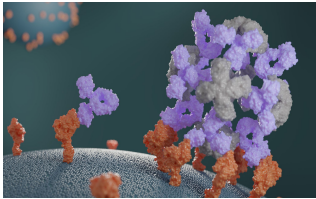
# Designing novel proteins with desired properties **is useful**



Plastic  
degrading  
enzymes



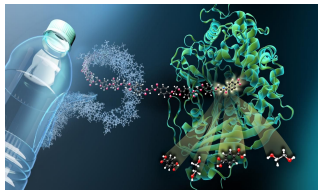
Vaccine  
development



Antibody  
therapeutics

*and much more...*

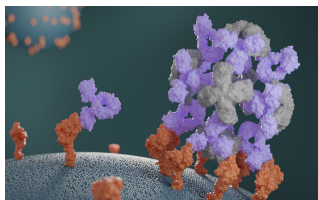
# Designing novel proteins with desired properties is hard



Plastic  
degrading  
enzymes



Vaccine  
development



Antibody  
therapeutics

*and much more...*

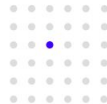
HUMAN  
PROTEINS

~20,000



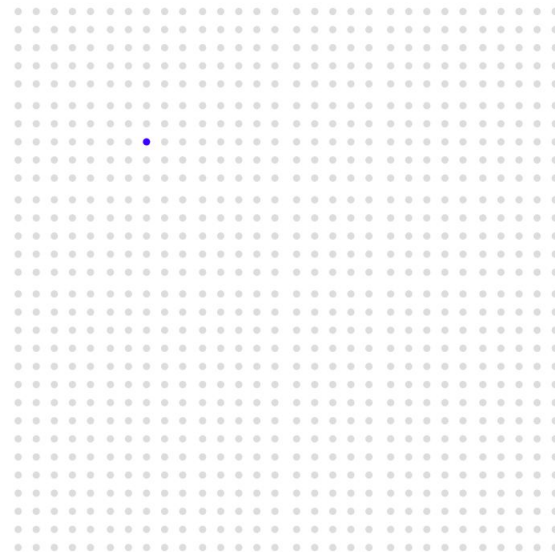
UNIQUE  
PROTEINS  
ON EARTH

~1,000,000,000,  
000,000,000

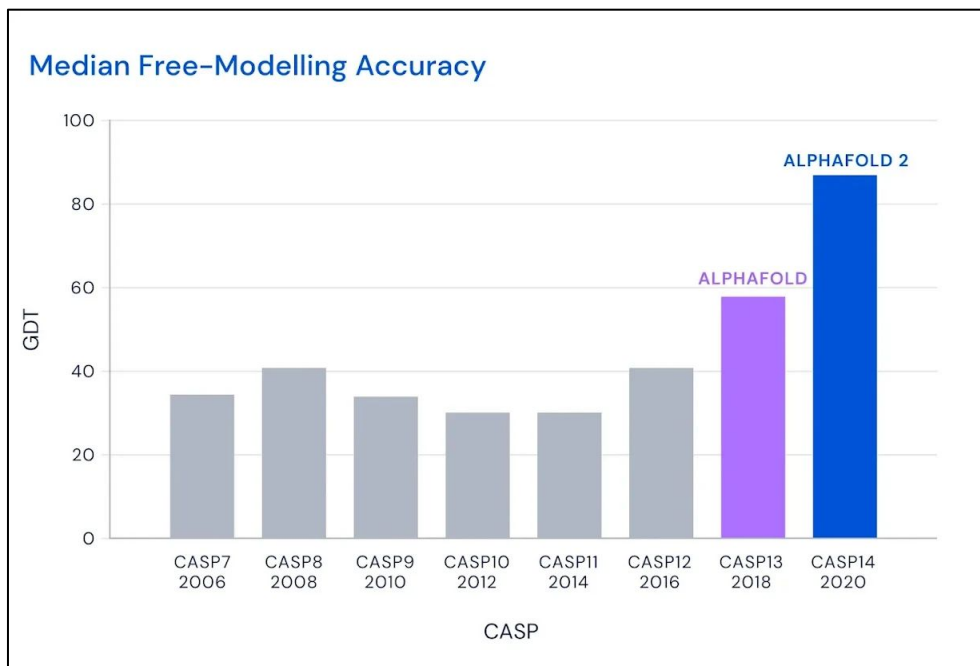
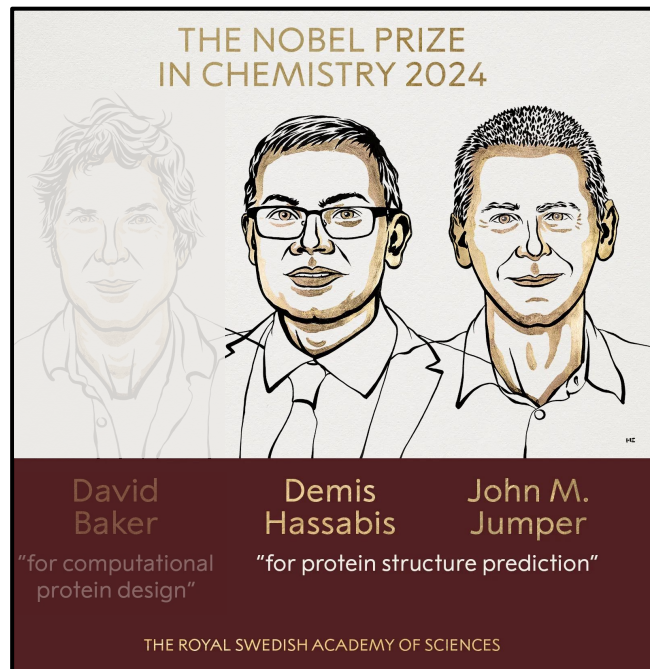


POSSIBLE PROTEINS

~10,000,000,000,000,000,000,000,000,000,  
000,000,000,000,000,000,000,000,000,000,  
000,000,000,000,000,000,000,000,000,000,  
000,000,000,000,000,000,000,000,000,000

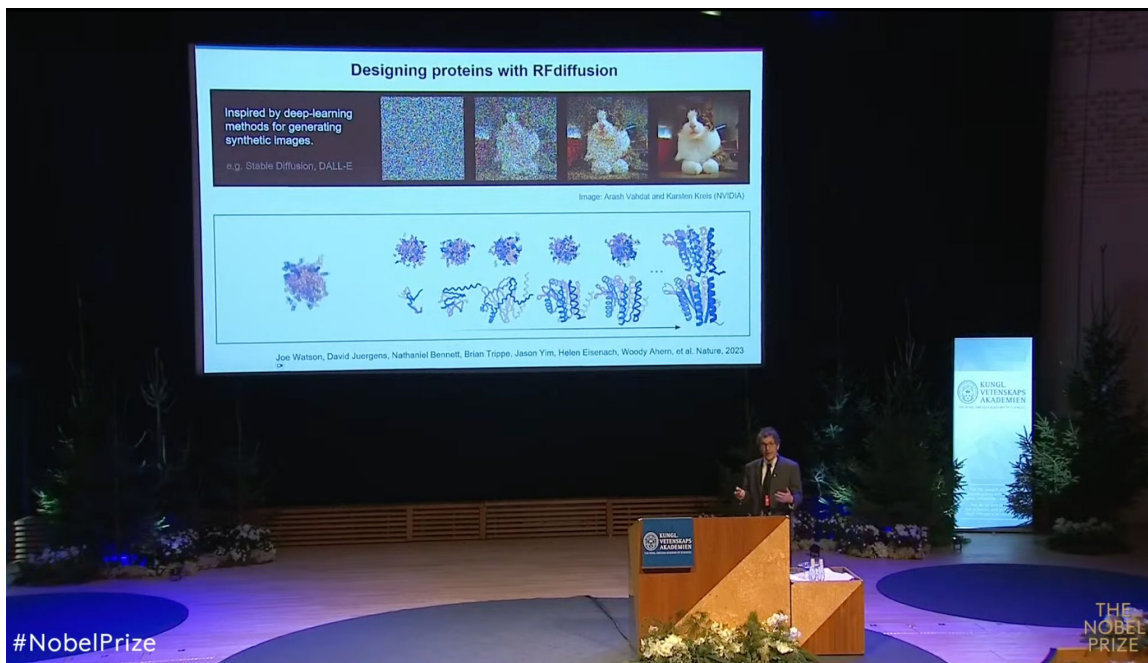
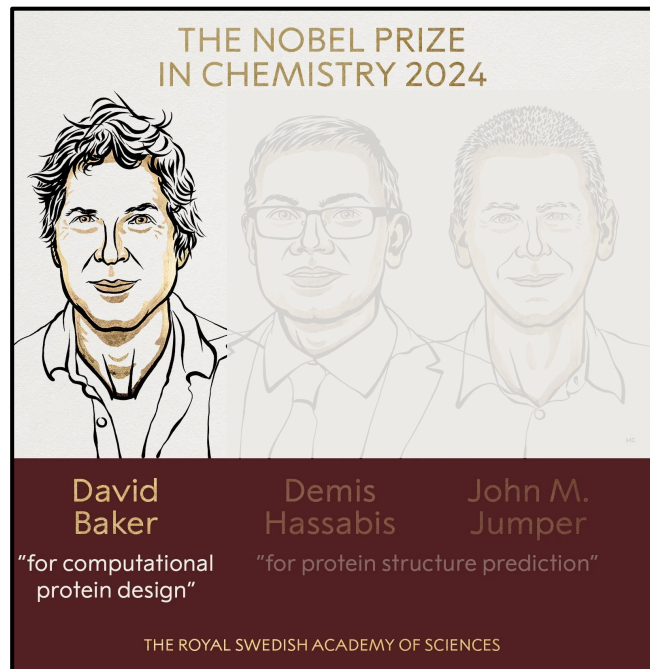


# The potential of deep learning for protein structure prediction





# The potential of deep learning for protein structure generation



# But what else would we need for drug discovery?

this talk

## Co-generation

→ can we simultaneously generate sequence and structure?

## Control

→ how can we specify complex and multi-objective constraints?

## Immunogenicity

→ can we achieve organism specificity?

## Deployment

→ can we speed up inference? what about ethics and biosafety concerns?

## Task generalization

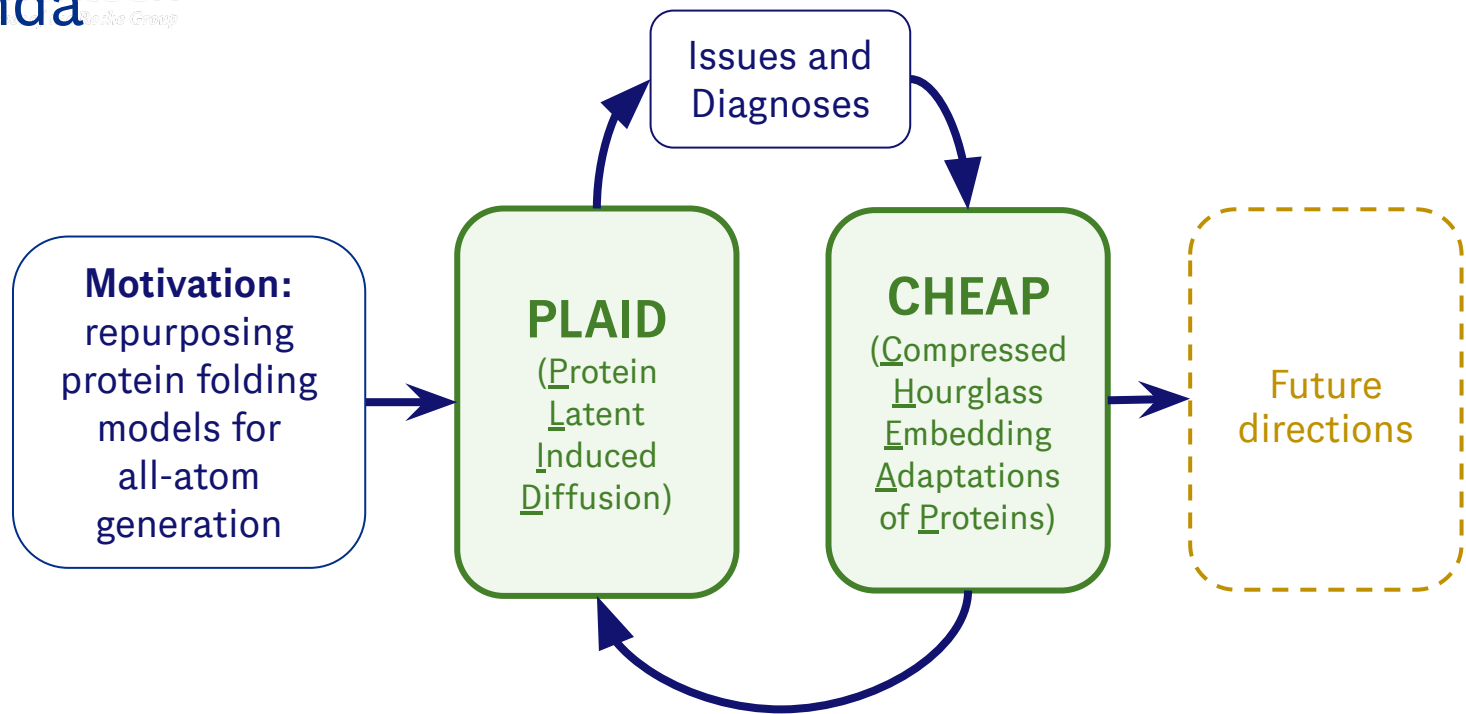
→ how can we generalize few-shot to new targets?

## Modeling complexes

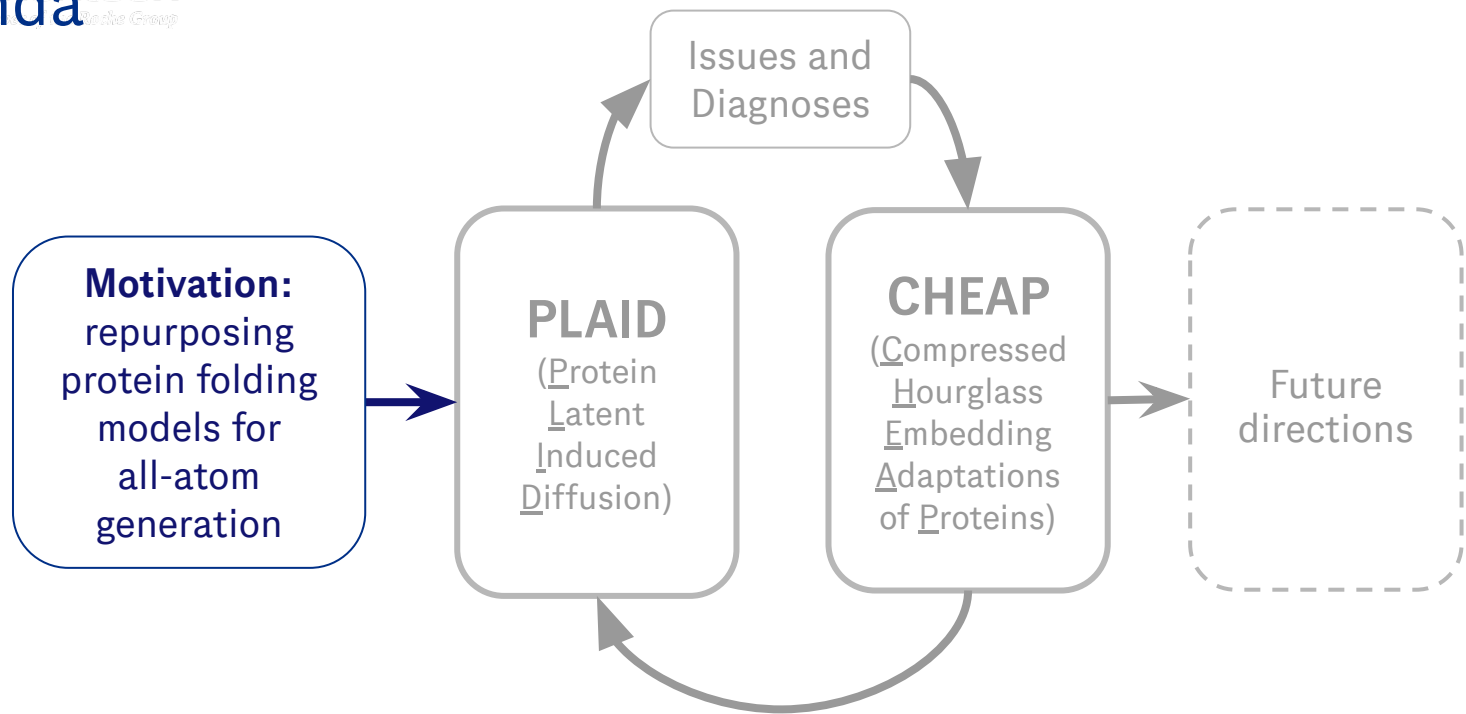
→ can we design small- & large-molecule binders?

other work / planned future work

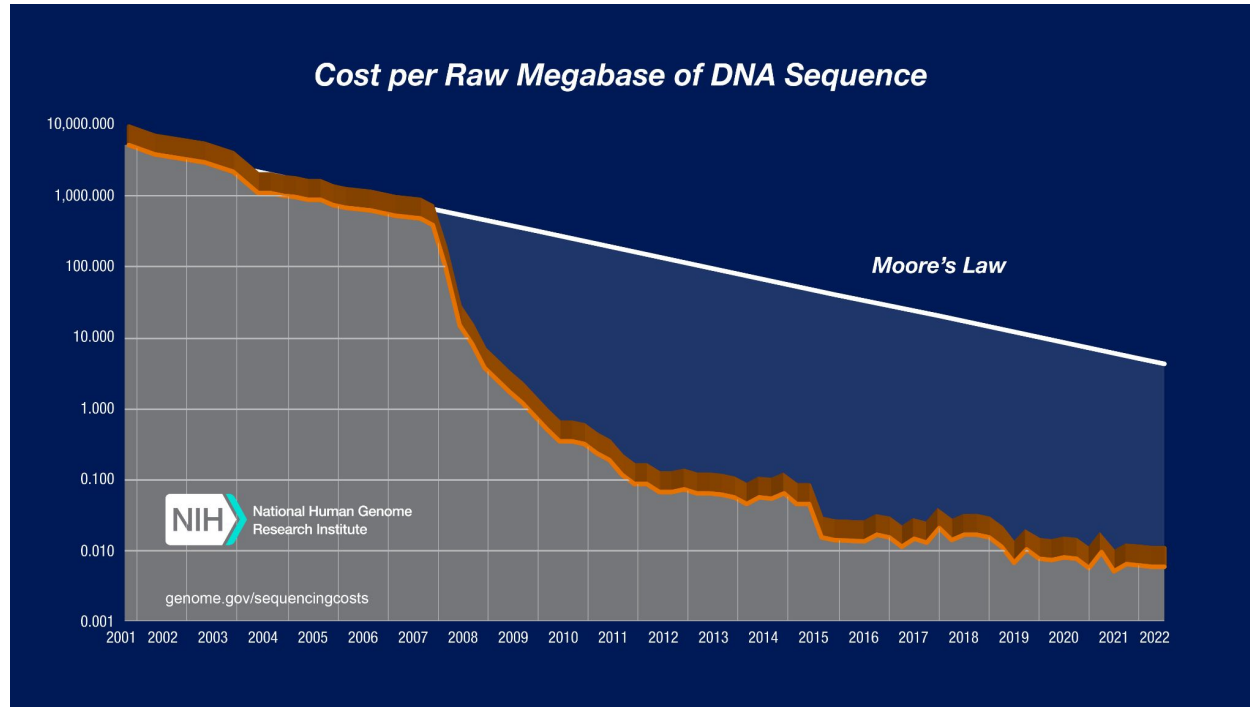
# Agenda



# Agenda

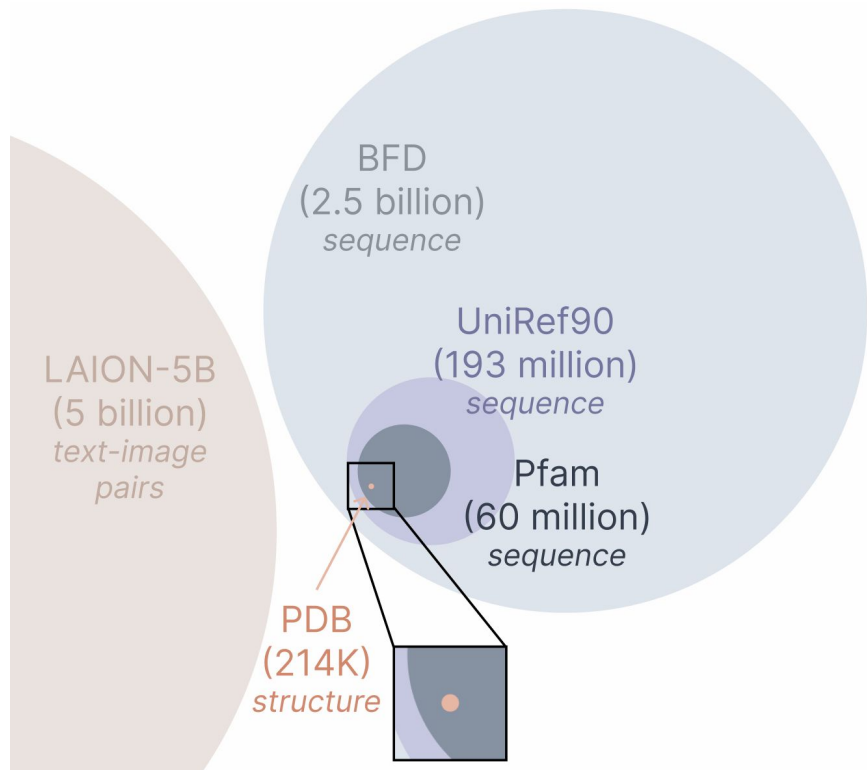


# Sequence data is cheaper to collect than structure

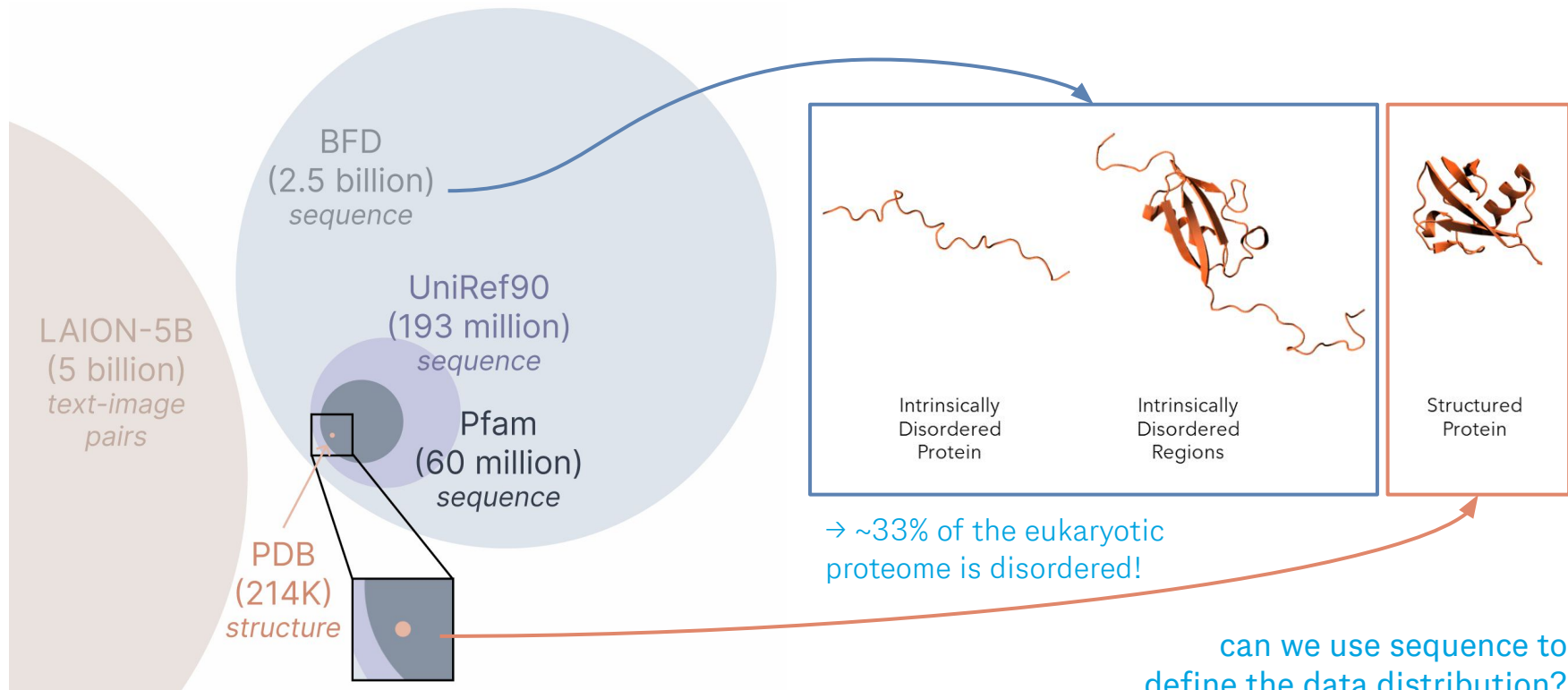


Source: <https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data>

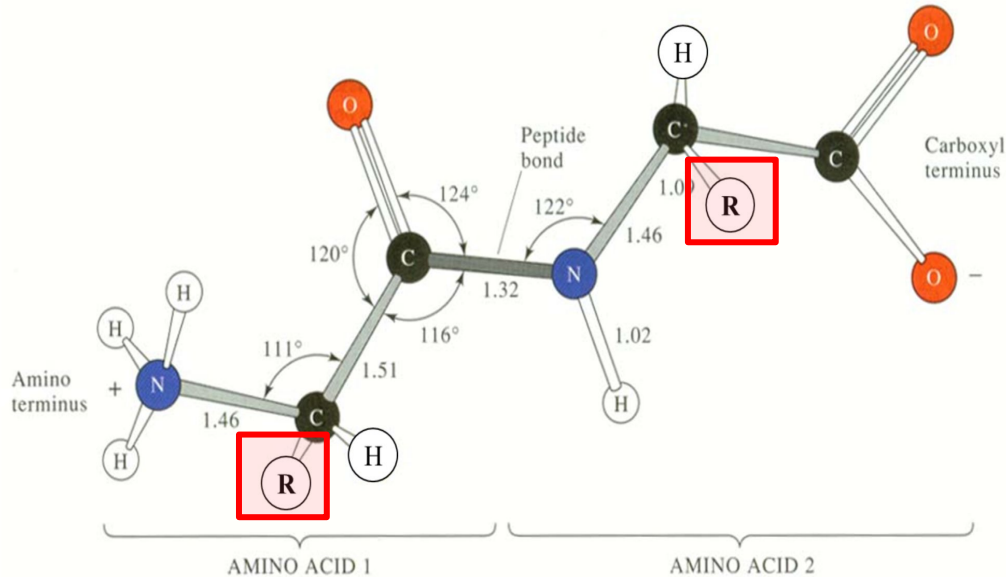
# Sequence data is more **abundant** than structure



# Sequence data has different coverage than structure



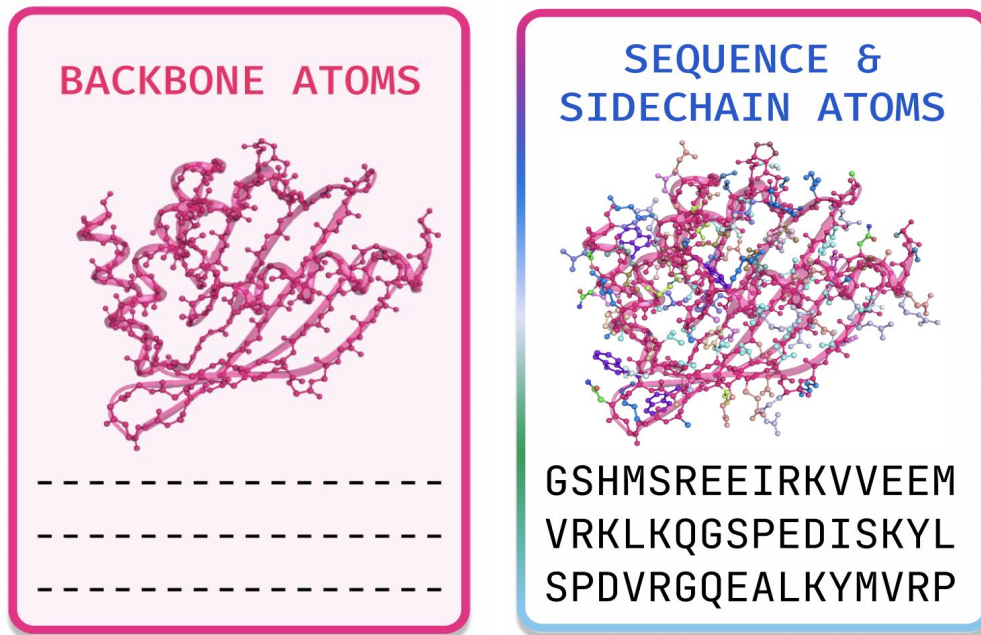
# Backbone structure vs. all-atom structure



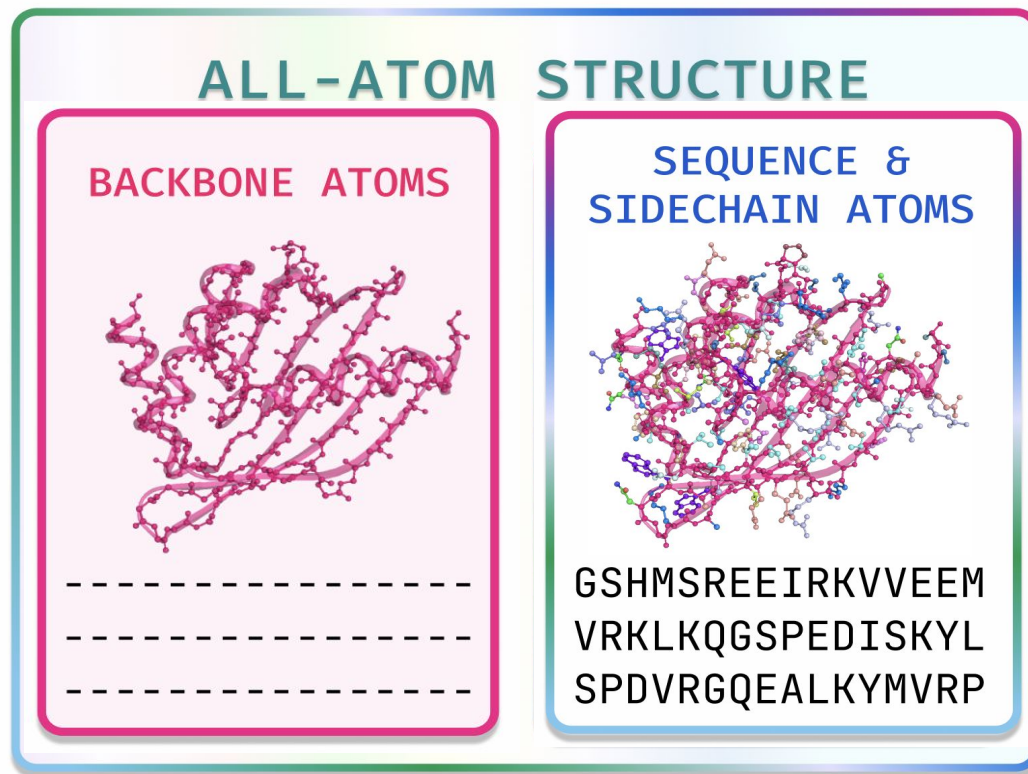
(order of t-shirts => protein sequence)



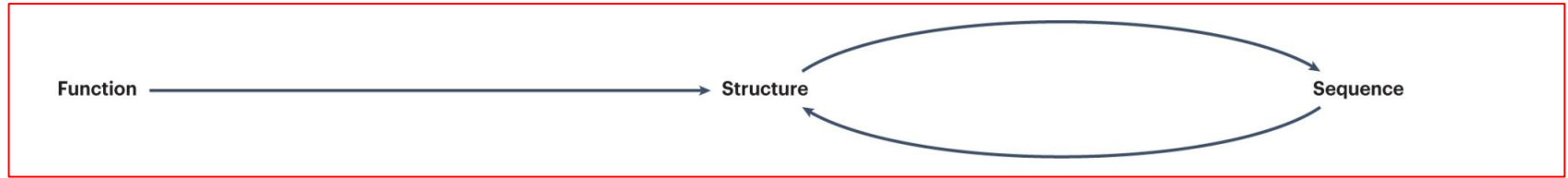
# Sidechain atoms generation require knowing the sequence



# All-atom design as a multimodal generation problem

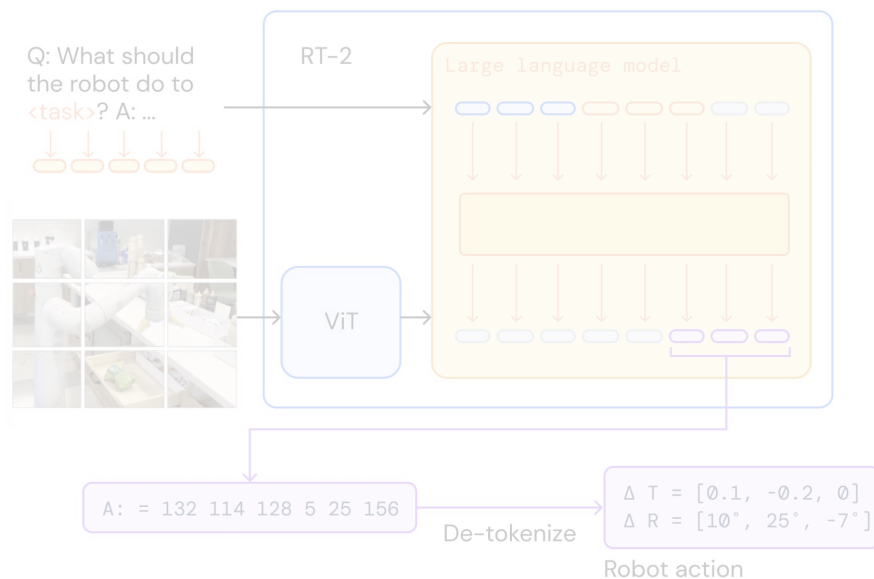


# Problem: Current all-atom methods iterate between sequence and structure design



Chu, A.E., Lu, T. & Huang, PS. **Sparks of function by de novo protein design**. Nat Biotechnol 42, 203–215 (2024).

# Motivation: Can we repurpose priors from pretrained models?



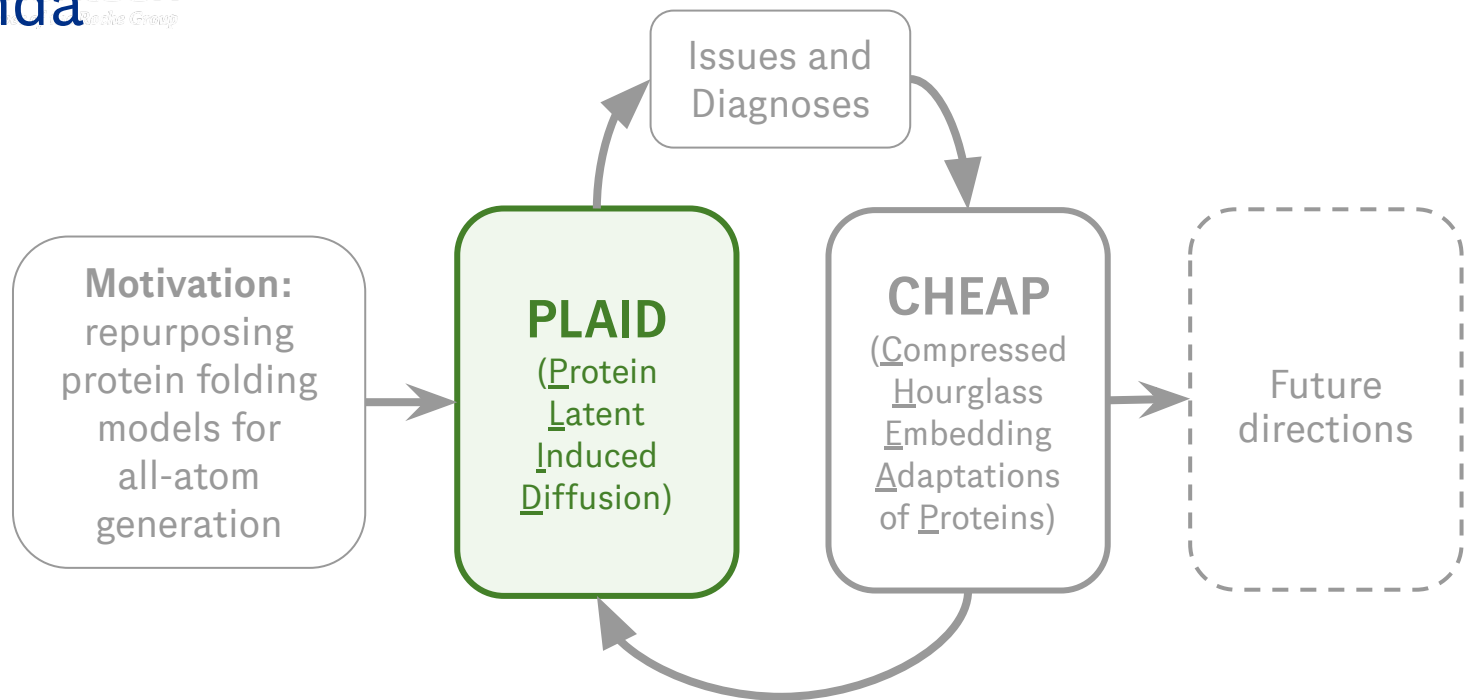
Vision-language models trained on internet-scale datasets capture useful priors for robotic tasks.

How can we apply this to biology?

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

Can we sample all-atom structure from the joint distribution  $p(\text{sequence}, \text{structure})$  and use priors from pretrained protein folding models?

# Agenda

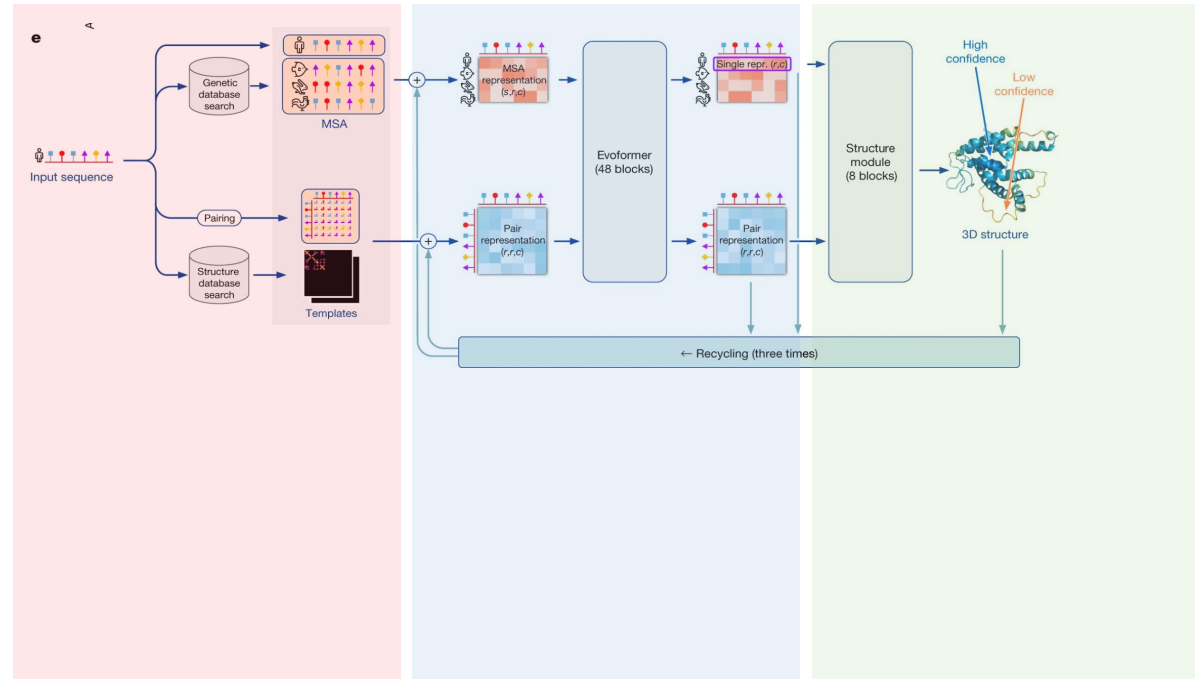


# The base components: protein folding model architectures



## AlphaFold2:

Uses an explicit retrieval step



harness additional  
sequence-based priors

learn structural features  
from sequence latents

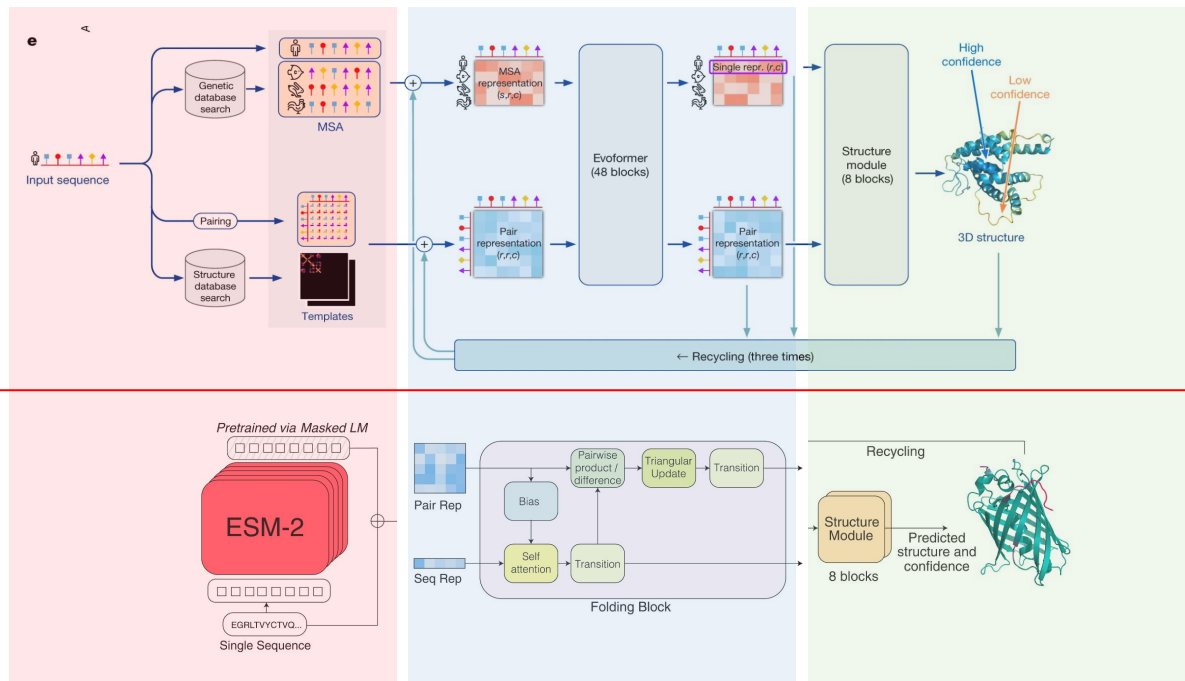
generate structures

# The base components: protein folding model architectures



## AlphaFold2:

Uses an explicit retrieval step



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     )
```

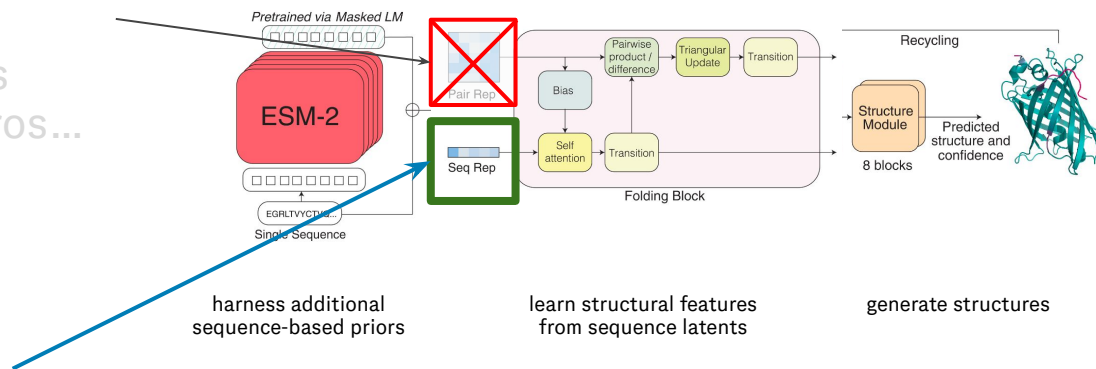


Observation: at inference,  
the pairwise input is  
initialized as zeros...

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



→ Sequence representation contains all information about the structure!

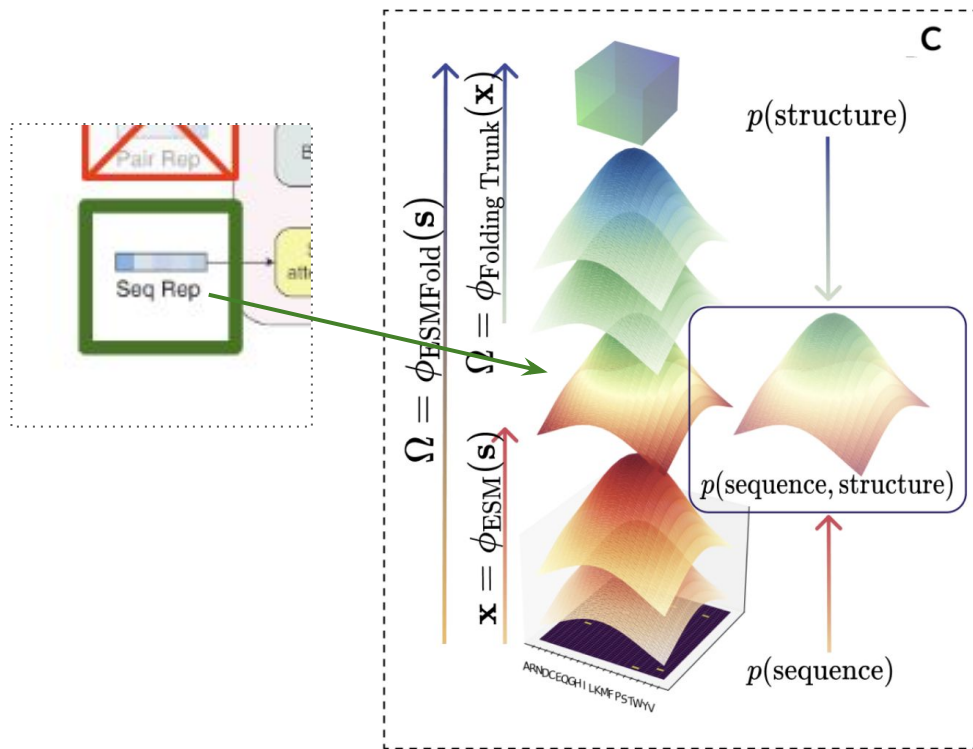


Observation: at inference, the pairwise input is initialized as zeros...

→ Sequence representation contains all information about the structure!



Generating this embedding would only require the sequence during training.



# A 5 second primer on diffusion models

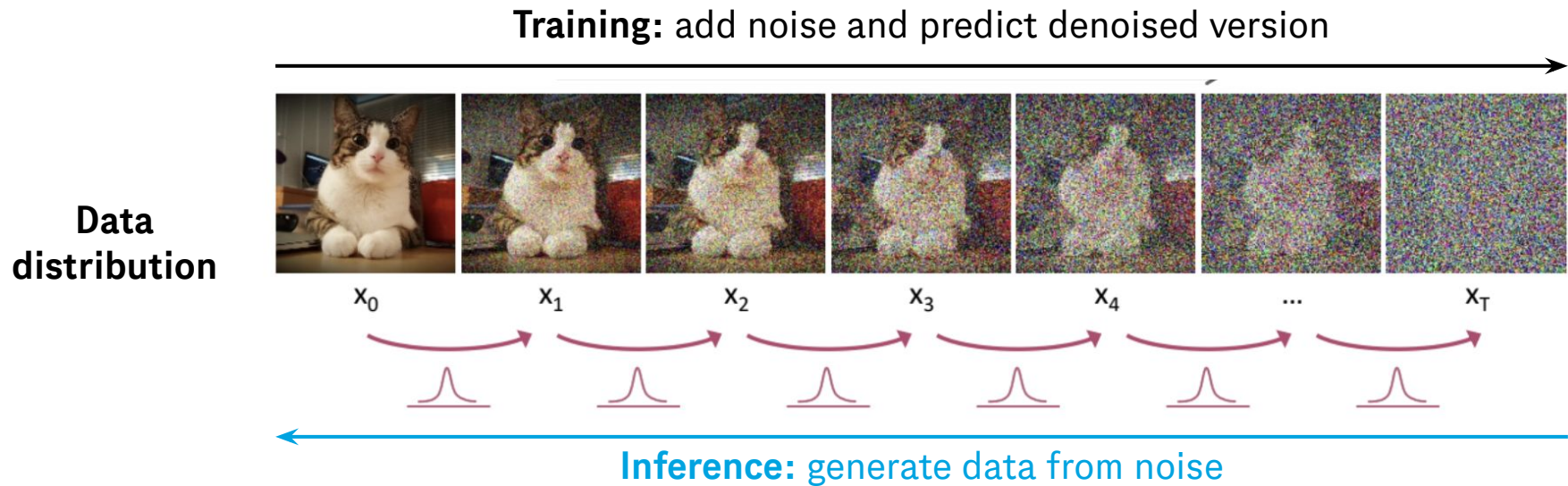
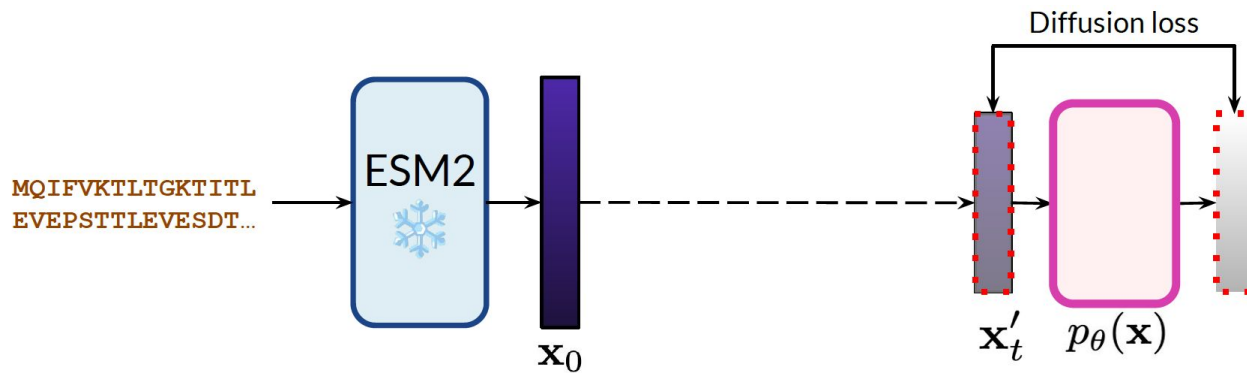
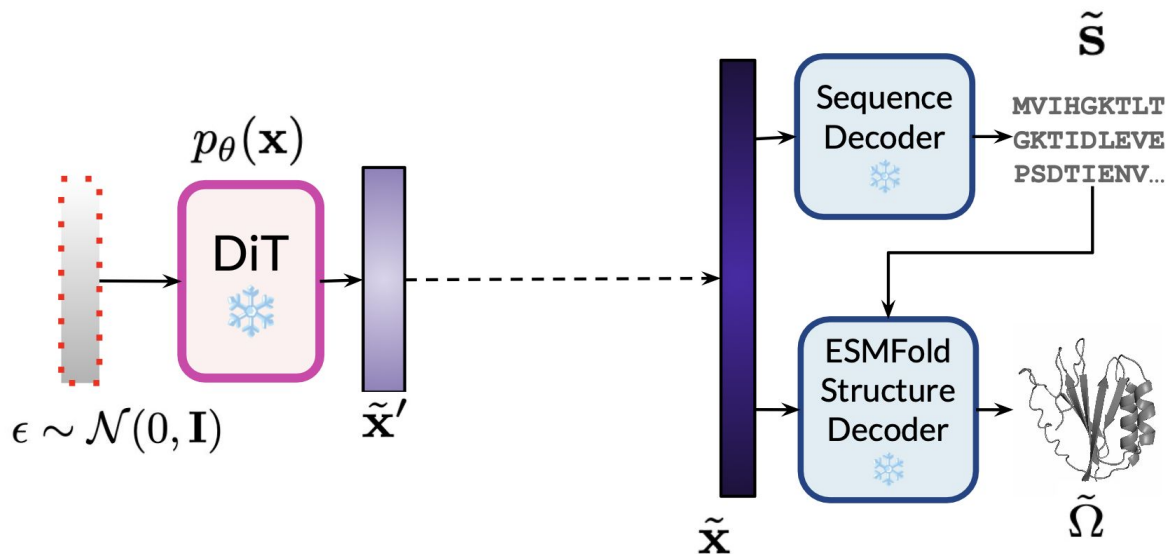


Image source: Arash Vahdat and Karsten Kreis

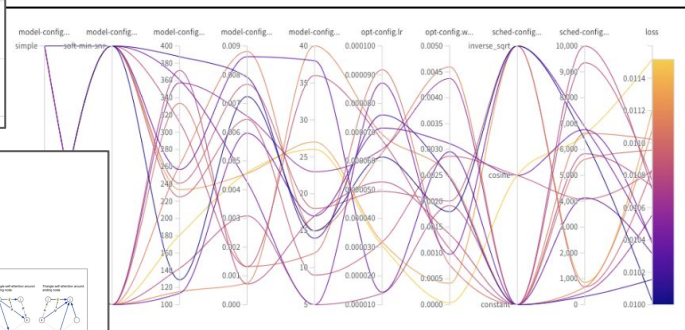
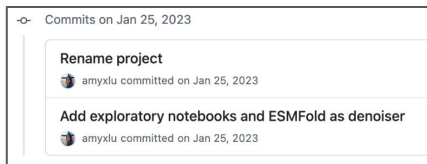
# PLAID v0.5: Training a latent diffusion model



# PLAID v0.5: Inference-time all-atom generation



# PLAID v0.5: Early attempts



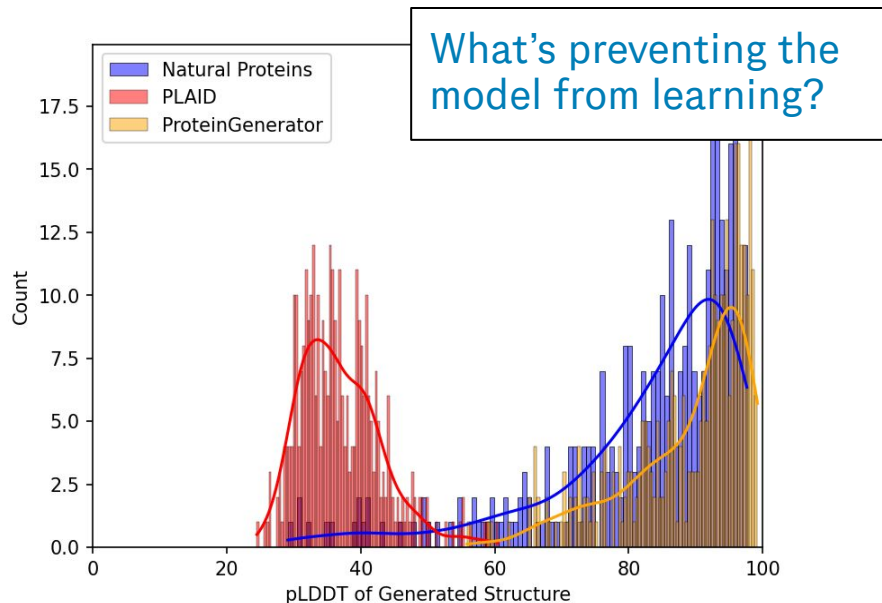
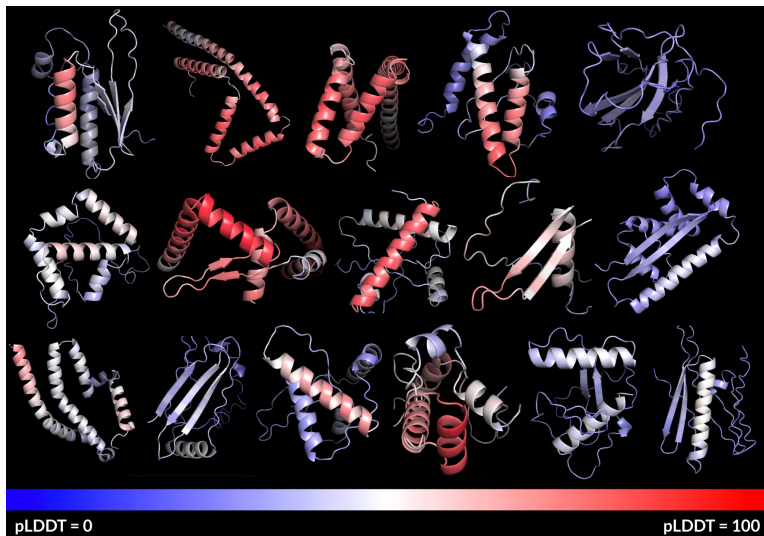
tried many things which did not work 😞

- adding a structure/sequence semantic reconstruction loss
- Finetuning ESMFold (à la RFDiffusion)
  - predict  $x_0$  instead of noise?
- Huber vs. MSE loss
- Gradient-based guidance towards properties
- Triangle self-attention
- ...

04/20/2023

- **Architecture:** cannot use UNet naively due to padding.
  - Vanilla transformer (# parameters: 51,448,832)
  - Vanilla transformer, large (126,109,696 parameters)
  - Invariant Point Attention:
    - Full 48 blocks for AlphaFold2 (~57,000,000 parameters)
    - Finetune first 4 blocks
    - Use first 4 blocks, but initialize randomly.
- **Loss:** Latent diffusion only, or also constrain for the fidelity of the two additional modalities?
  - Reconstruction loss only
  - Reconstruction loss + structure loss + sequence loss
- **Dataset:** If we don't need the structure loss, then we can train on a larger dataset w/ better coverage of the natural distribution of proteins.
- **Tuning hyperparameters:**
  - Beta schedule, LR, LR scheduler, batch size & grad accumulation, grad norm clipping, etc.

# PLAID v0.5: Early attempts



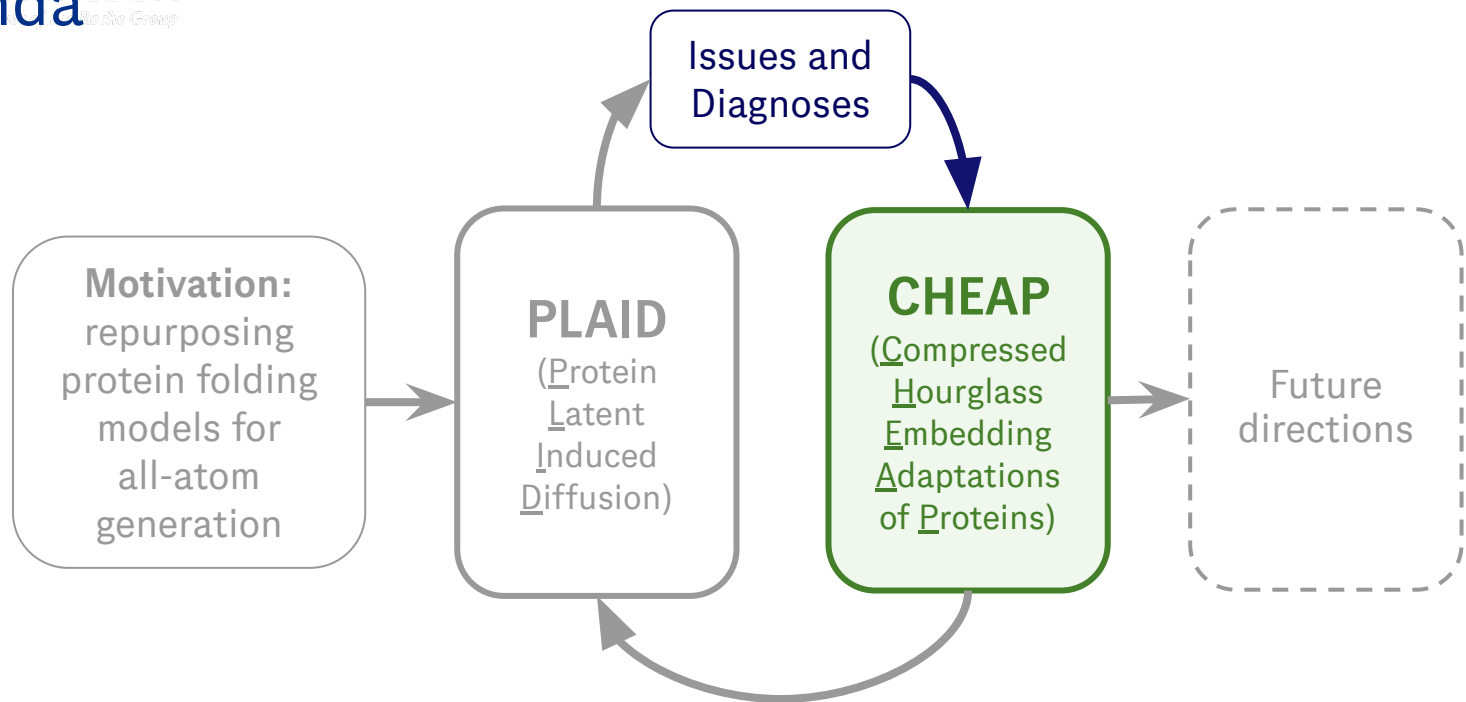
## PLAID v0.5: Generating Protein Sequence and Structure Without Structural Training Data

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

ICML 2024 Workshop on Machine Learning for Life and Material Sciences



# Agenda



# Issues and hypotheses

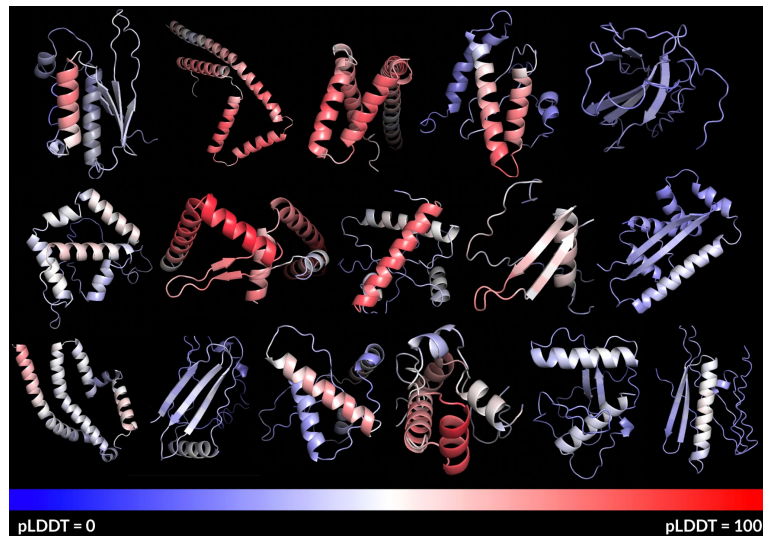
- Latent space requires regularization

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

Rombach et al. [High-Resolution Image Synthesis with Latent Diffusion Models](#), CVPR 2022

# Issues and hypotheses

- Latent space requires regularization
- Overcome  $O(L^2)$  memory constraints and increase protein length to 512



# Issues and hypotheses

- Latent space requires regularization
- Overcome  $O(L^2)$  memory constraints and increase protein length to 512
- Large latent space corresponds to **high-resolution** image generation
  - Rombach et al. latent space:  
 $H \times W \times 4 = 64 \times 64 \times 4$
  - Ours:  
 $L \times 1024 = 512 \times 1024$

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



# Addressing the hypotheses: latent space regularization

## Issues and hypotheses -> CHEAP

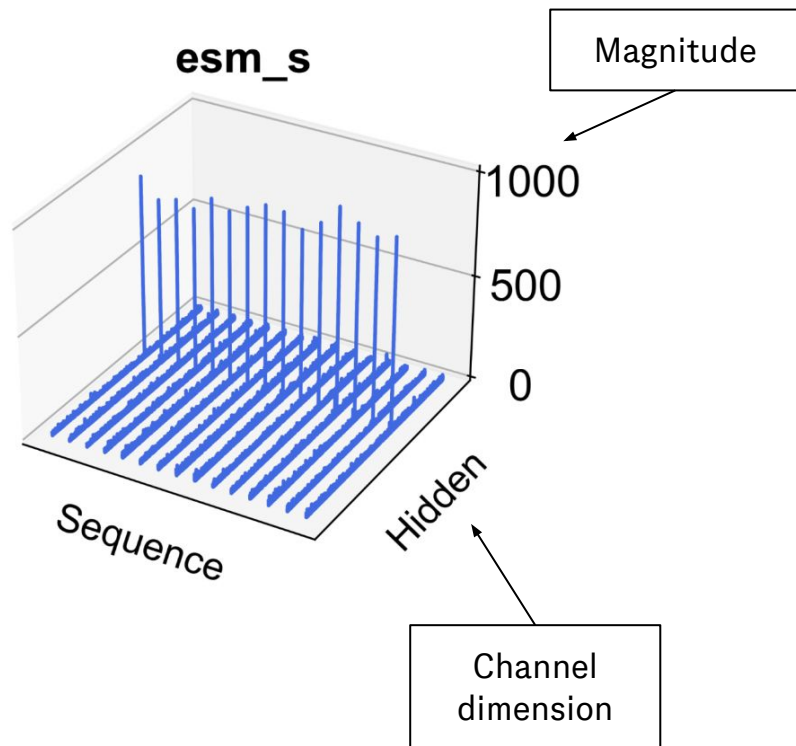
- Latent space requires regularization

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

Rombach et al. [High-Resolution Image Synthesis with Latent Diffusion Models](#), CVPR 2022

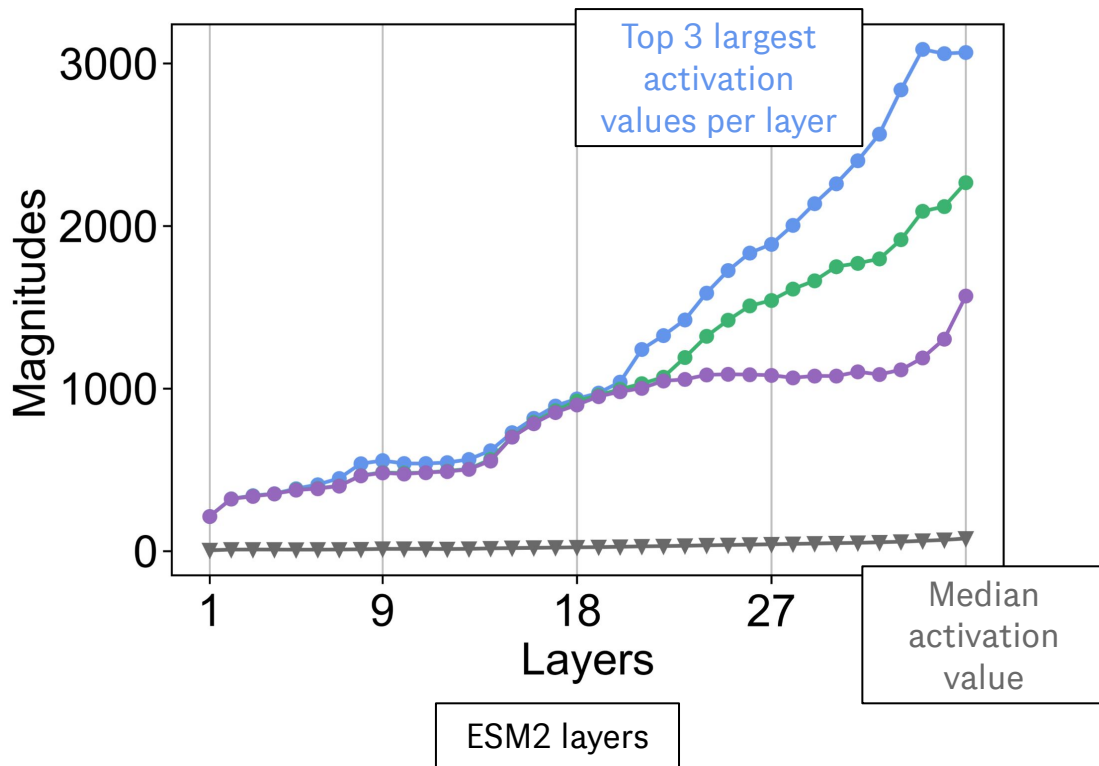
(1) Are the problems coming from the ESMFold latent space itself?

# ESMFold latent space exhibits pathologically large values



Latent space will require regularization for diffusion to work.

# ESM-Fold ESM2 latent space exhibits pathologically large values



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

→ a pervasive issue across LLMs, ViTs, etc.

[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

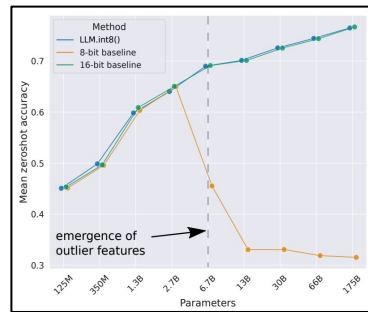
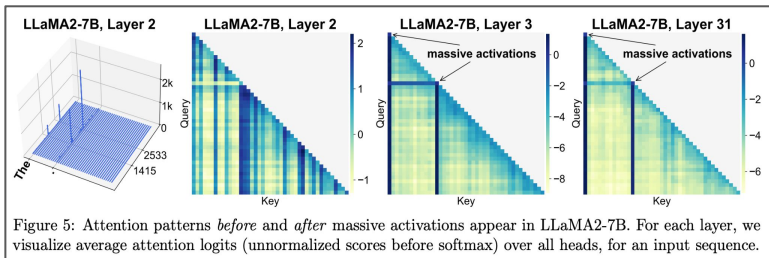
## LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale

Tim Dettmers<sup>Λ\*</sup>

Mike Lewis<sup>†</sup>

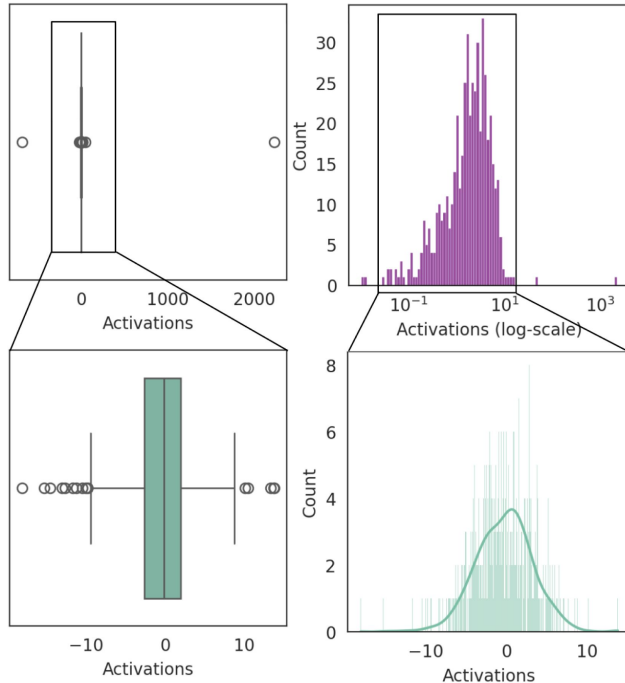
Younes Belkada<sup>§‡</sup>

Luke Zettlemoyer<sup>†λ</sup>

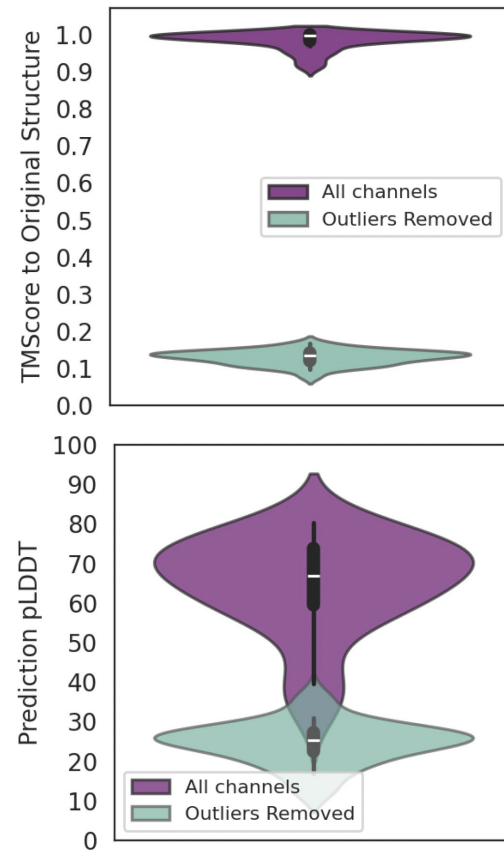
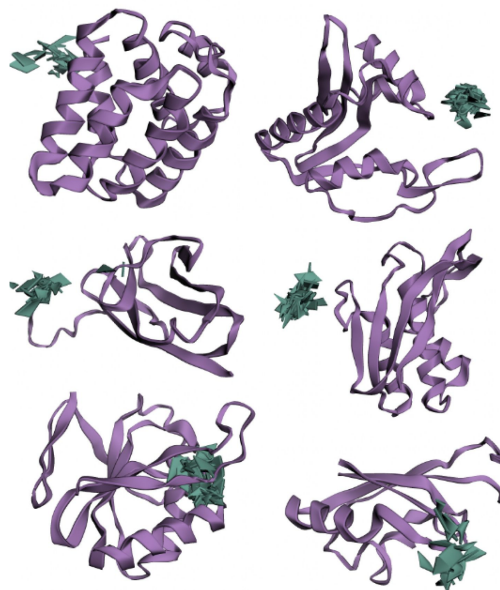
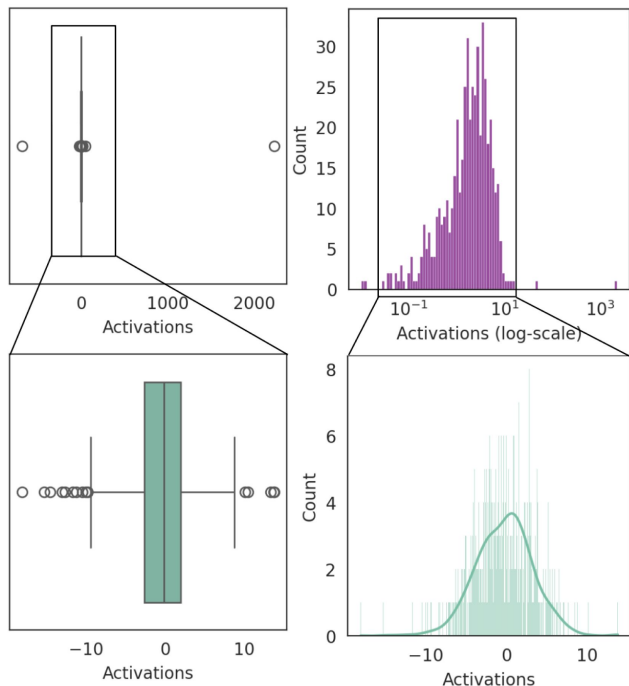




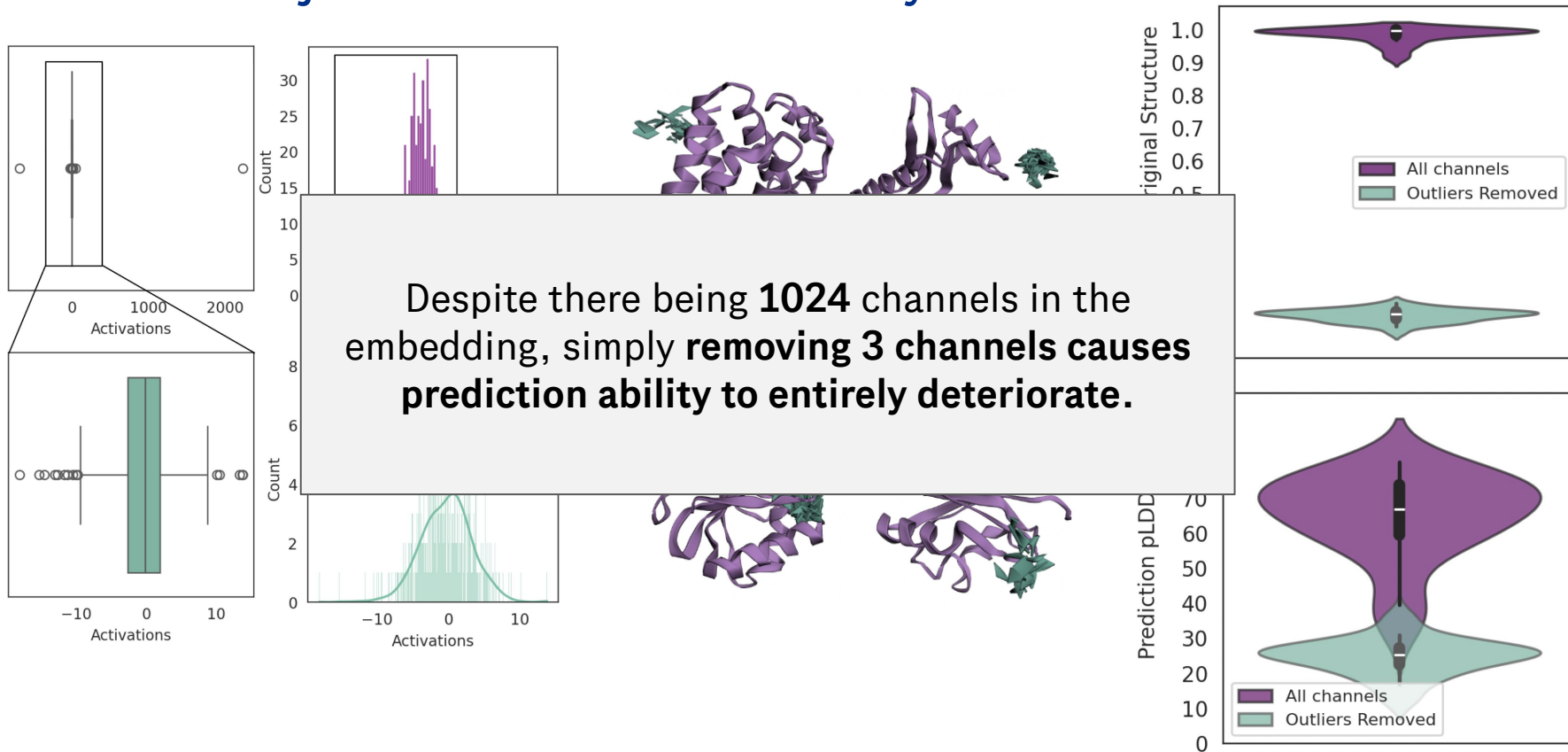
# What if we just remove these wacky channels?



# What if we just remove these wacky channels?



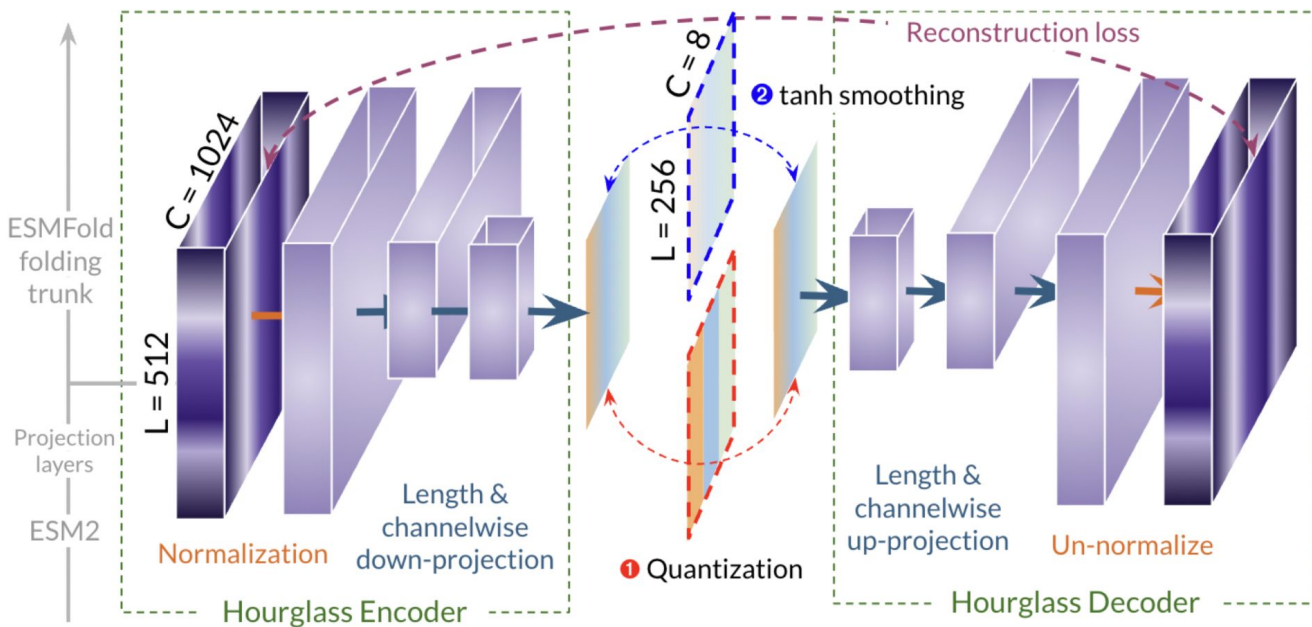
# What if we just remove these wacky channels?



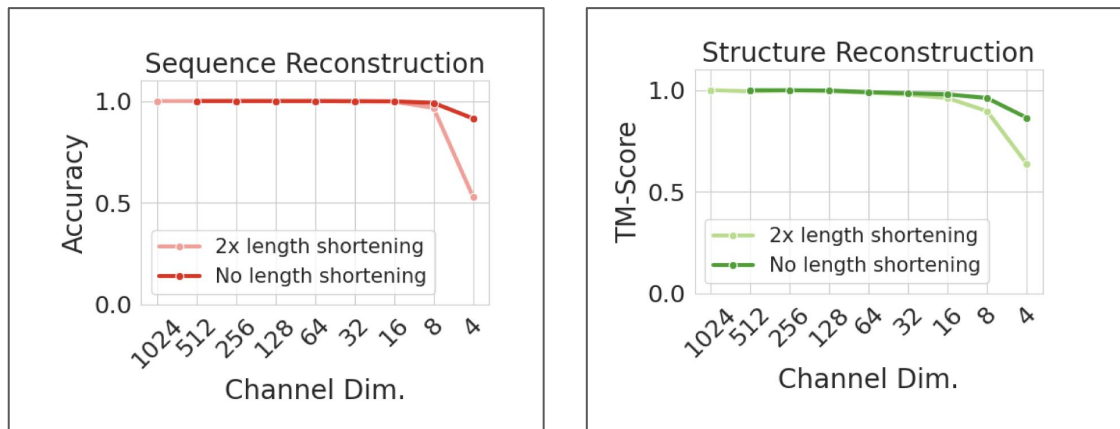
# An autoencoder for protein embedding compression

Simple fix for massive activations:  
standardize each  
channel independently.

$$\mathbf{x}' = \frac{\mathbf{x} - \mathbf{x}_{\min}}{\mathbf{x}_{\max} - \mathbf{x}_{\min}}$$

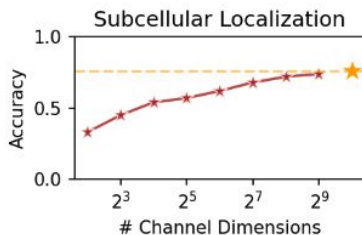
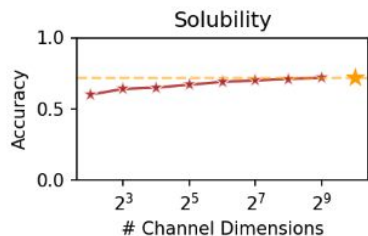


# Turns out the latent space is highly compressible!



Sequence information is easier to retain than structure.

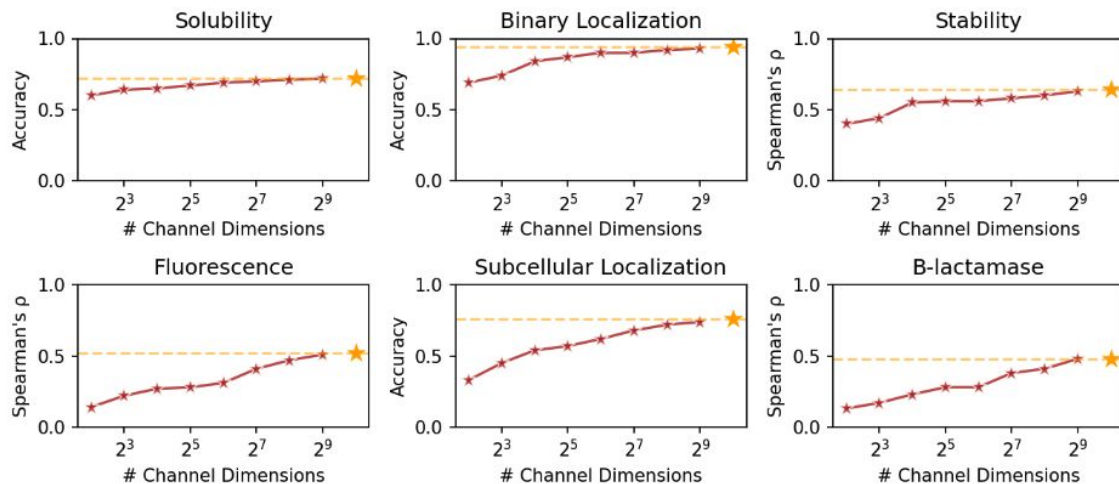
# What about function information?



Performance degradation with  
compression is more gradual...

...for some functions.

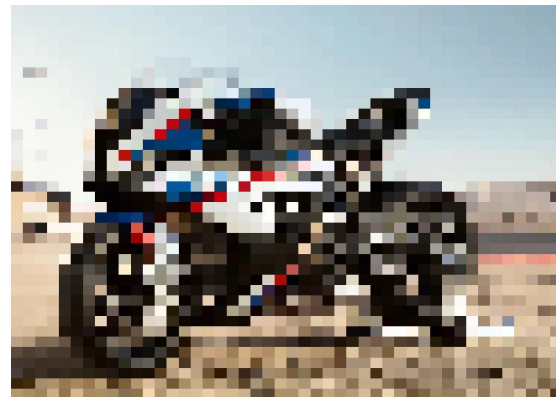
# What about function information?



Performance degradation with  
compression is more gradual...

...for some functions.

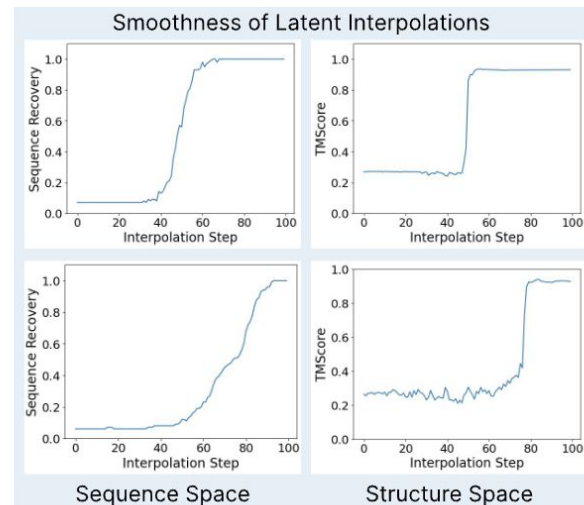
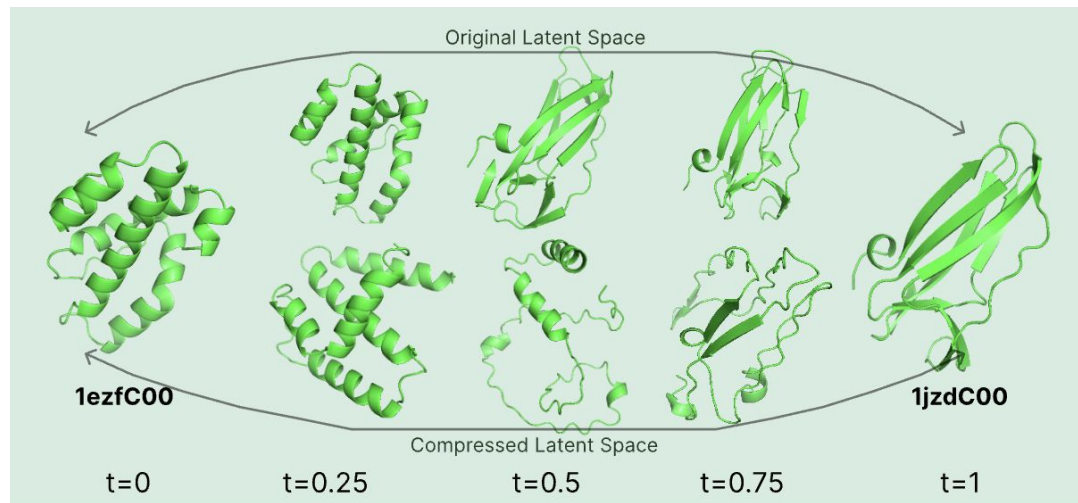
# Intuition: what is the speed of this motorcycle?



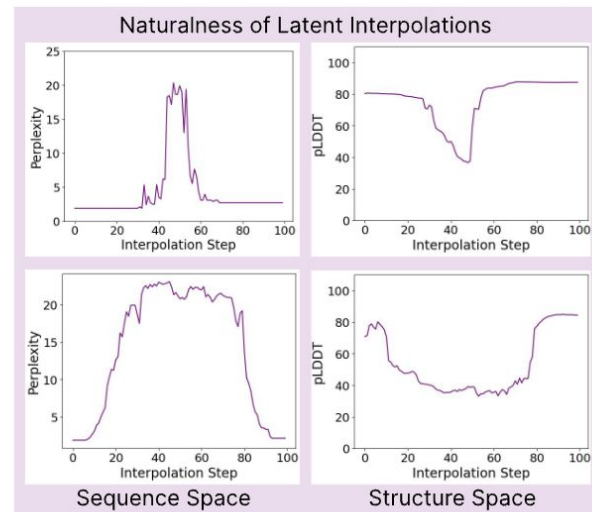
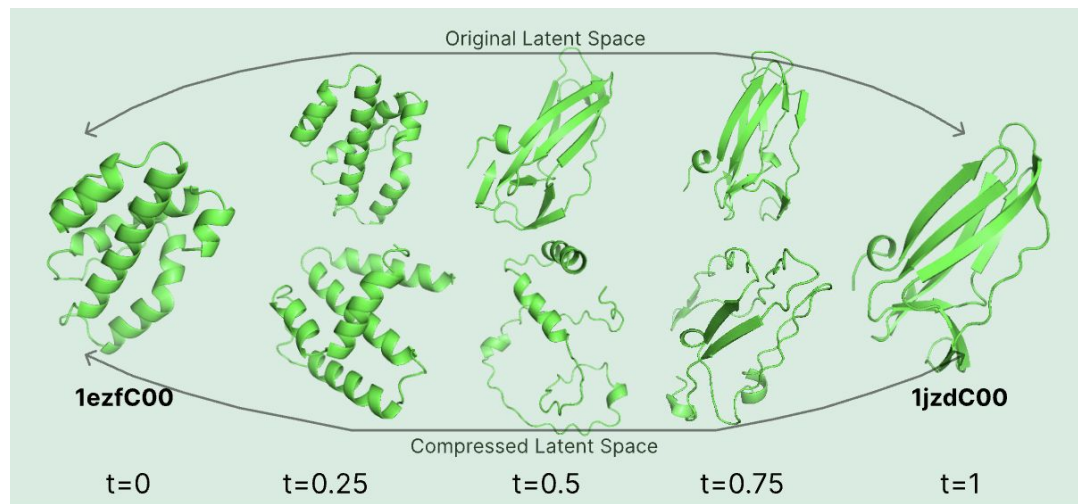
→ BMW S1000RR: 188 mph



# Linear interpolation in the latent space

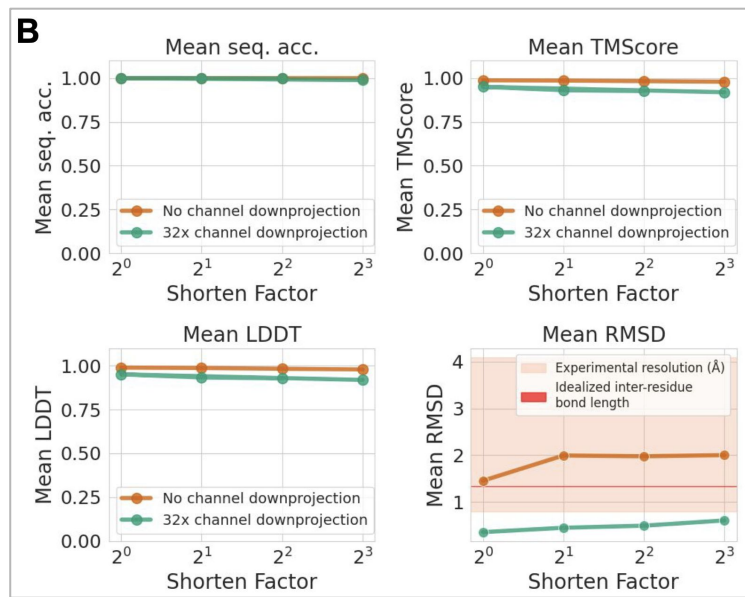
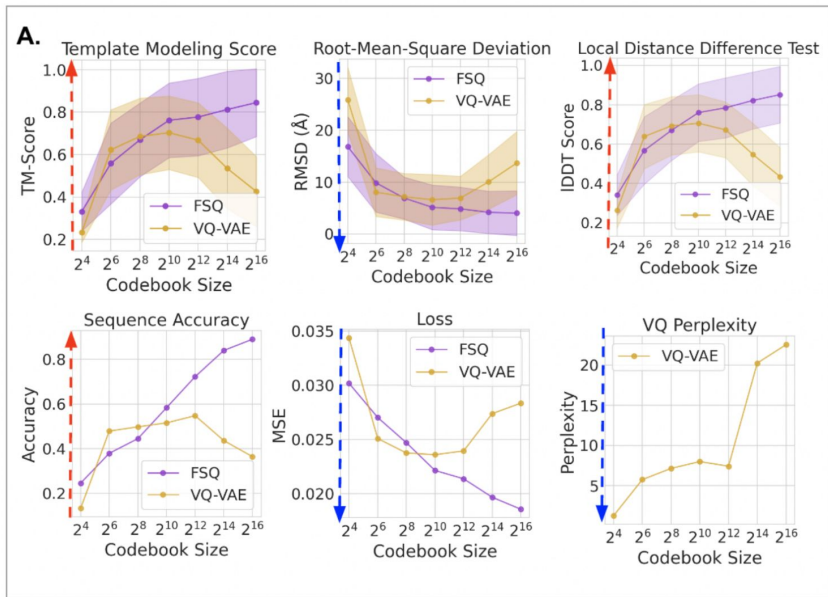


# Linear interpolation in the latent space

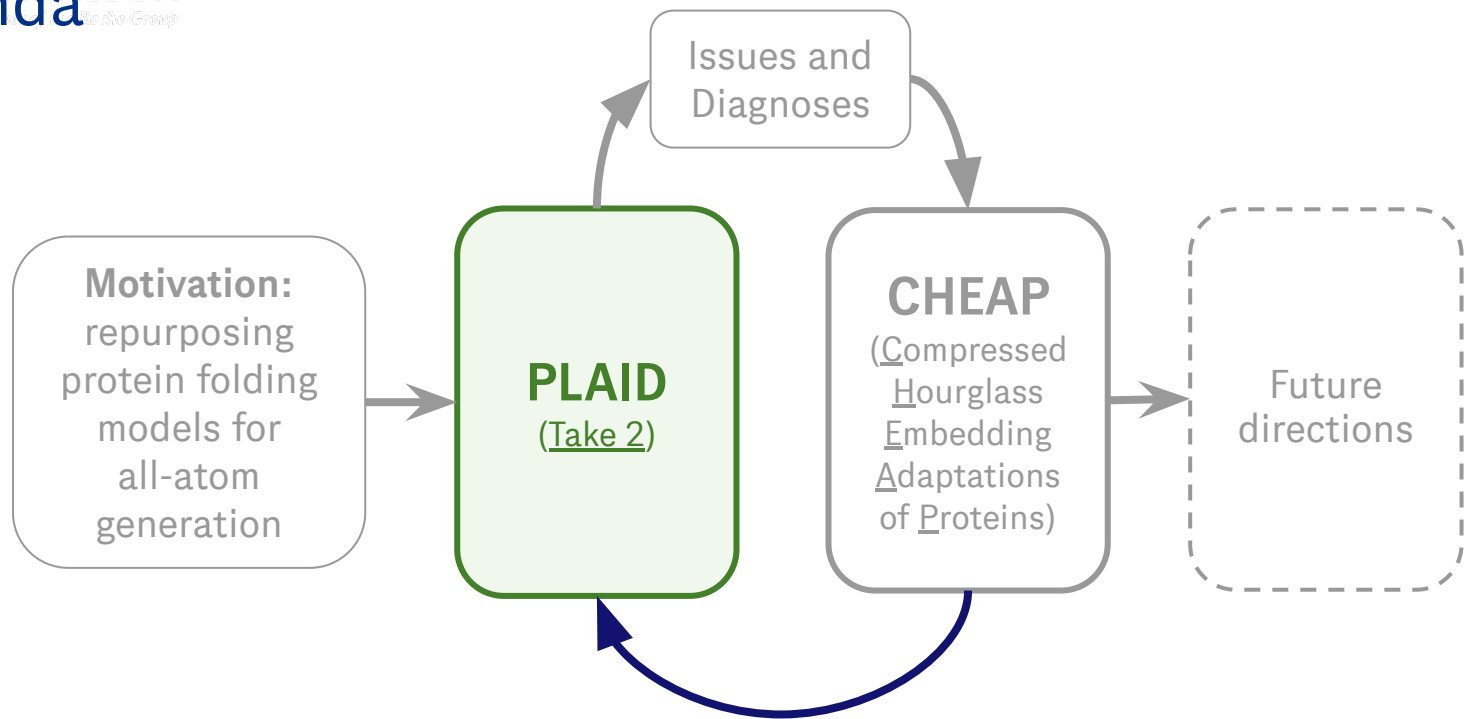


Protein language model latent spaces are less rugged than true fitness landscapes!

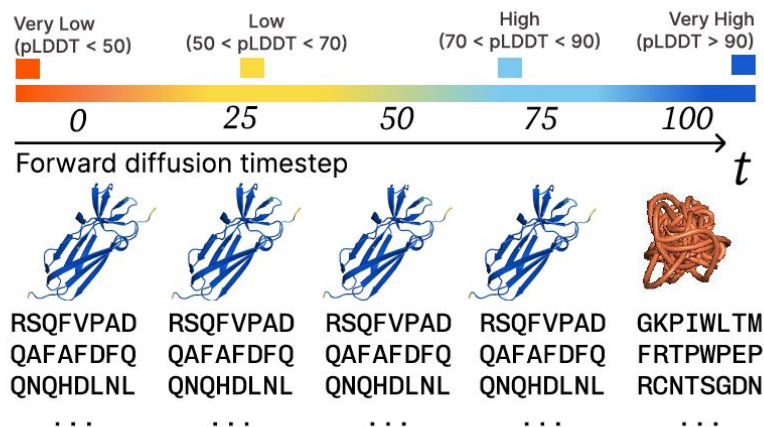
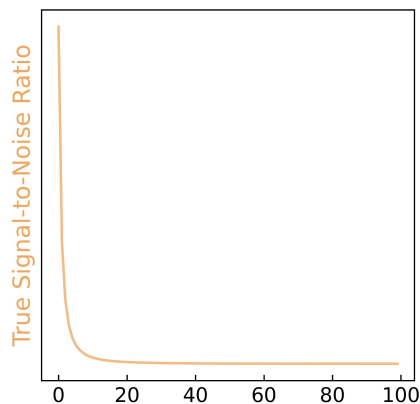
# More results: [bit.ly/cheap-proteins](https://bit.ly/cheap-proteins)



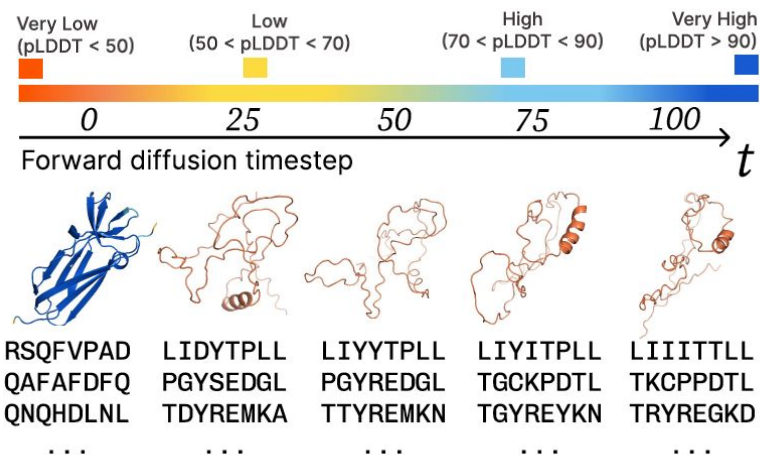
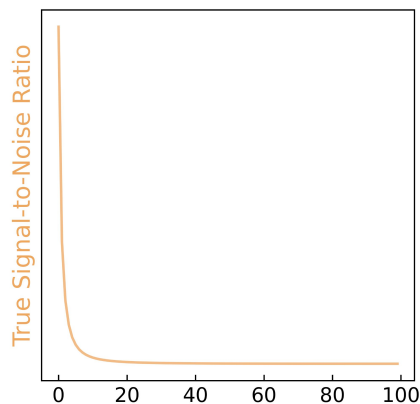
# Agenda



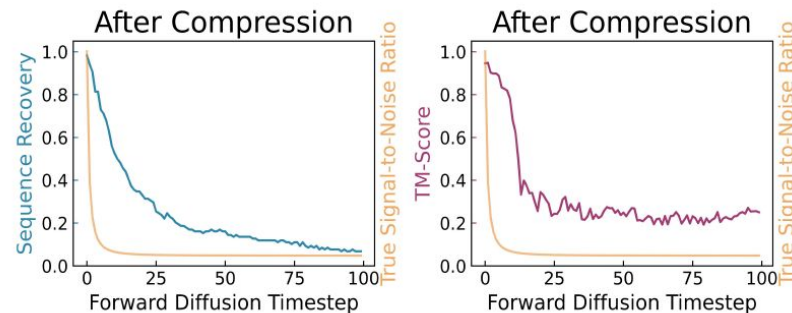
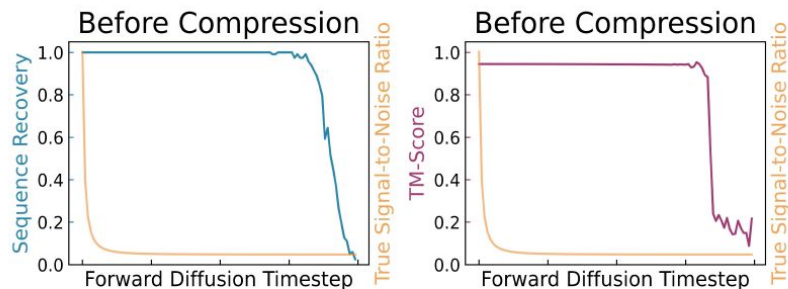
# Noising the **original** latent space does not affect the structure...



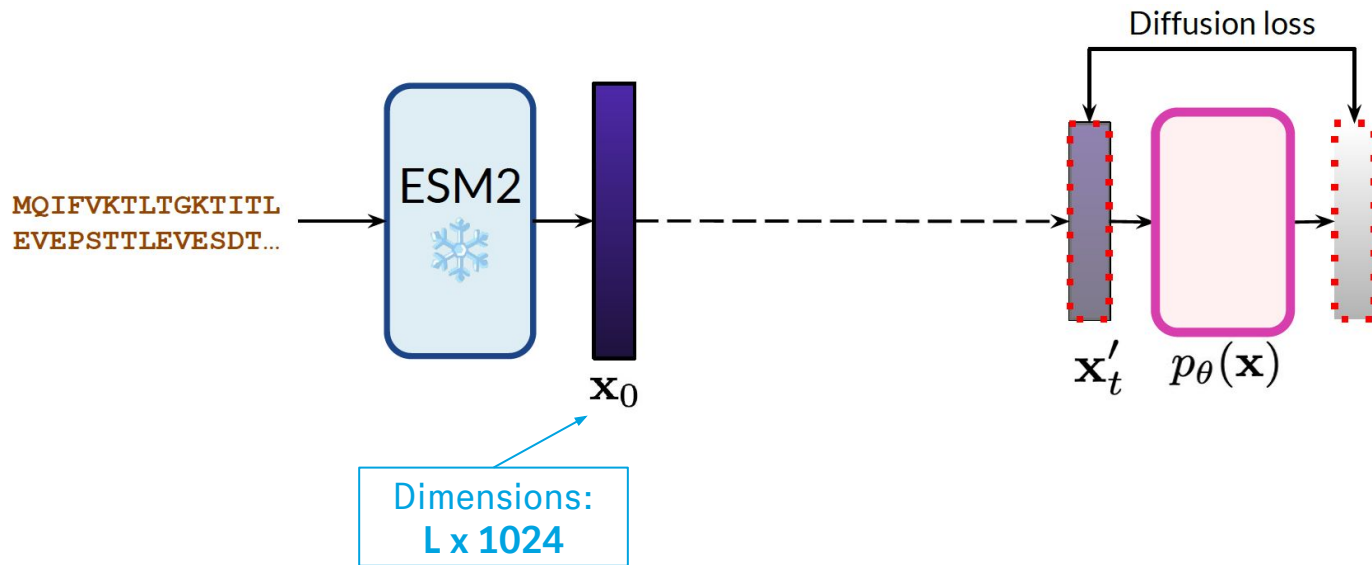
...noising the compressed latent space does map to corrupted structures



...noising the **compressed** latent space does map to corrupted structures

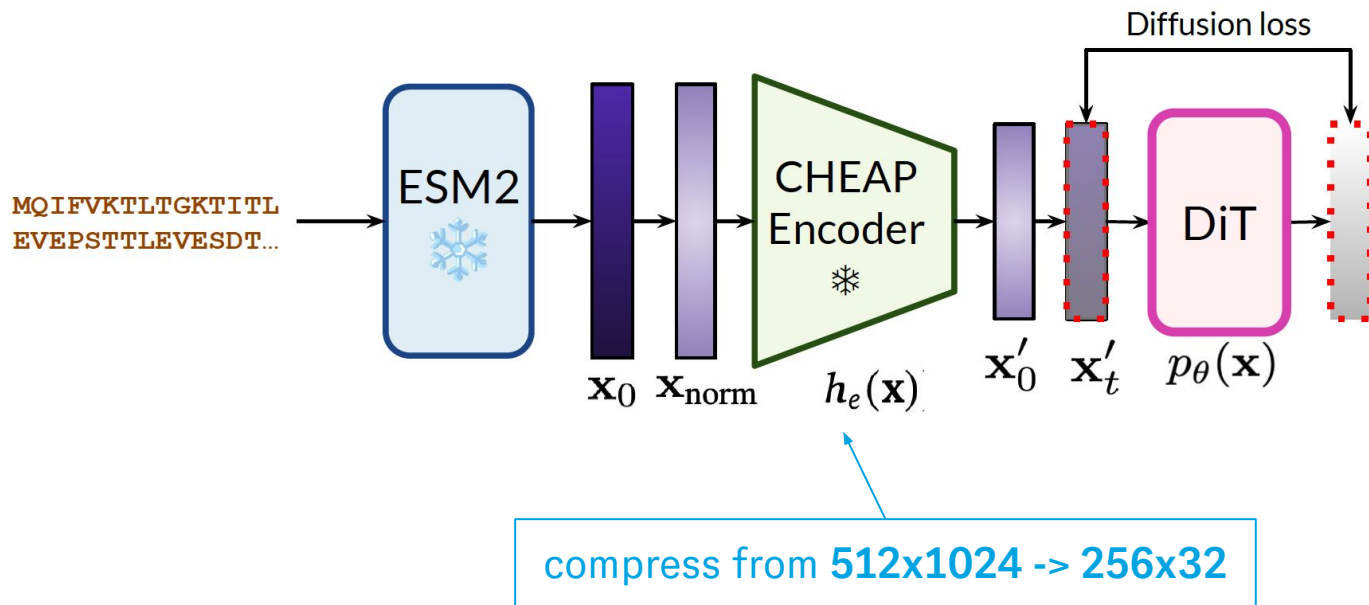


# Training the PLAID latent diffusion model...

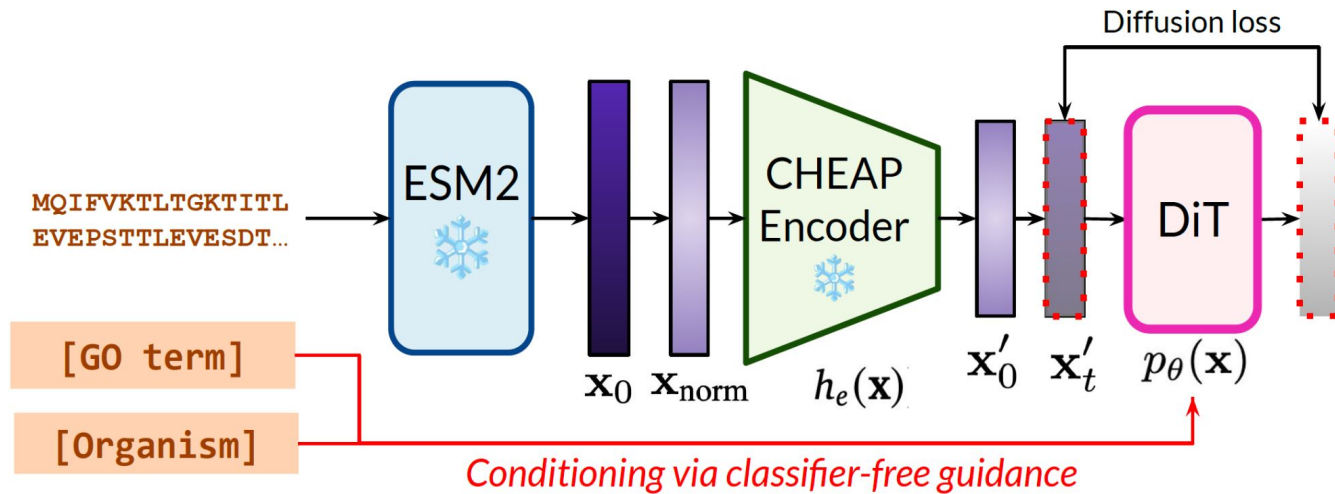




...but add embedding compression with CHEAP

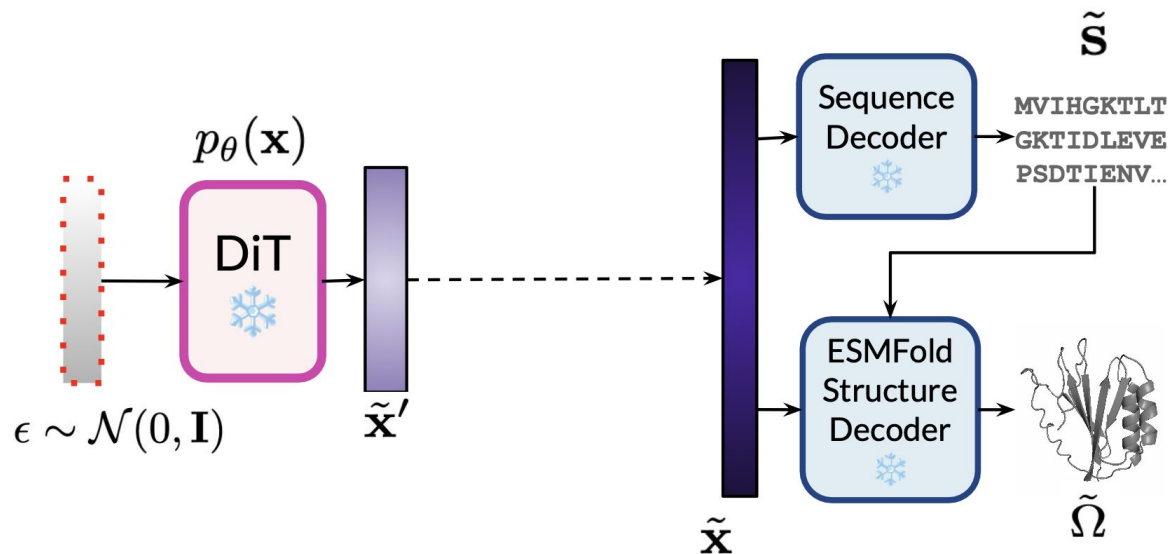


# Adding compositional function + taxonomic conditioning

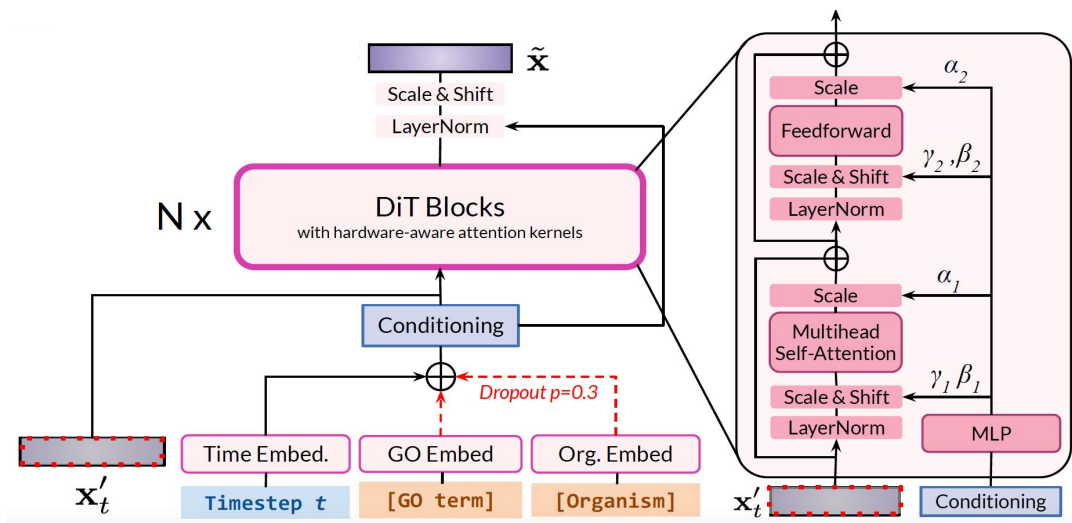


Sequence databases have more sample-annotation pairs!

# PLAID Inference with CHEAP decoder

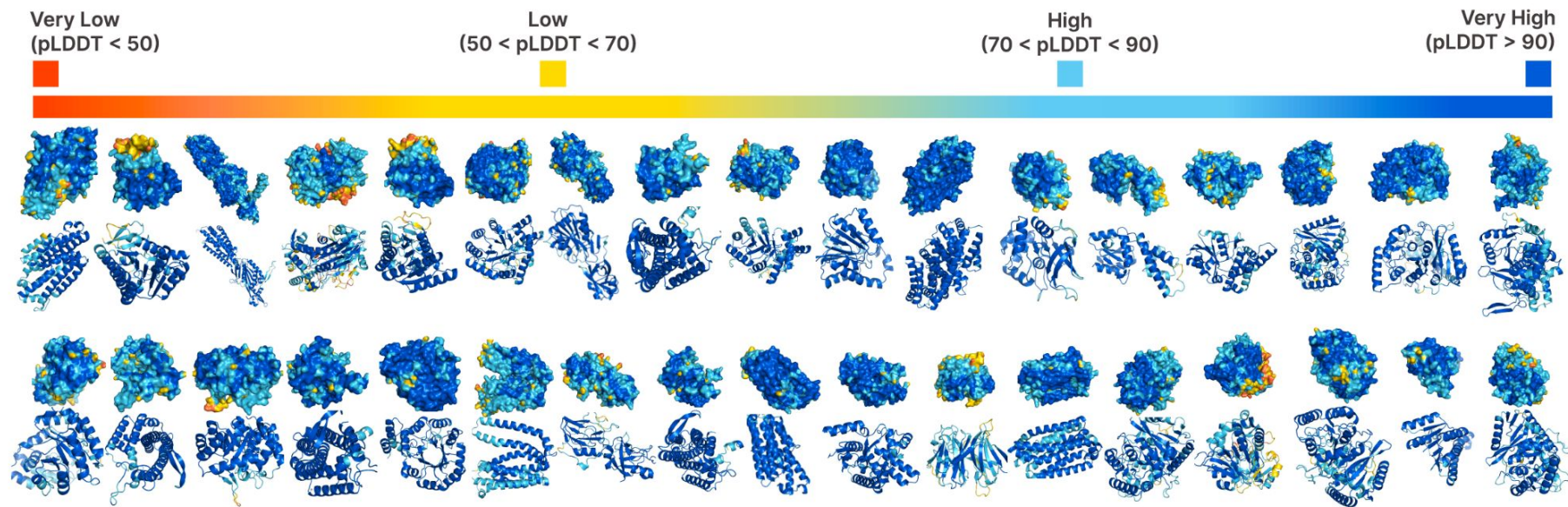


# Scalable Architecture using Diffusion Transformers



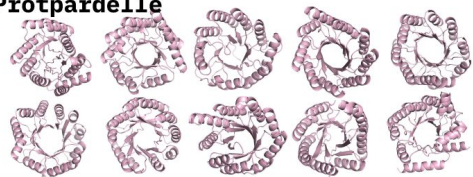
Latent diffusion allows for flexible architecture choice.

# PLAID unconditionally generates diverse all-atom structures



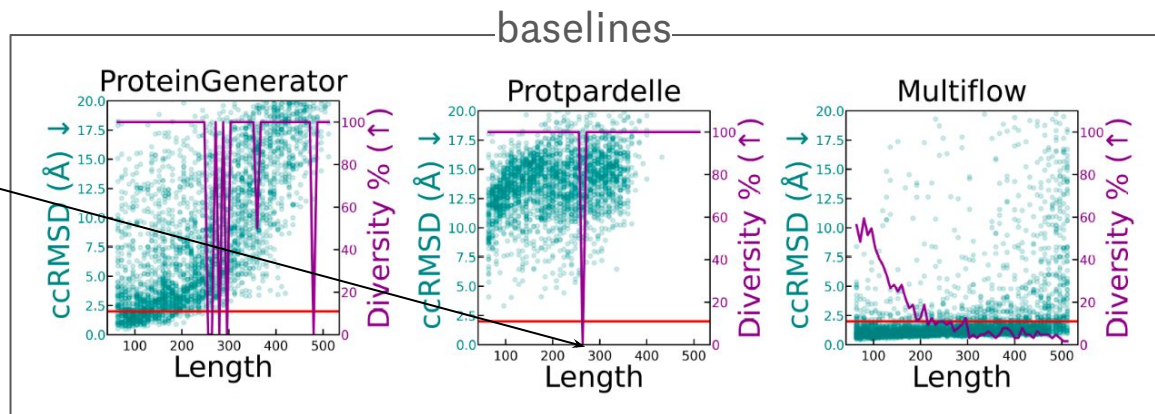
# PLAID unconditionally generates diverse, high-quality folds

## Protpardelle



## Protpardelle

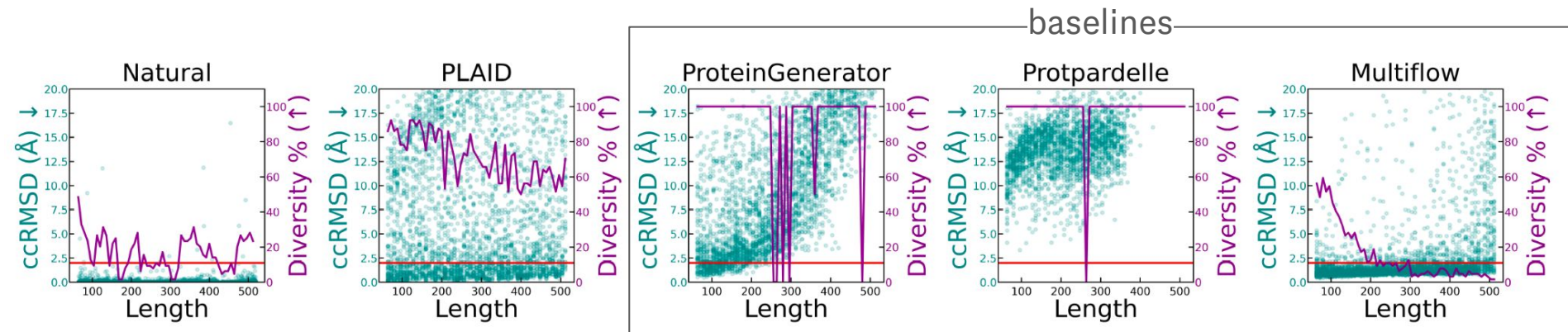
```
>len600_samp97
AGGGGGGGGGGGGGGGGGGGGLGLGLLLPPAGL...
>len600_samp98
PPPPGGAGGGGAAAAAGGSPGGPPGGGGGGGGGGGG...
>len600_samp99
PPGPALPPSPGGVPPPPPLPPPLPGGAPPAGGGLL...
```



**purple: diversity (↑)**  
(# of foldseek clusters /  
# of samples)

**teal: quality (↓)**  
(ccRMSD between generated structure and  
predicted structure of generated sequence)

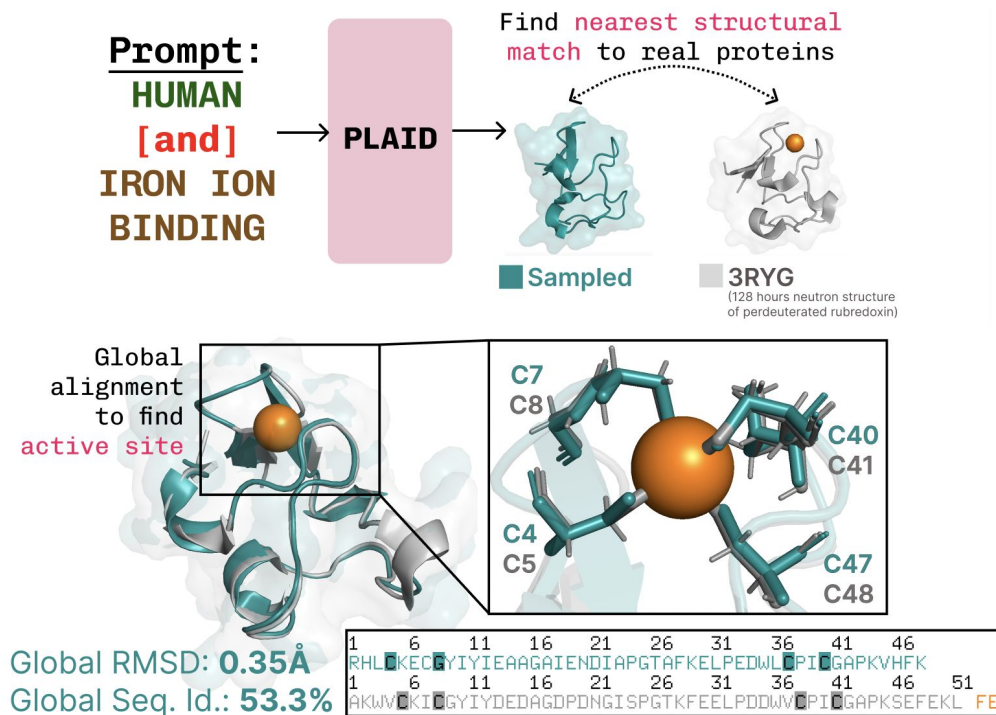
# PLAID unconditionally generates diverse, high-quality folds



PLAID better balances diversity and quality, especially at longer sequence lengths.



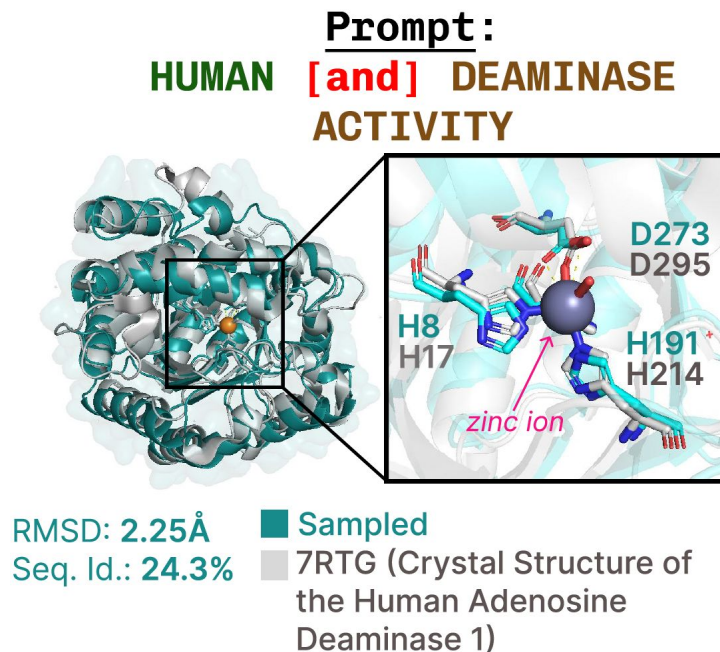
# Function-prompted generations learn active site sidechains



PLAID not only learns that cysteines coordinate the iron ion, but also the sidechain positioning...



# Function-prompted generations learn active site sidechains

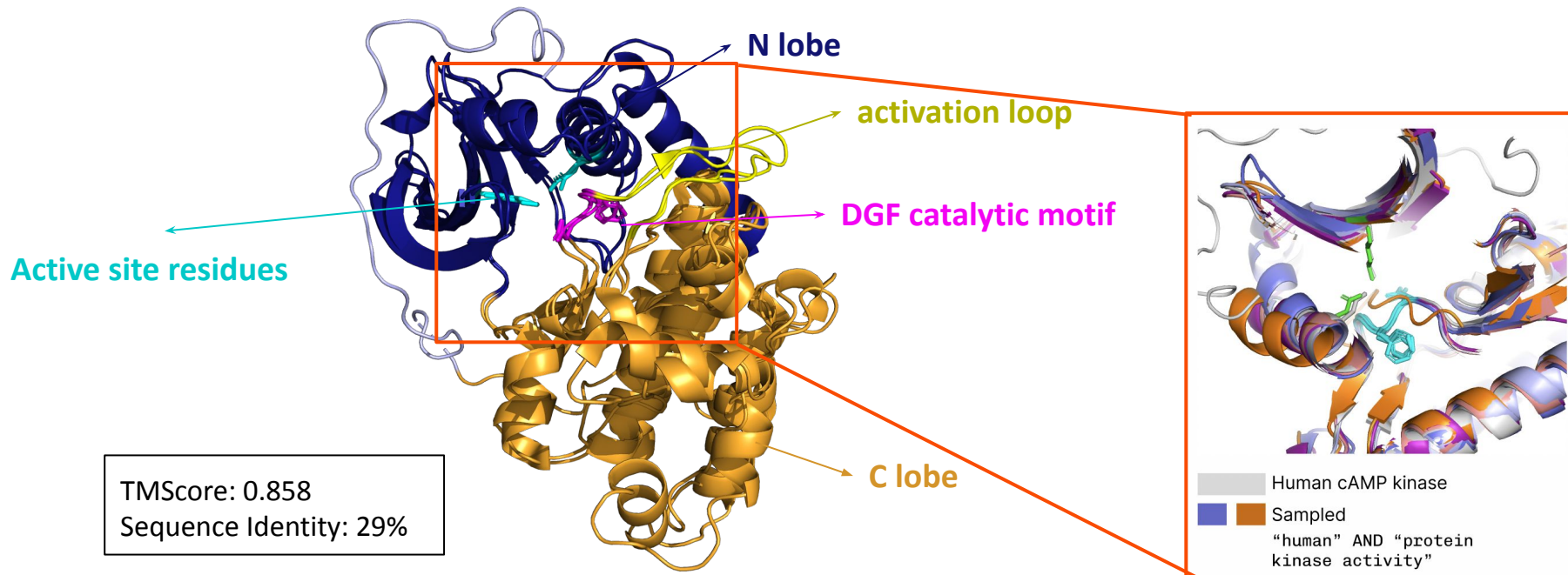


...despite these key residues not being adjacent in the sequence.

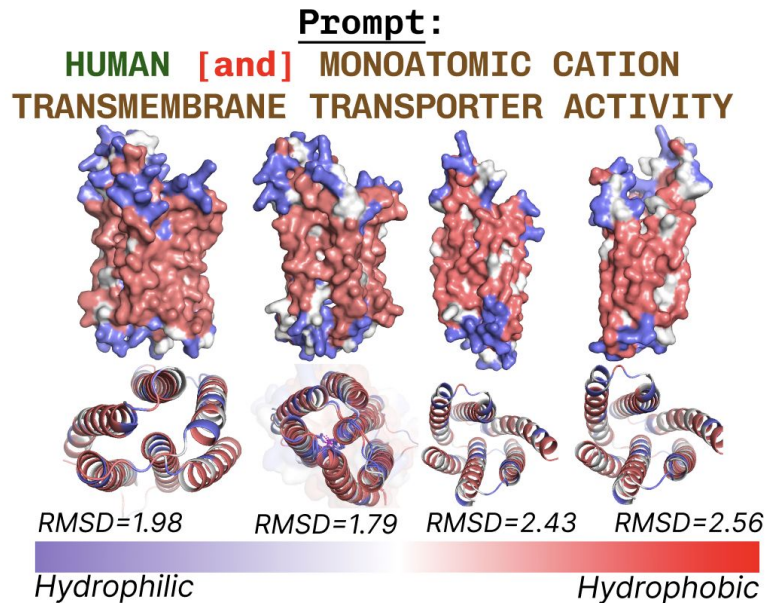
# Examining active site conservation

prompt: "human" AND "protein kinase activity"

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



# Transmembrane proteins exhibit expected hydrophobicity patterns

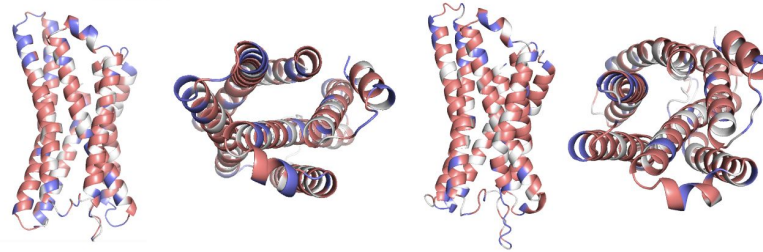
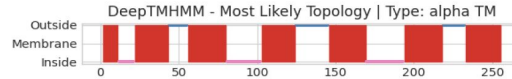


Hydrophobic residues are found at the core, as expected.

# Transmembrane proteins exhibit expected numbers of helices

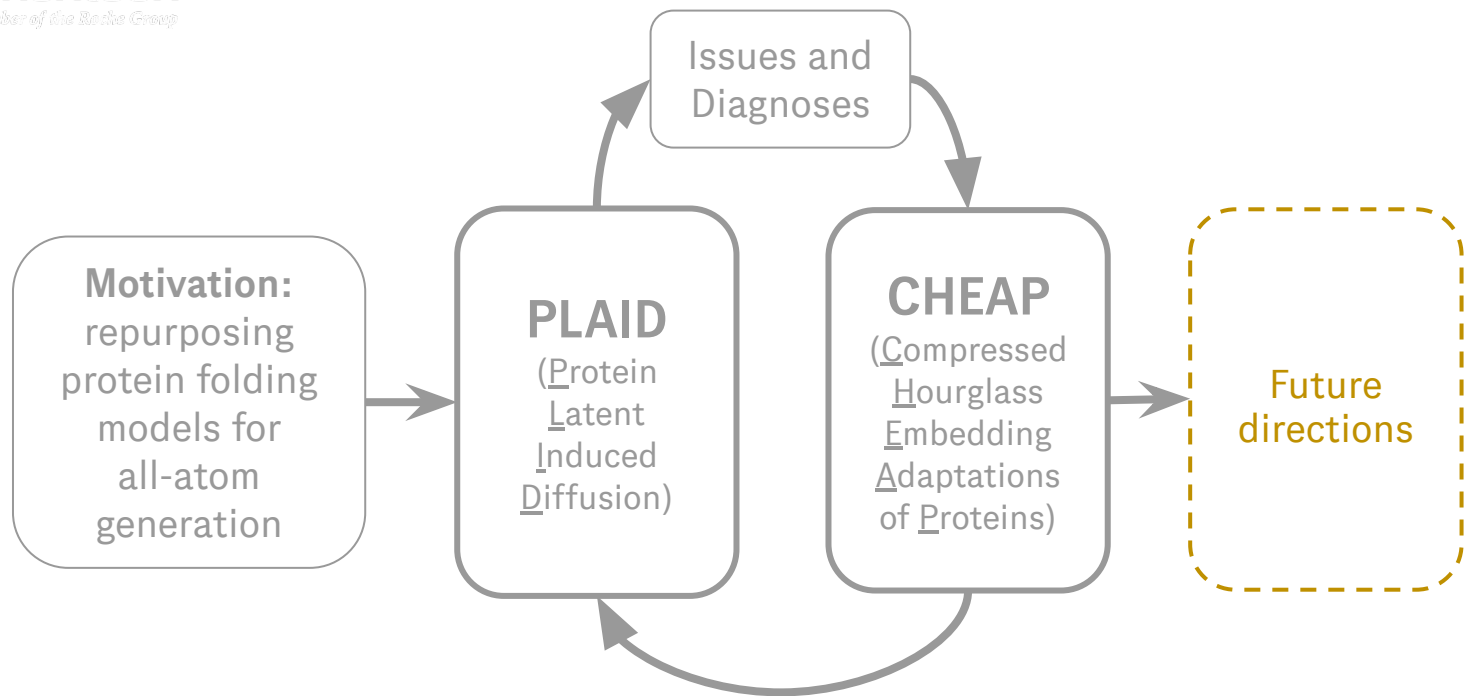
**Prompt:**  
**HUMAN [and]**  
**G PROTEIN-COUPLED RECEPTOR ACTIVITY**

GNVLVLIMMILKQREVKSMPNV  
WVFNLAISDLLFLLSTPLL VVK  
MSDTSWNLGLSPCKITTFLFL  
NLYSSVFFLAQLSLDRYLTVRQ  
VRSN . . .

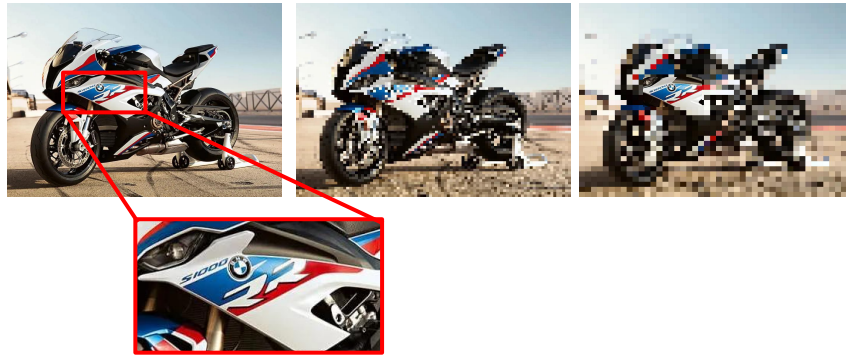


GPCRs have the expected 7-transmembrane topology, both when analyzing the sequence and structure.

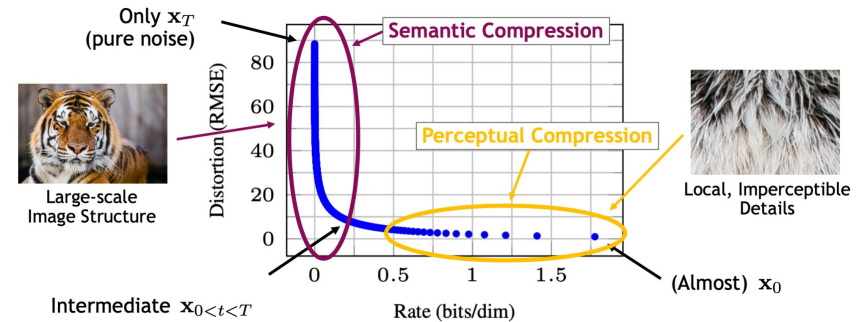




# Biological data selection from an information theoretic perspective



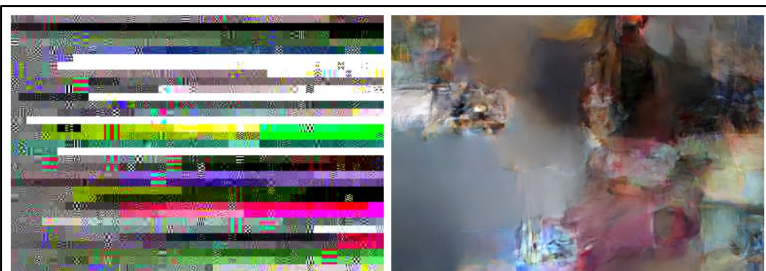
...what level of compression is optimal?



what constitutes semantic vs. perceptual compression for proteins? what level of detail do we need for drug discovery?



# Connection between CHEAP, PLAID, and lossy compression



**Figure 1.1:** Compression as generative modeling. *Left:* A sample drawn from the probabilistic model underlying JPEG, which betrays an assumption of independence among neighboring 8 by 8 pixel blocks (except for the DC components within each row). *Right:* A sample generated by a recent neural compression model by Minnen *et al.* [132].

Representation learning => data  
compression => generative models!

## Compression with Flows via Local Bits-Back Coding

Jonathan Ho  
UC Berkeley  
jonathanho@berkeley.edu

Evan Loh  
UC Berkeley  
evan.lohn@berkeley.edu

Pieter Abbeel  
UC Berkeley, covariant.ai  
pabbeel@cs.berkeley.edu



## Denoising Diffusion Probabilistic Models

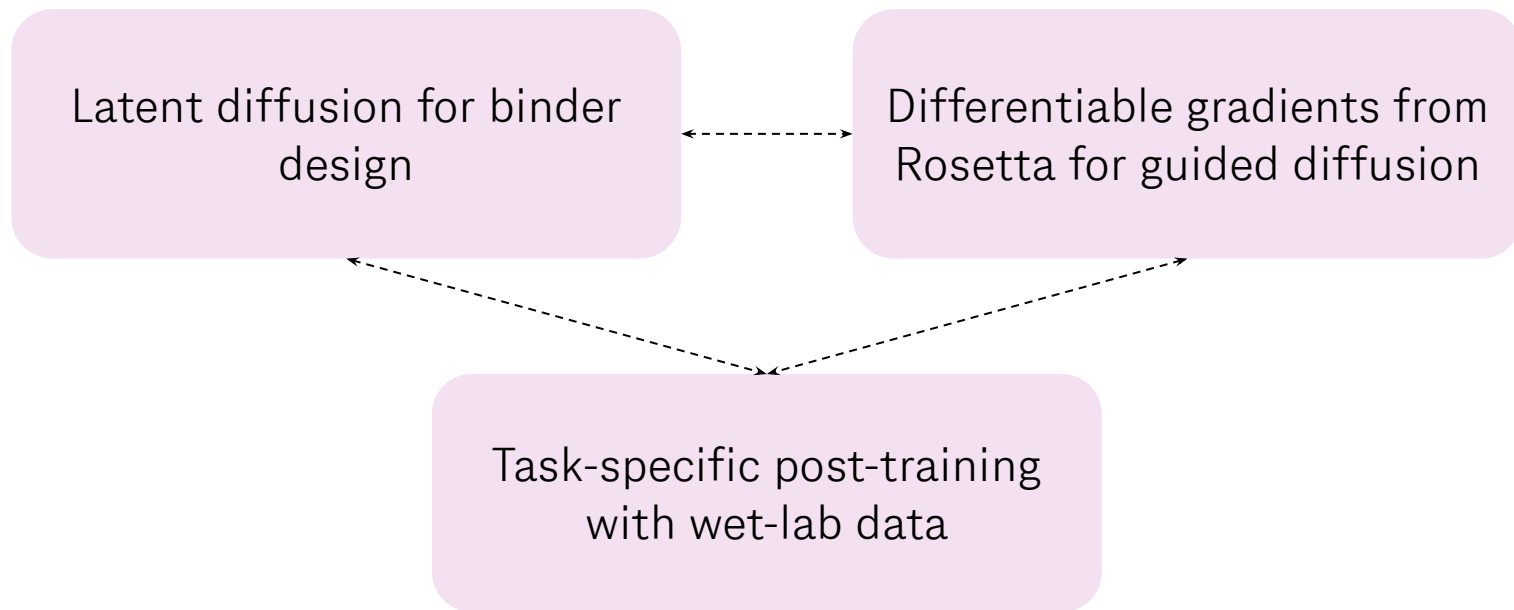
Jonathan Ho  
UC Berkeley  
jonathanho@berkeley.edu

Ajay Jain  
UC Berkeley  
ajayj@berkeley.edu

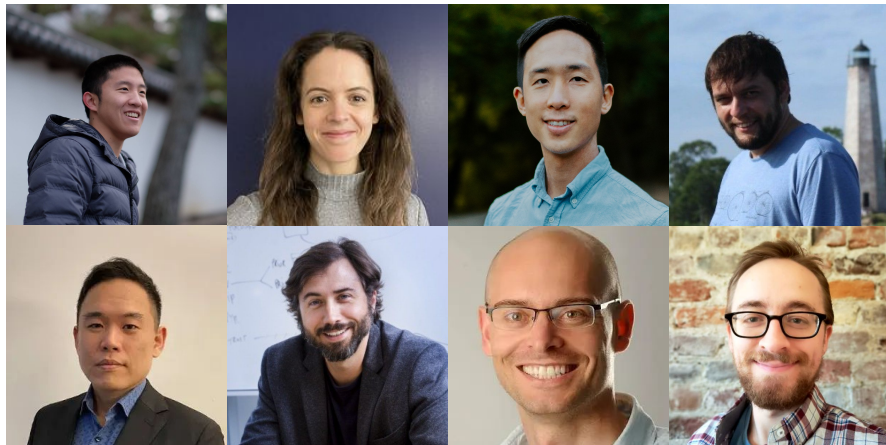
Pieter Abbeel  
UC Berkeley  
pabbeel@cs.berkeley.edu



# Future directions



# Thanks!



## Berkeley

Amy X. Lu  
Wilson Yan  
Pieter Abbeel

## Microsoft Research

Kevin Yang

## Prescient Design

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



@amyxlu



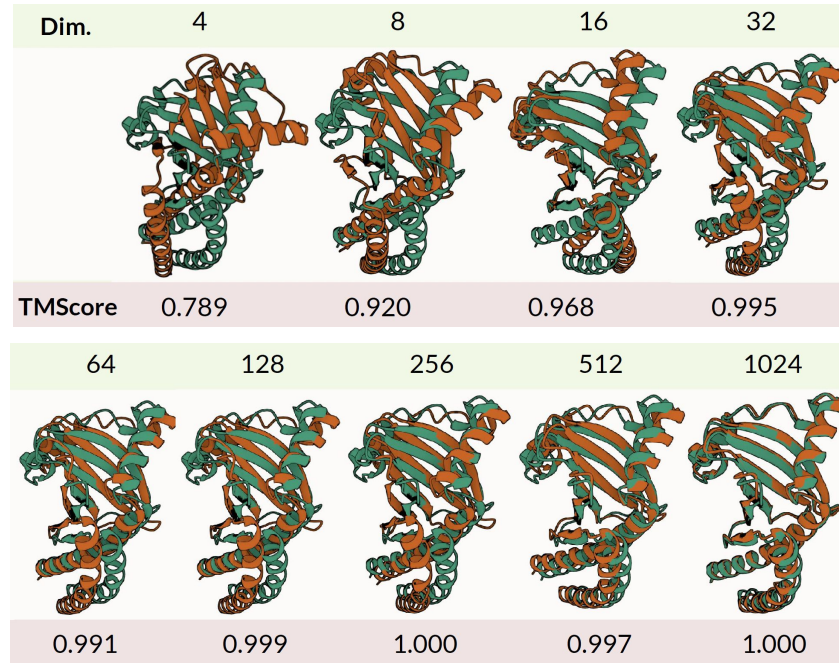
amyxlu.github.io



amyxlu@berkeley.edu

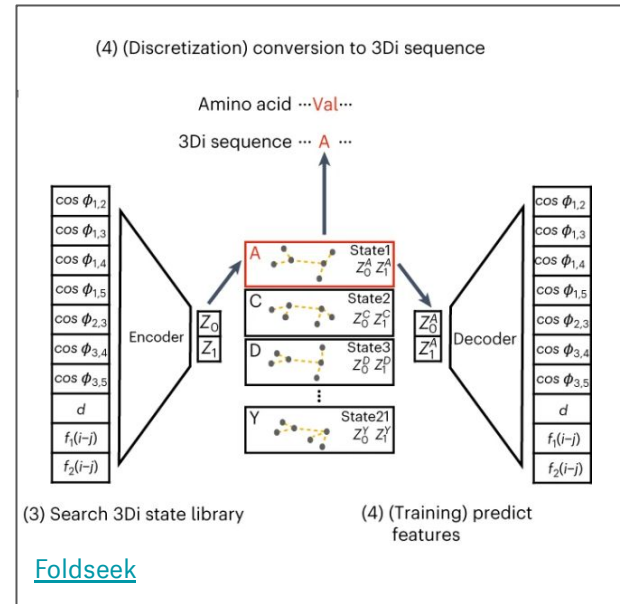
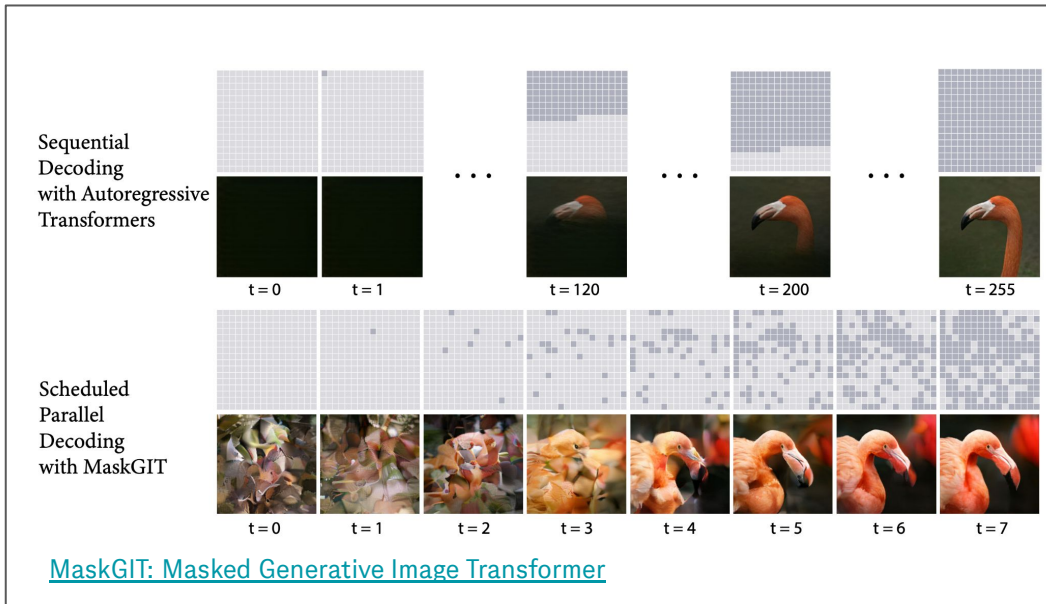
# Appendix

# Secondary structure is retained even at 256x compression

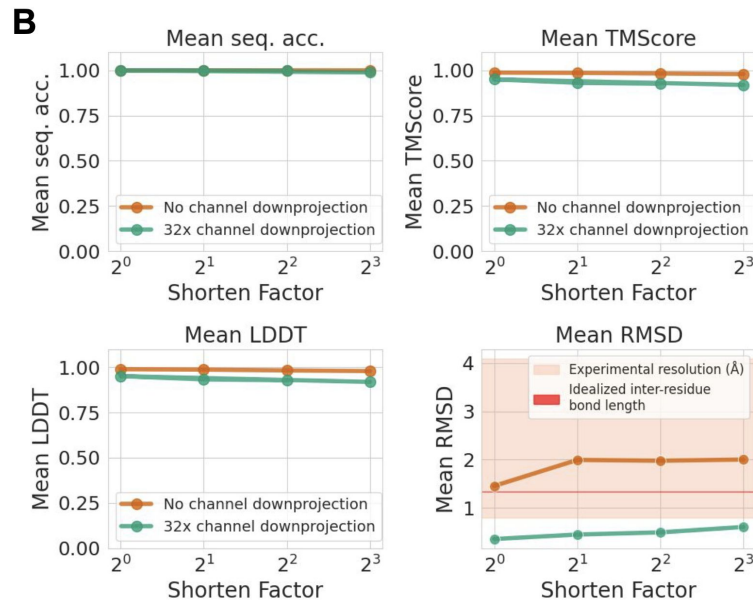


# Side note: why tokenized representations?

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

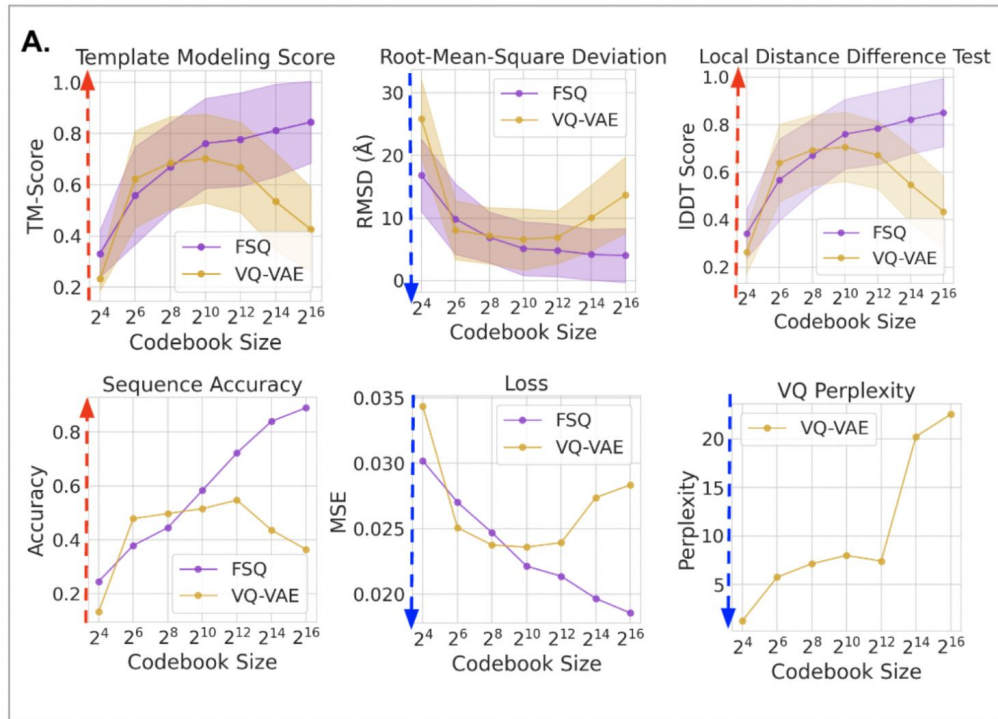


We can compress lengthwise and channelwise:

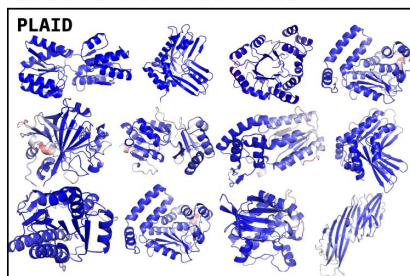
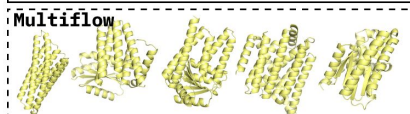
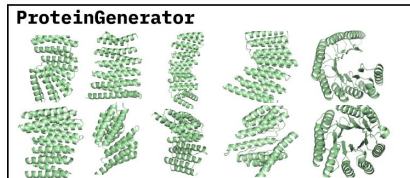
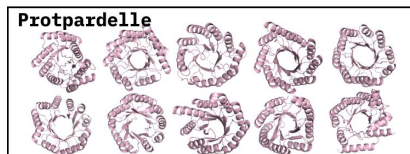


What does this mean for how structural information is shared across residue positions?

# All-atom structural tokenizer, obtained from sequence alone



# PLAID unconditionally generates diverse, high-quality folds



## Protpardelle

```
>len600_samp97
AGGGGGGGGGGGGGGGGGGGGGGGGGGGGLGLGLLLPPAGL...
>len600_samp98
PPPPGGAGGGGAAAALAGGSPGGPPGGGGGGGGGGGG...
>len600_samp99
PPGPALPPSPGPGGVPPPPPLPPPLPGGAPPAGGGLL...
```

## ProteinGenerator

```
>len600_000097
GAAGLTAAAVVGAAAAAGAAAAALAAAGAGAAAA...
>len600_000098
AGAAAGAAAAAGAAAAAGAGGGAGGAAAAAG...
>len600_000099
VAAQAVQGAIAAAAAAATAALGLTAAGIAAPLLALV...
```

## Multiflow

```
>len600_sample_97
LLGGLLGGLLGGAAGGAGAGAAAAGGGAVGVVAGAVT...
>len600_sample_98
ADAATLTVGGGGTGGGGGAGGALGGAAGGGGRVTLVV...
>len600_sample_99
AGGGAGLAGGAGGAGGAAAAAGAGGGGAAAA...
```

## PLAID

```
>len600_sample97
PDMGTVLGLAHSVGHLDKTPDLSVADLETNLALLAAH...
>len600_sample98
FEMFDDKGGDLWERAASSGQLLIDVAYLANGLRDGAT...
>len600_sample99
NGGQARGTDDPLTHALQTLFQSAALDQSLQGDPENAV...
```



# Examining sampling hyperparameters

