# Generative Models for Real-World Drug Discovery

**Amy X. Lu**
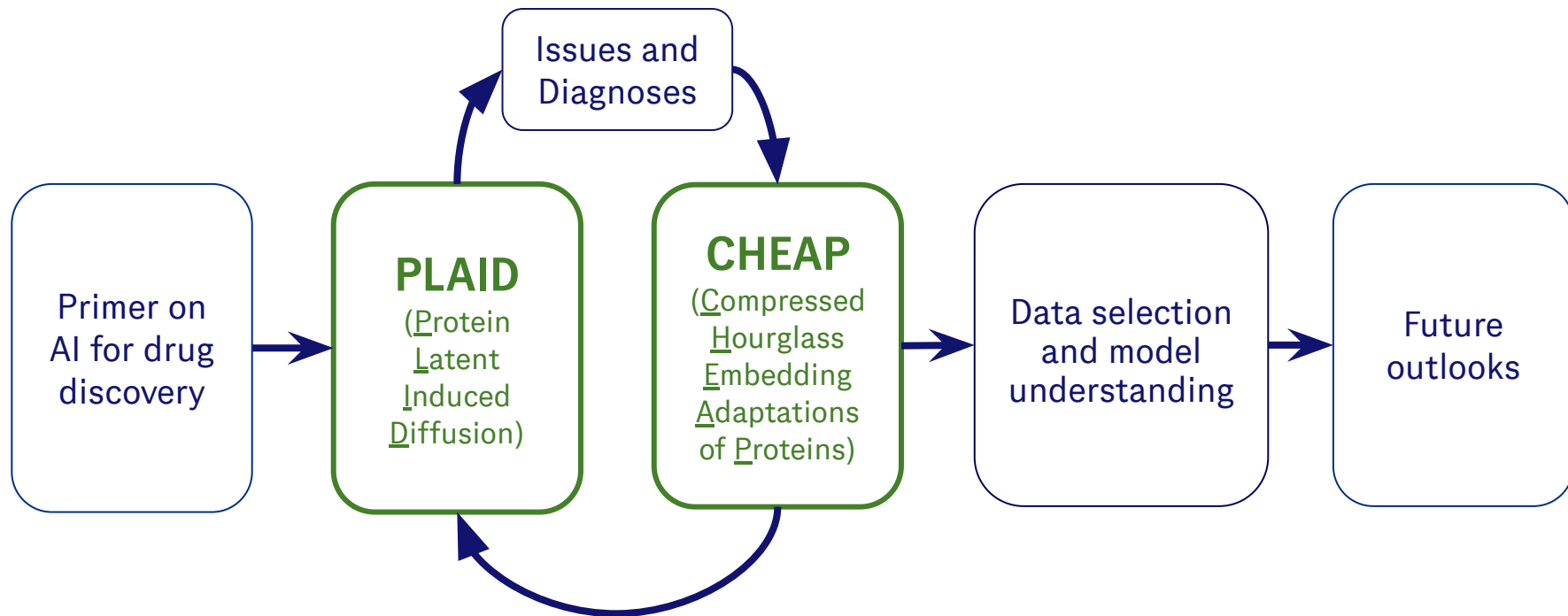April 30th, 2025
PhD Dissertation Talk
BAIR Seminar

→ biology as a data modality for generative modeling
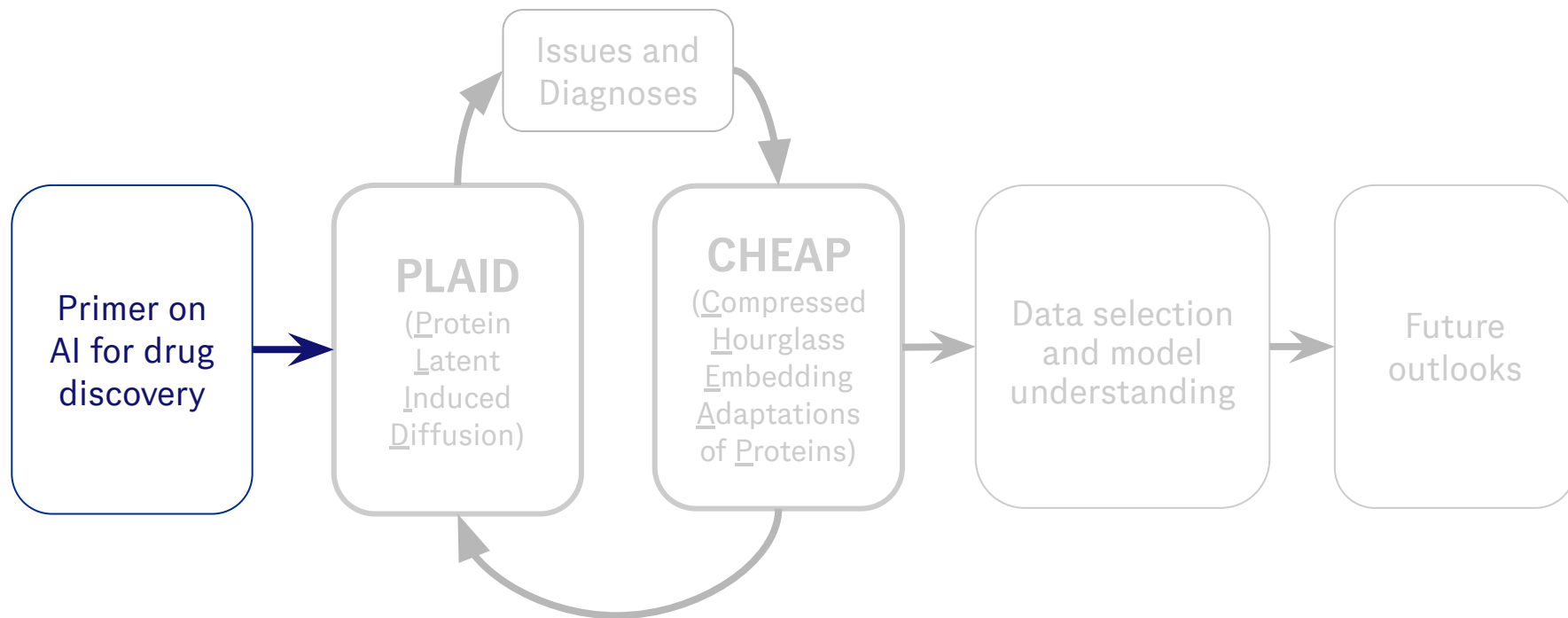
# Generative Models for Real-World Drug Discovery

→ evaluations/tasks anchored around drug discovery for protein design

# Agenda

# Agenda

# A simplified look at modern drug discovery…



**GLOBAL HEALTH**

## From Jan. 2020: China Identifies New Virus Causing Pneumonialike Illness

The new coronavirus doesn't appear to be readily spread by humans, but researchers caution that more study is needed.

# A simplified look at modern drug discovery...



(disease identification)

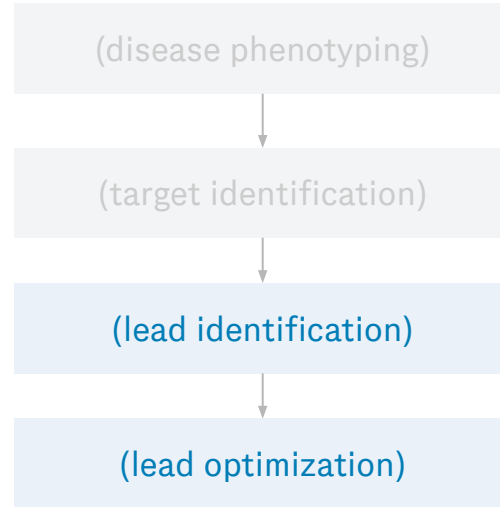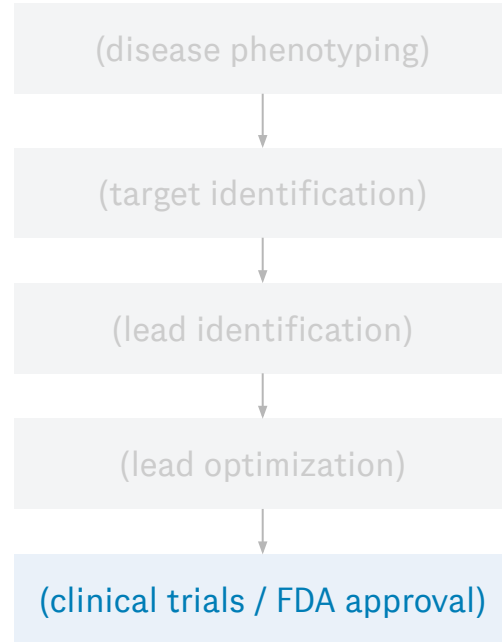# A simplified look at modern drug discovery...



(disease identification)
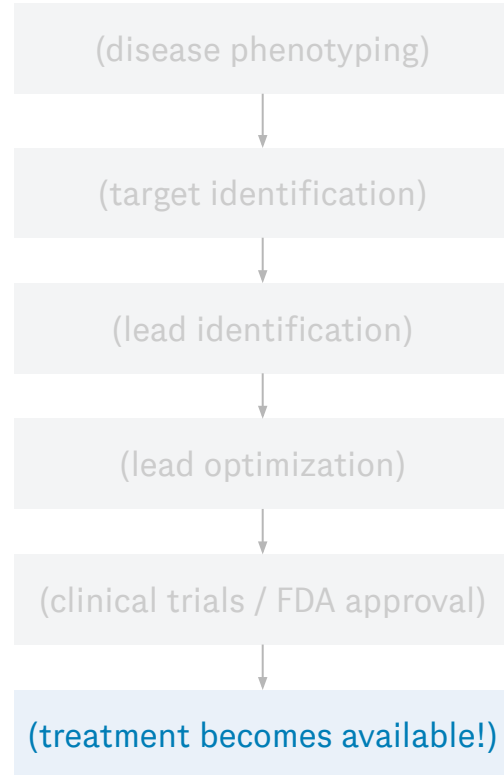
↓

(target identification)

# A simplified look at modern drug discovery…



(disease phenotyping)

(target identification)
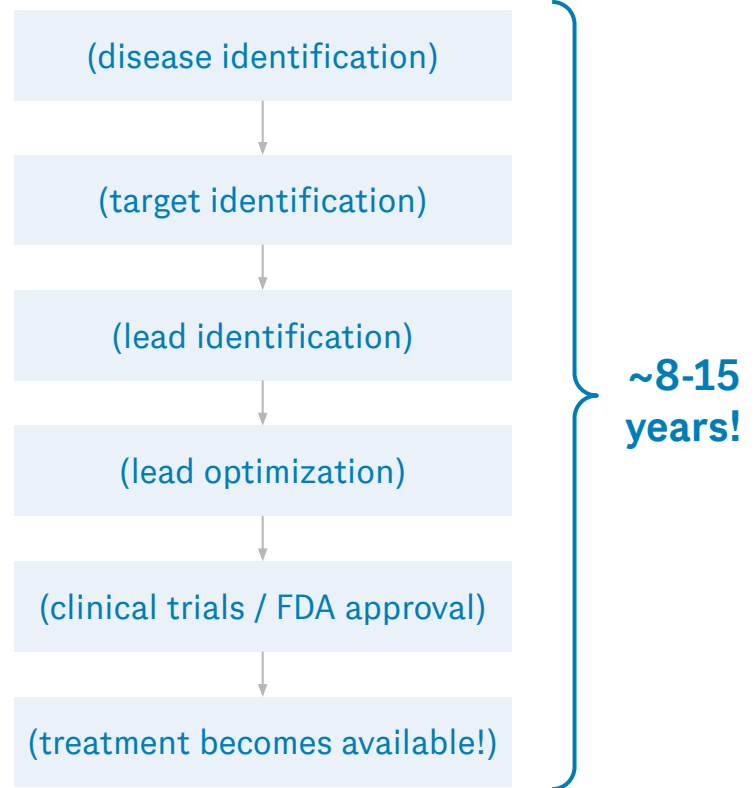
(lead identification)

(lead optimization)

# A simplified look at modern drug discovery...

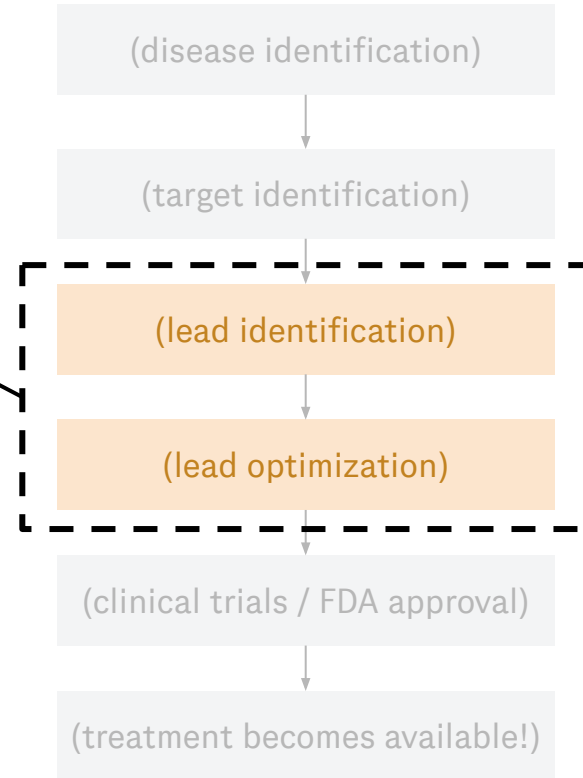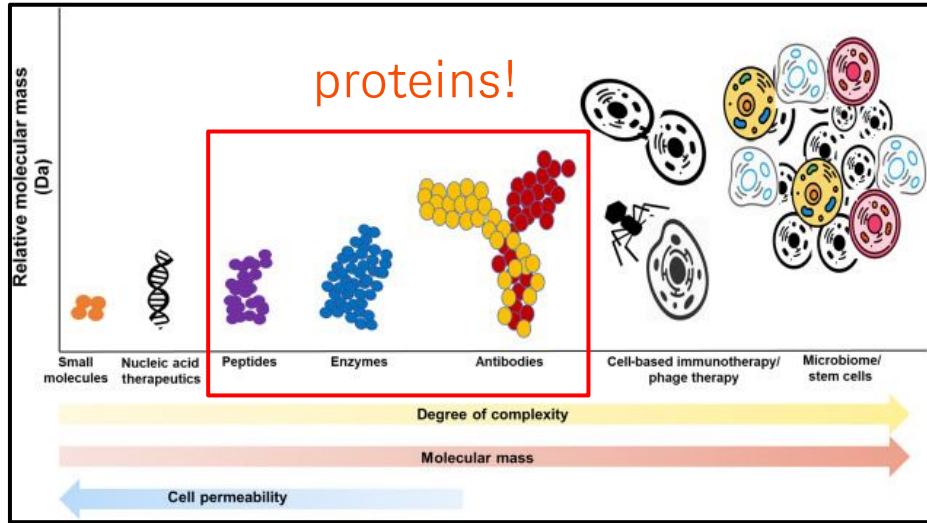# A simplified look at modern drug discovery...



(disease phenotyping)

↓

(target identification)

↓

(lead identification)

↓

(lead optimization)

↓

(clinical trials / FDA approval)

↓

(treatment becomes available!)

# Drug discovery is time-consuming



(disease identification)

↓

(target identification)

↓

(lead identification)

↓

(lead optimization)

↓

(clinical trials / FDA approval)

↓

(treatment becomes available!)

**~8-15 years!**

# Accelerating drug discovery with AI?



(disease identification)

(target identification)

(lead identification)

(lead optimization)

(clinical trials / FDA approval)

(treatment becomes available!)

# Accelerating protein design with AI?



proteins!

(disease identification)

(target identification)

(lead identification)

(lead optimization)

(clinical trials / FDA approval)
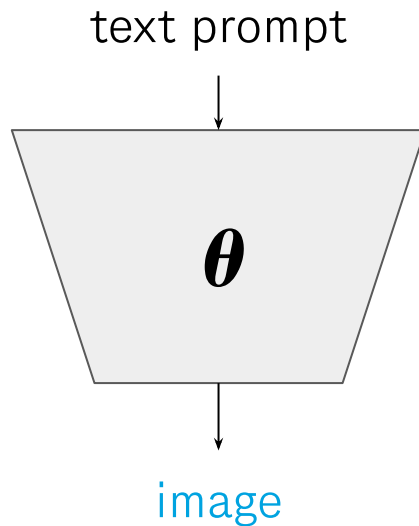
(treatment becomes available!)

?

# Accelerating protein design with AI?

Please generate a comic about the following. Make it as clear as possible what is happening in the scene; conveying clarify is more important than artistic quality. Make it very clear, like it's a children's book and a 5 year old should be able to understand what's going on
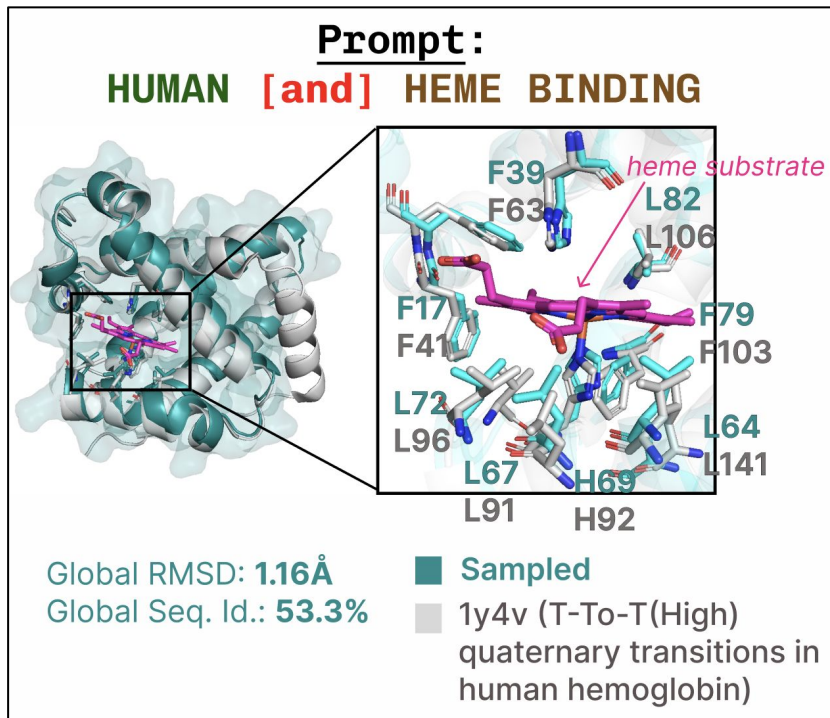
* A car component is loose and there's a weird sound
* we open the car hood to look to identify what's wrong. Turns out there is a loose bolt
* we go through a bunch of tools and hammers and wrenches to find the right tool that will tighten this up
* after trying a bunch, we finally find the right wrench!
* we take the car out for a test drive to make sure we didn't break anything else
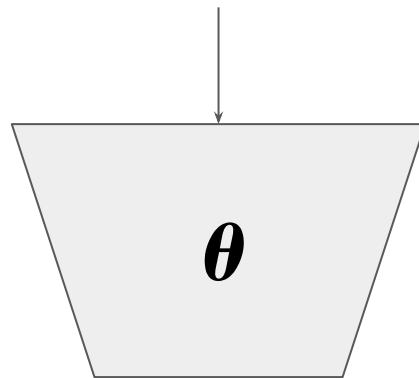* cool, everything works, now we sell this wrench to other people whose cars are also making this sound.

Image created

text prompt

$$\theta$$

image

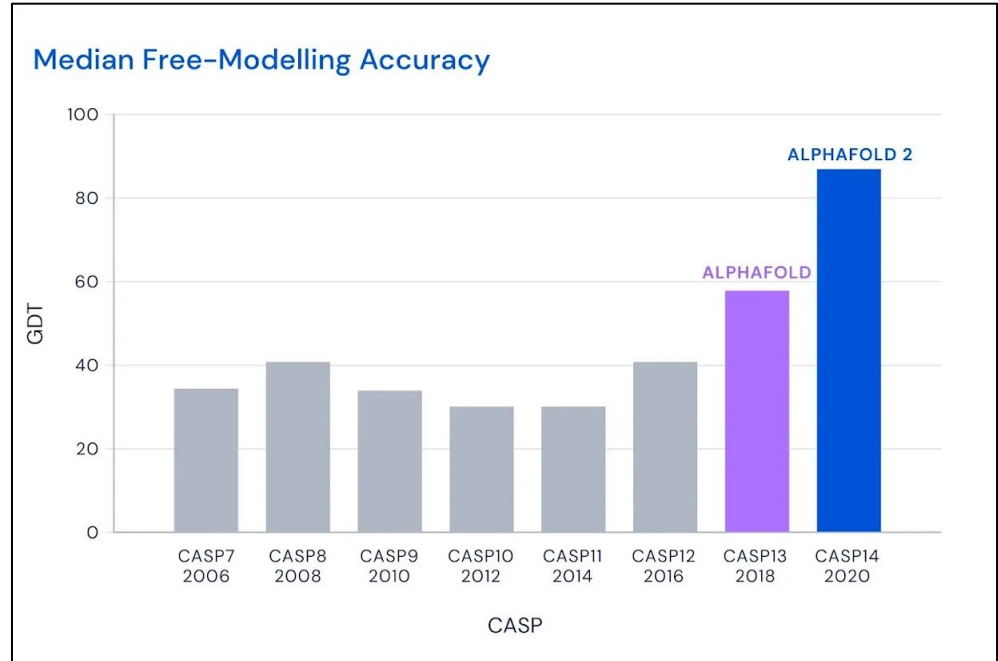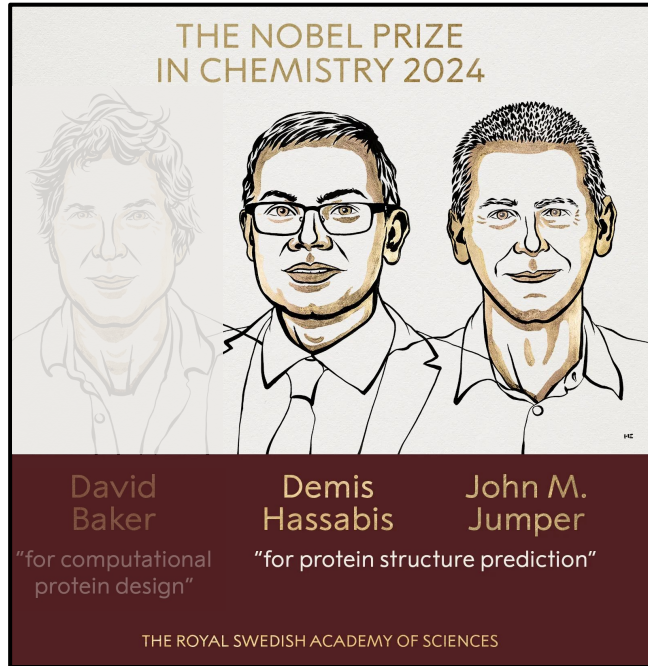# Accelerating protein design with AI?
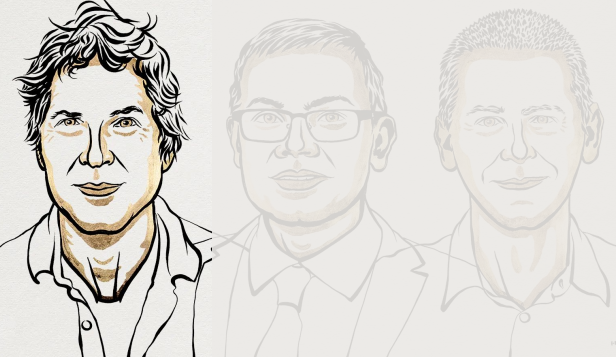
# The potential of deep learning for protein structure prediction

# The potential of deep learning for protein structure generation

# What else might 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 & antigen expression**
→ can we achieve organism specificity?

**Deployment**
→ can we speed up inference to improve "shots on goal"?

**Biosecurity**
→ how should we measure and prevent the potential for dual use?

**Data curation**
→ how should we collect data for model (pre)training in costly acquisition regimes?

other PhD works

# What else might we need for drug discovery?

this talk

**Co-generation**
→ can we simultaneously generate sequence and structure?

**Control**
→ how can we specify complex functions and constraints?

**Immunogenicity**
→ can we achieve organism specificity?

**Compressed protein representations**
*(Cell Patterns, to appear)*

**Latent diffusion for all-atom generation**
*(in submission)*

# What else might we need for drug discovery?

other PhD research: **deployment & model understanding**

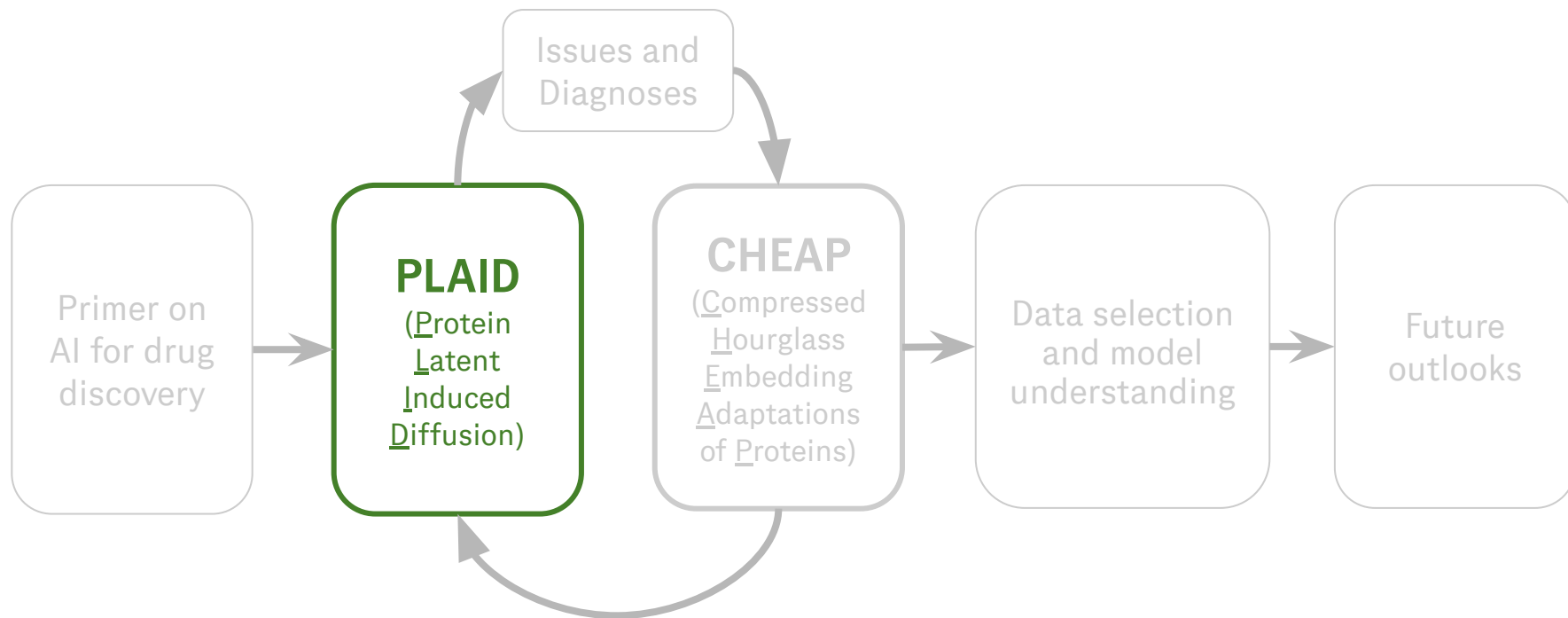| | | | | |
|---|---|---|---|---|
| Model-based optimization for protein engineering (Kolli et al., 2022) | Dense passage retrieval for homology search (Boger et al., 2023) | Guided diffusion with differentiable biophysical energies *(unpublished)* | Effect of training data compositions on protein language model likelihoods *(Gordon et al., 2024)* | Evo2 biosecurity and inference optimization (Brixi et al., 2025) |

**Compressed protein representations**
*(Cell Patterns, to appear)*
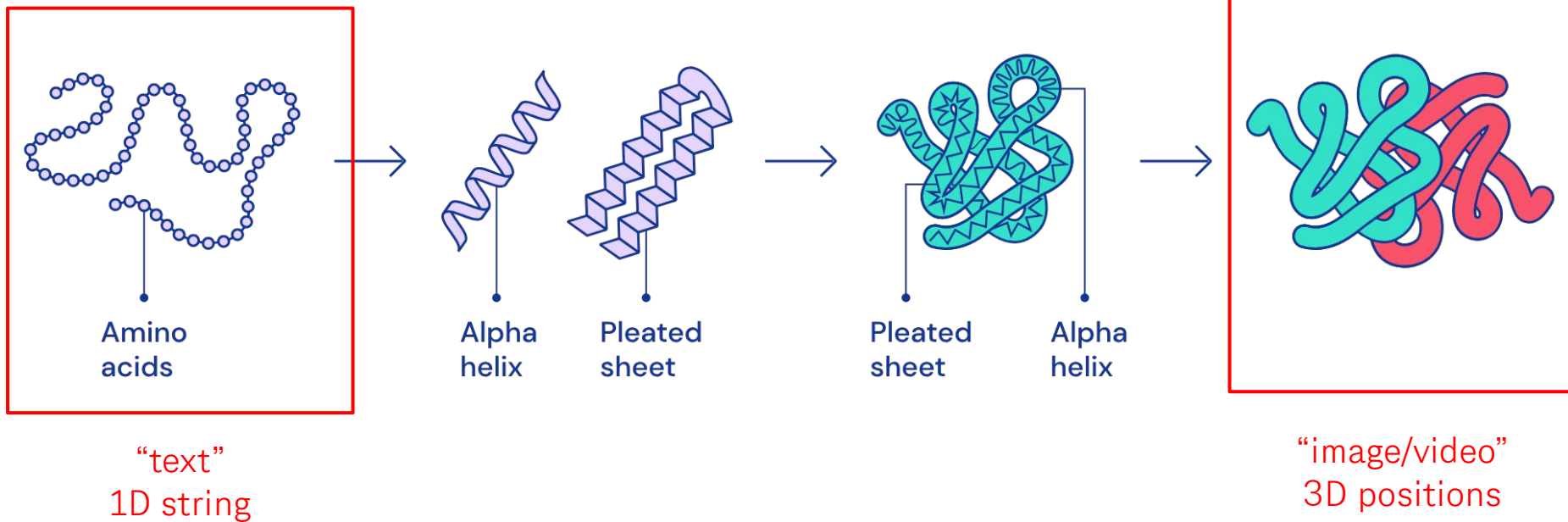
**Latent diffusion for all-atom generation**
*(in submission)*

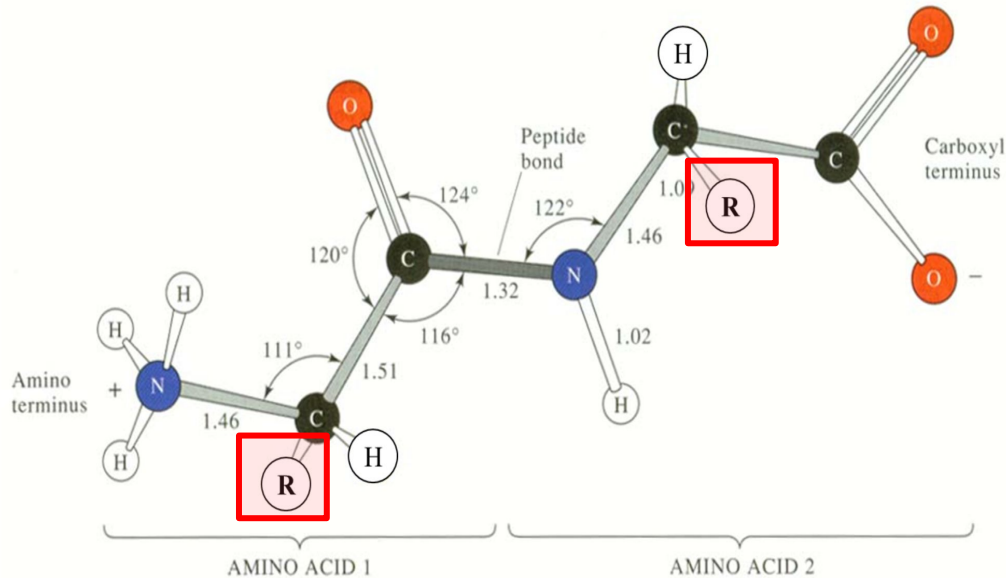Biological data selection from an information theoretic perspective
*(in progress)*

# Agenda



- Primer on AI for drug discovery
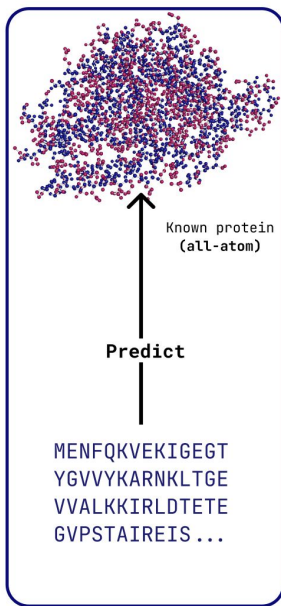- **PLAID** (<u>P</u>rotein <u>L</u>atent <u>I</u>nduced <u>D</u>iffusion)
- Issues and Diagnoses
- **CHEAP** (<u>C</u>ompressed <u>H</u>ourglass <u>E</u>mbedding <u>A</u>daptations of <u>P</u>roteins)
- Data selection and model understanding
- Future outlooks

# What exactly is a protein?



Amino acids → Alpha helix, Pleated sheet → Pleated sheet, Alpha helix →

"text"
1D string

"image/video"
3D positions

# Backbone structure vs. all-atom structure





GBYR...

(order of t-shirts => protein sequence)

# The co-generation problem



Known protein (all-atom)

**Predict**

MENFQKVEKIGEGT
YGVVYKARNKLTGE
VVALKKIRLDTETE
GVPSTAIREIS...

**AlphaFold / ESMFold**
Prediction only

Designed protein (backbone only)

**Sample**

Feasible protein structures

**Previous generative methods**
Backbone-only

sidechains are crucial for mediating function!

A

B

Met161   Tyr87

Trp185

Ser160

His237

Trp159

Asp206

GVPSTAIREIS...

**AlphaFold / ESMFold**
Prediction only

Designed protein (backbone only)

protein structures

**Previous generative methods**
Backbone-only

# Sidechain atoms generation require knowing the sequence

# All-atom design as a multimodal generation problem

# All-atom design as a multimodal generation problem



ALL-ATOM STRUCTURE

BACKBONE ATOMS

SEQUENCE & SIDECHAIN ATOMS

GSHMSREEIRKVVEEM
VRKLKQGSPEDISKYL
SPDVRGQEALKYMVRP

e.g. ESMFold

p(structure | sequence) p(sequence)

sample from

p(sequence | structure) p(structure)

e.g. ProteinMPNN

sample from

Goal:  p(sequence, structure)

sample from

# Motivation: Can we repurpose priors from pretrained models?



RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control

Vision-language models trained on internet-scale datasets capture useful priors for decision making tasks.
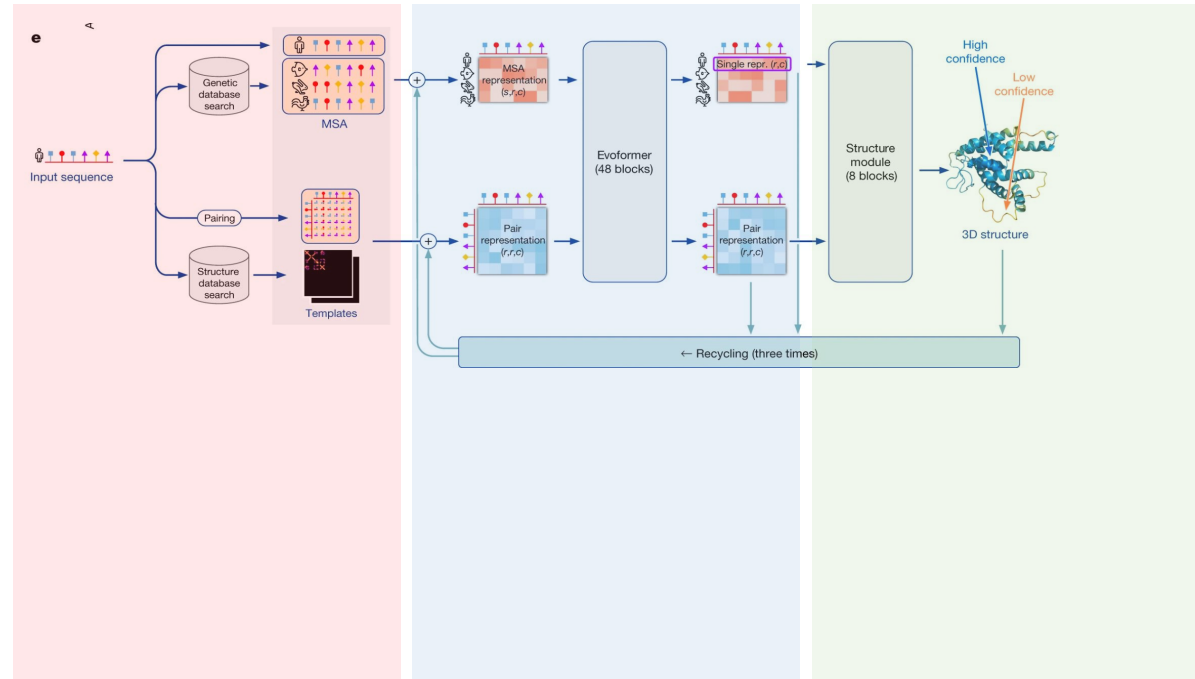
Can we apply this to biology?

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

# The base components: protein folding model architectures

**AlphaFold2:**

Uses an explicit retrieval step

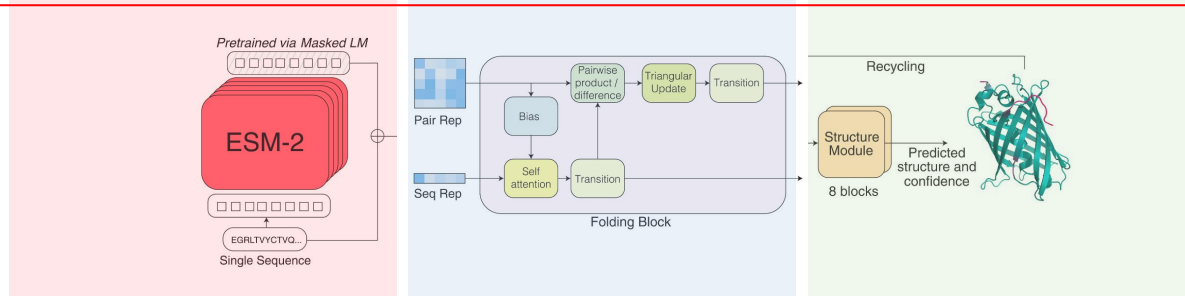# The base components: protein folding model architectures

**AlphaFold2:**
Uses an explicit retrieval step

**ESMFold:**
Replaces retrieval step with a **language model**
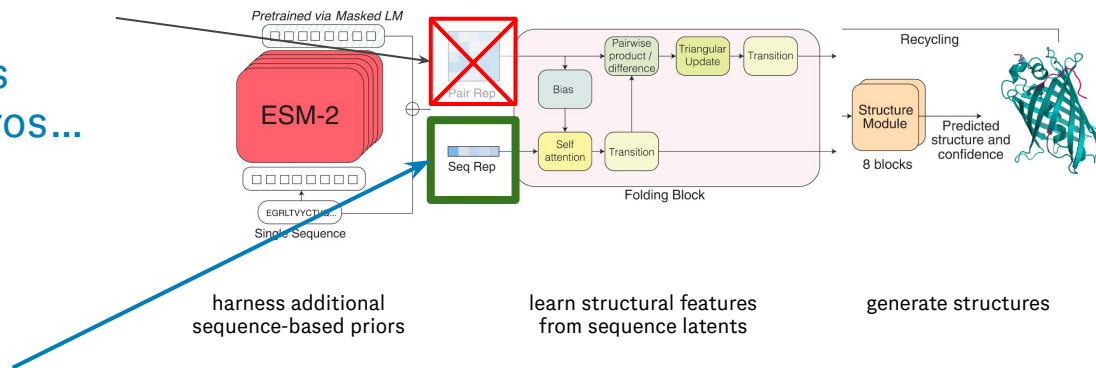


harness additional sequence-based priors

learn structural features from sequence latents

generate structures

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

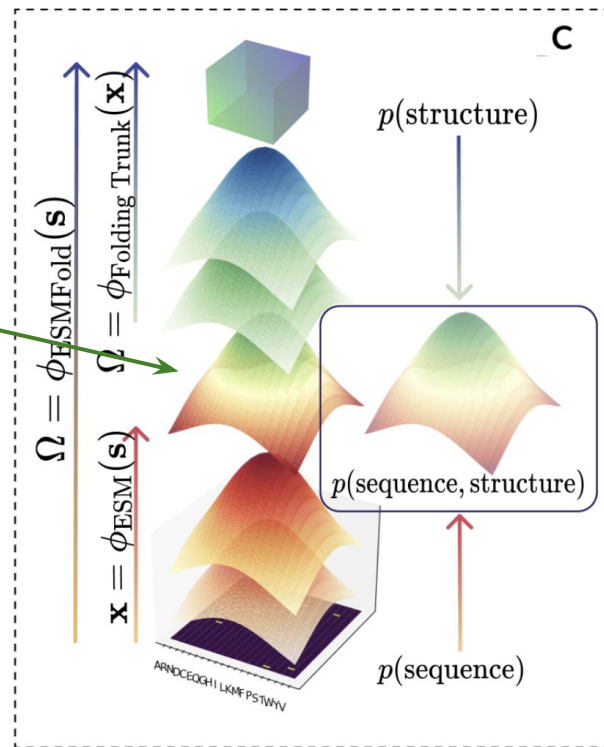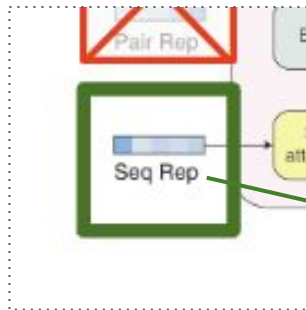→ Sequence representation contains all information about the structure!

Pretrained via Masked LM

ESM-2

Single Sequence

Pair Rep

Seq Rep

Bias

Pairwise product / difference

Triangular Update

Transition

Self attention

Transition

Folding Block

Recycling

Structure Module

8 blocks

Predicted structure and confidence

harness additional sequence-based priors

learn structural features from sequence latents

generate structures

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 (!)**
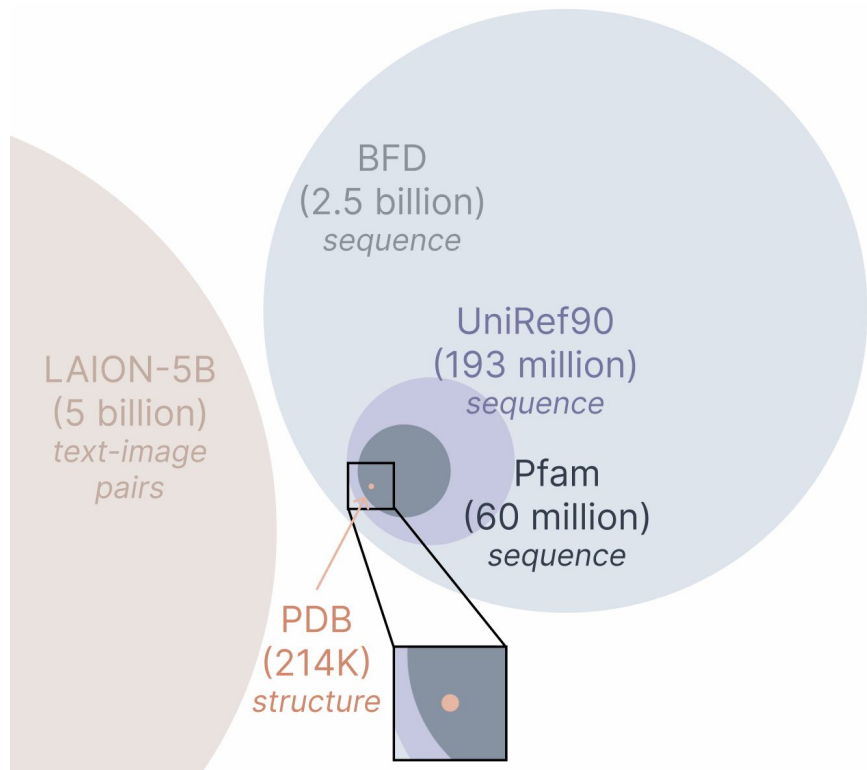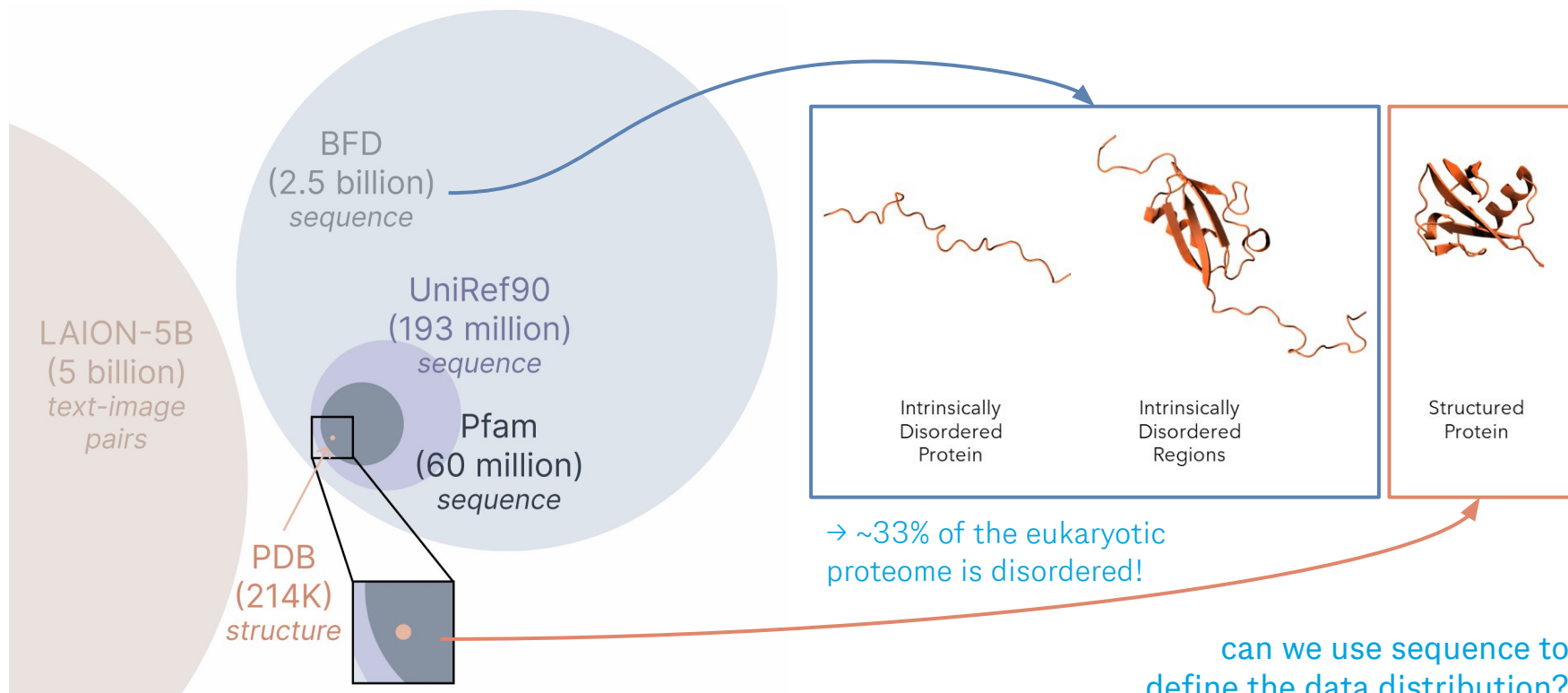
# Sequence data is cheaper to collect than structure
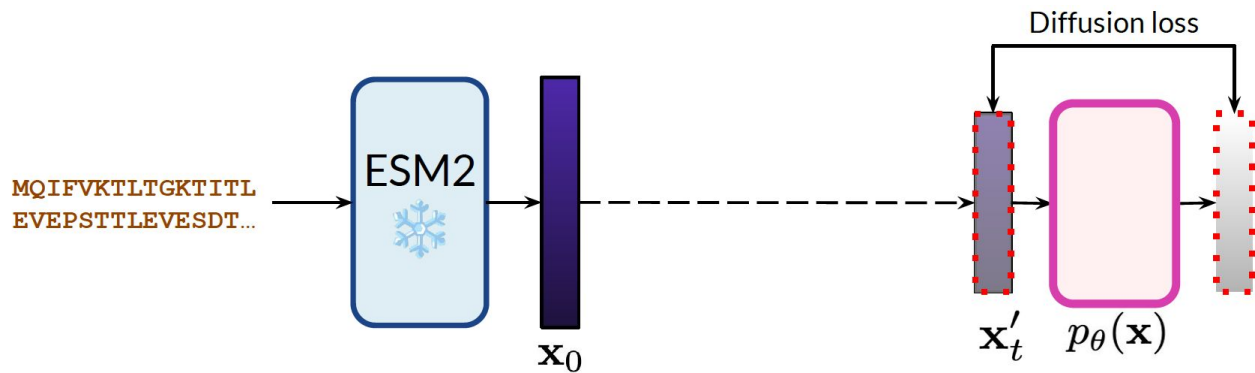


Source: https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data

# Sequence data is more abundant than structure



LAION-5B
(5 billion)
*text-image pairs*

BFD
(2.5 billion)
*sequence*

UniRef90
(193 million)
*sequence*

Pfam
(60 million)
*sequence*

PDB
(214K)
*structure*

# Sequence data has different coverage than structure



BFD
(2.5 billion)
*sequence*

UniRef90
(193 million)
*sequence*

Pfam
(60 million)
*sequence*

LAION-5B
(5 billion)
*text-image pairs*

PDB
(214K)
*structure*

Intrinsically Disordered Protein

Intrinsically Disordered Regions

Structured Protein

→ ~33% of the eukaryotic proteome is disordered!

can we use sequence to define the data distribution?

# PLAID v0.5: Training a latent diffusion model

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

# PLAID v0.5: Early attempts



What's preventing the model from learning?

**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*
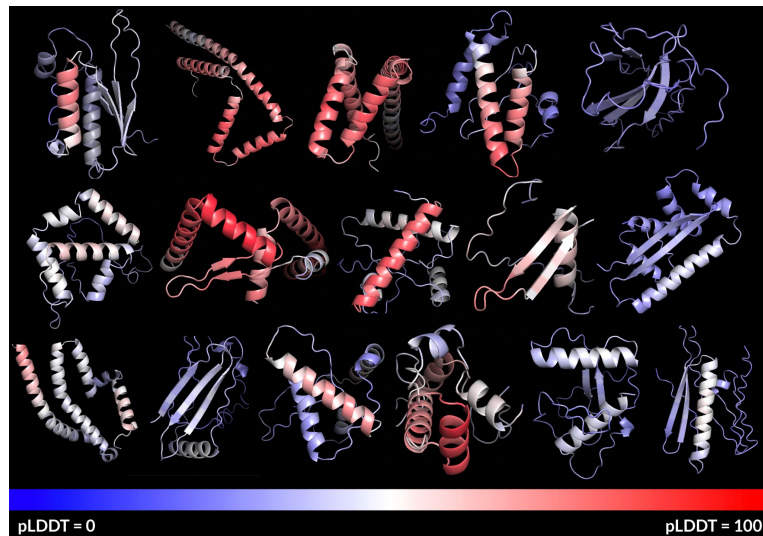
# Agenda

# 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:
    HxWx4 = 64 x 64 x 4
  - Ours:
    Lx1024 = 512 x 1024



G. NCSN++ (Song et al., 2021) FFHQ-$1024^2$ Reference Samples

# ESMFold latent space exhibits pathologically large values



Latent space will require regularization for diffusion to work.

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



Top 3 largest activation values per layer

Median activation value

ESM2 layers

# What if we just remove these wacky channels?

# What if we just remove these wacky channels?

# Addressing the hypotheses: embedding compression

### Issues and hypotheses

- Latent space requires regularization
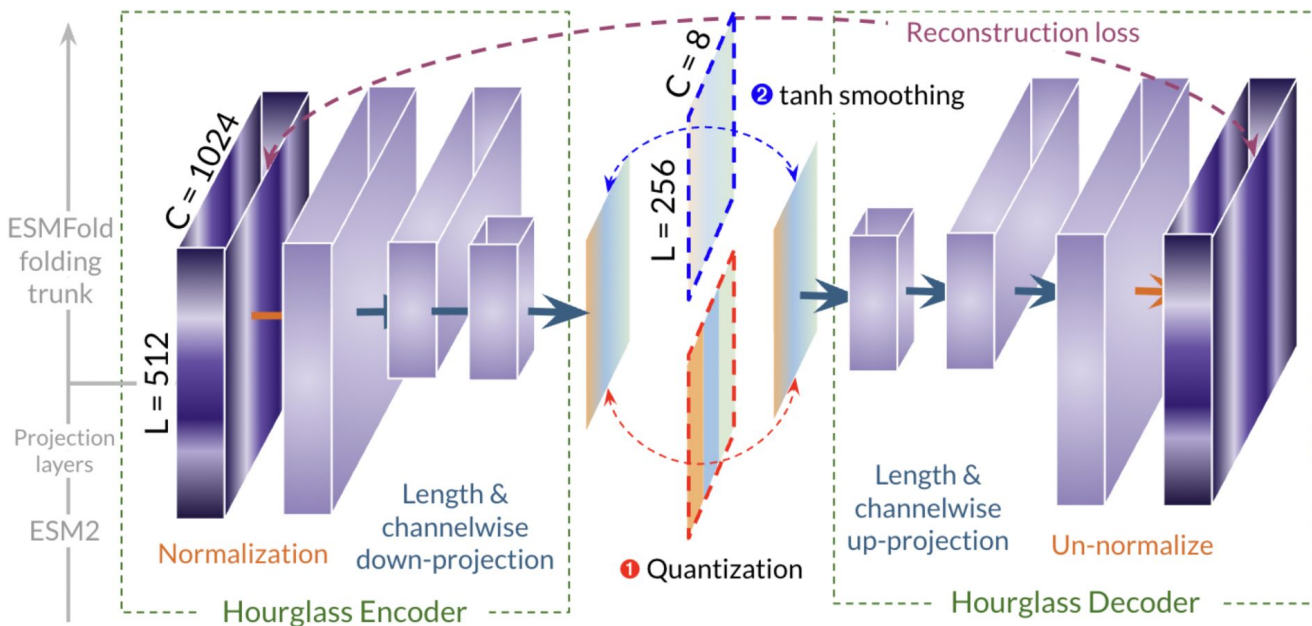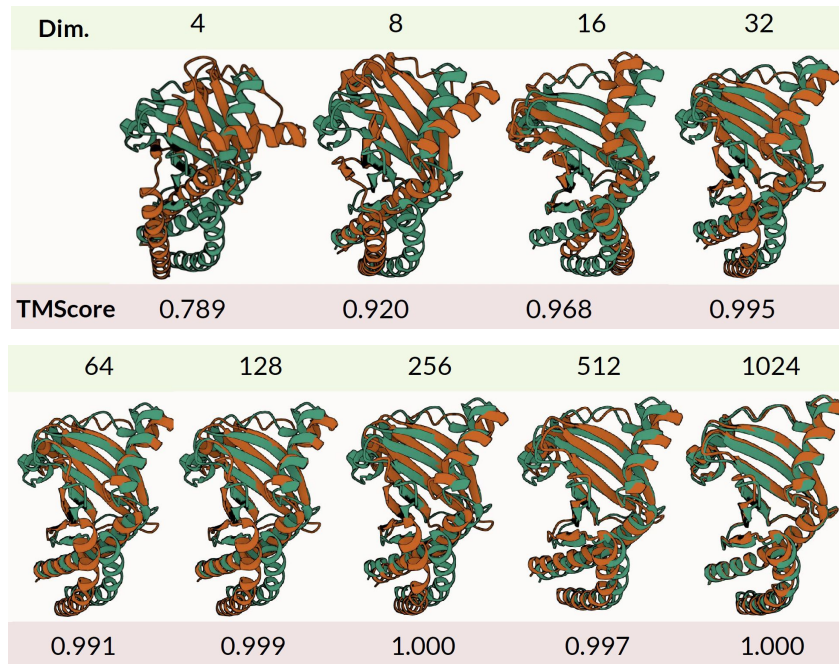- Overcome $O(L^2)$ memory constraints and increase protein length to 512
- Large latent space corresponds to high-resolution image generation
  - Rombach et al. latent space: HxWx4 = 64 x 64 x 4
  - Ours: Lx1024 = 512 x 1024

G. NCSN++ (Song et al., 2021) FFHQ-$1024^2$ Reference Samples

Diffusion models in their naive formulation often fail for 1024 x 1024 resolution generation.

Since not all channels are necessary, can we compress the embedding?

### Issues and hypotheses

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

pLDDT = 0    pLDDT = 100

Can we also reduce the protein length?

# An autoencoder for protein embedding compression

# Turns out the latent space is highly compressible!



| Dim. | 4 | 8 | 16 | 32 |
|---|---|---|---|---|
| TMScore | 0.789 | 0.920 | 0.968 | 0.995 |

| 64 | 128 | 256 | 512 | 1024 |
|---|---|---|---|---|
| 0.991 | 0.999 | 1.000 | 0.997 | 1.000 |

# Turns out the latent space is highly compressible!

# What about function information?



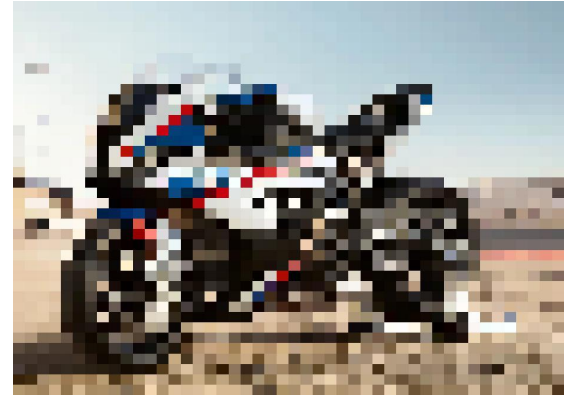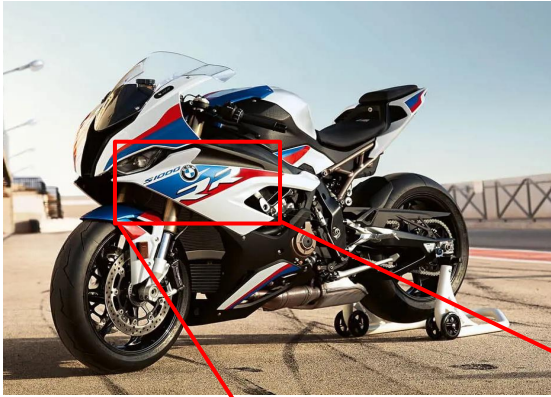Performance degradation with compression is more gradual…

…for some functions.

# What about function information?



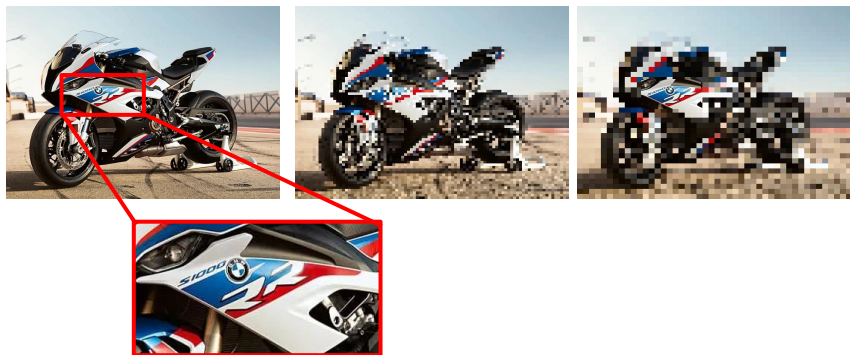Performance degradation with compression is more gradual…

…for some functions.

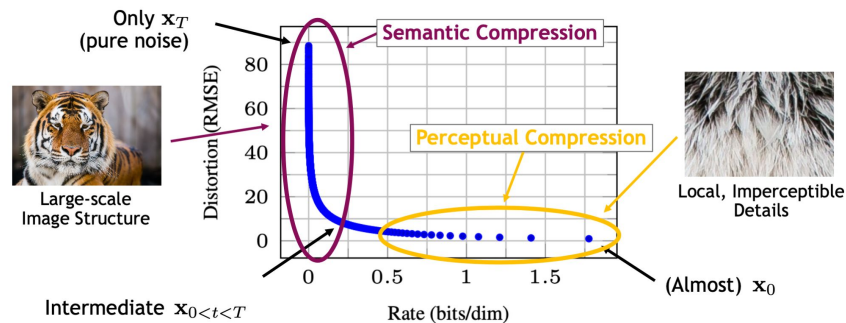# Intuition: what is the speed of this motorcycle?



→ BMW S1000RR: 188 mph

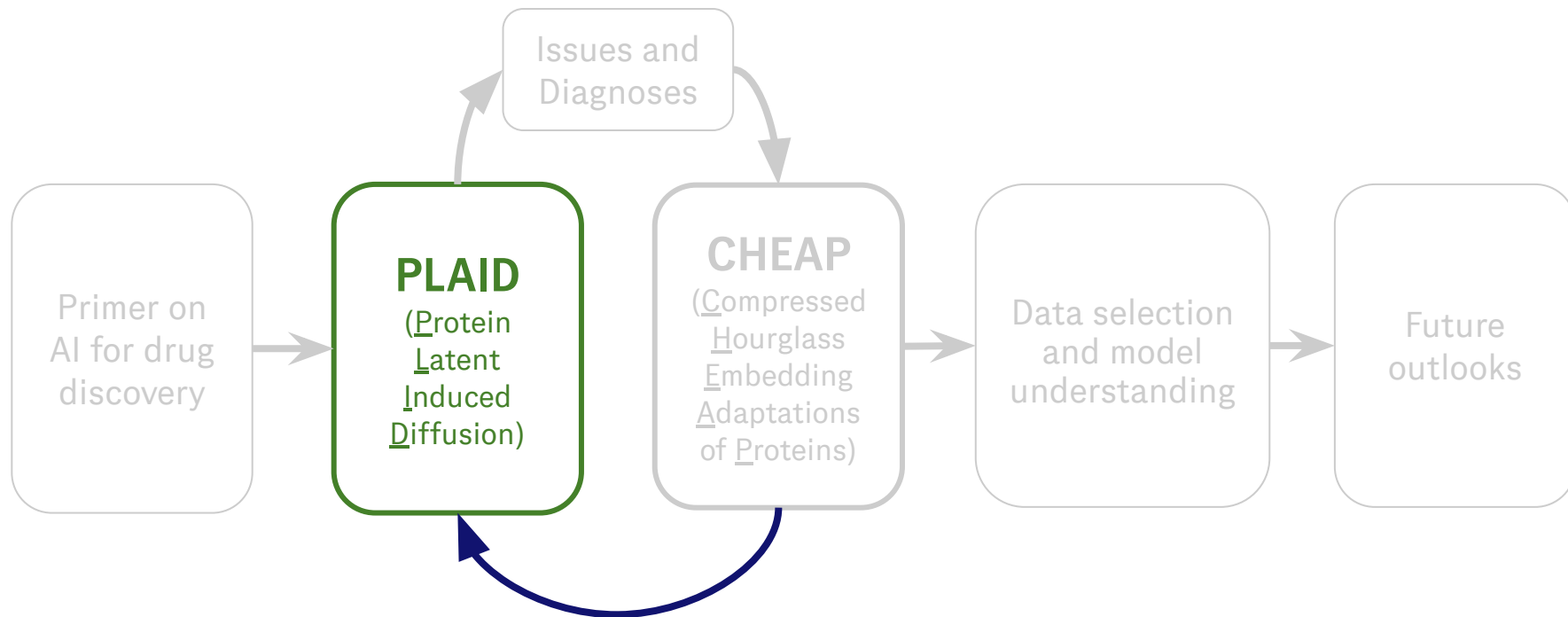# Intuition: what is the speed of this motorcycle?
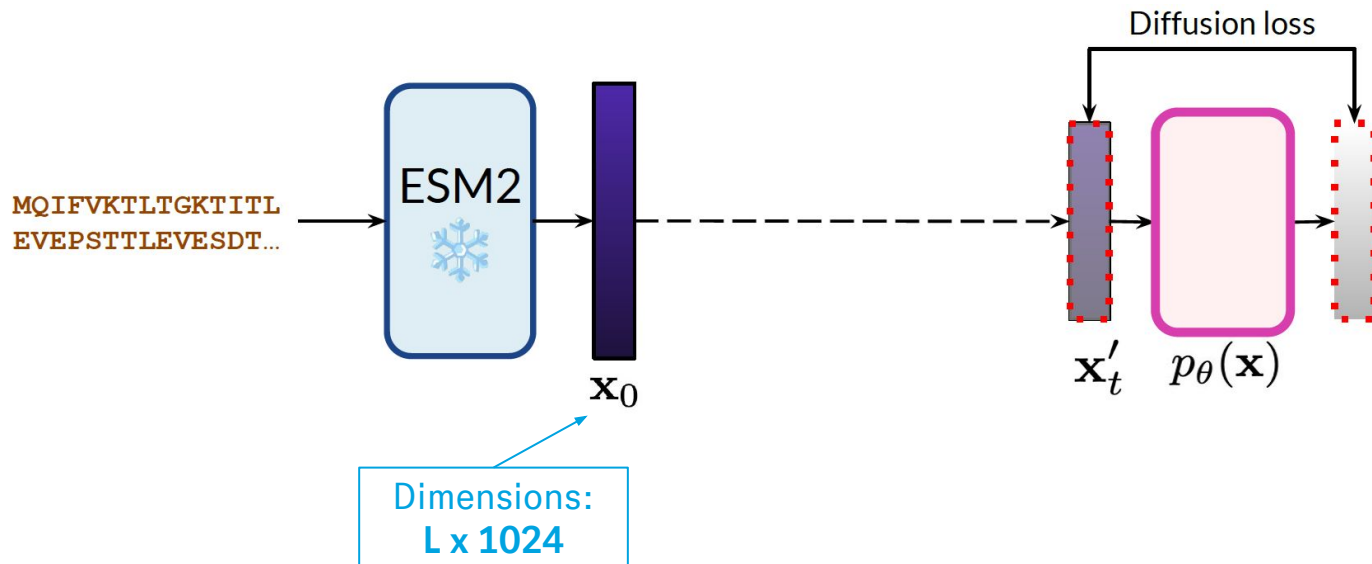


...what level of compression is optimal?



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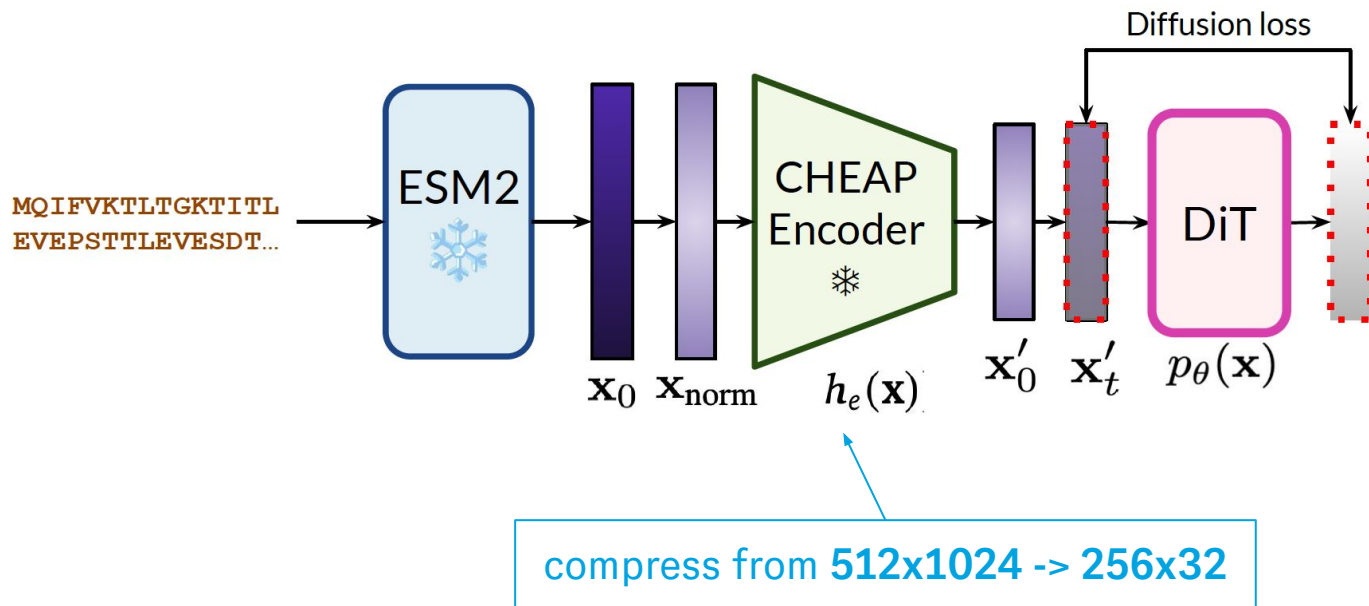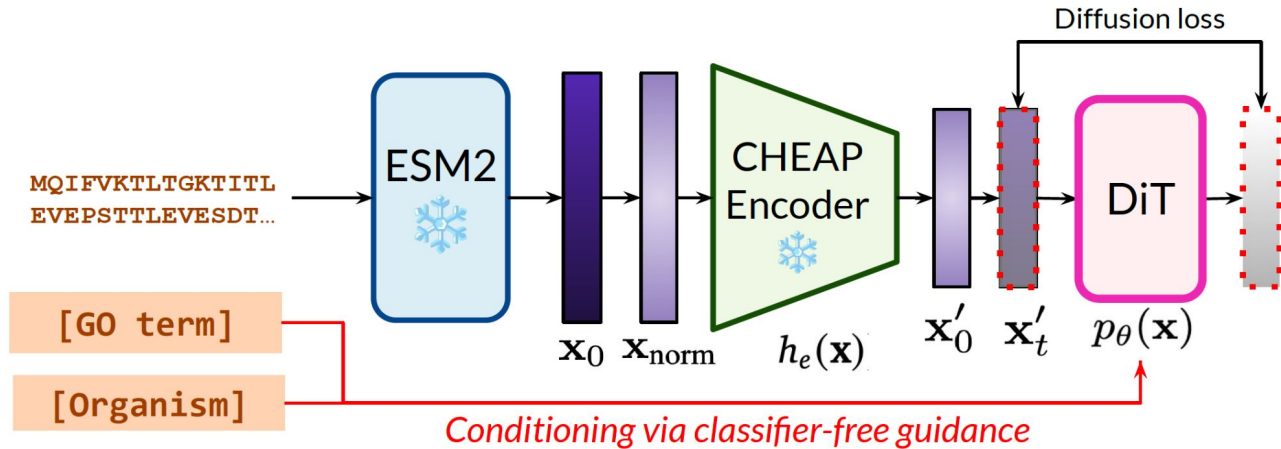
# Agenda

# Training the PLAID latent diffusion model…

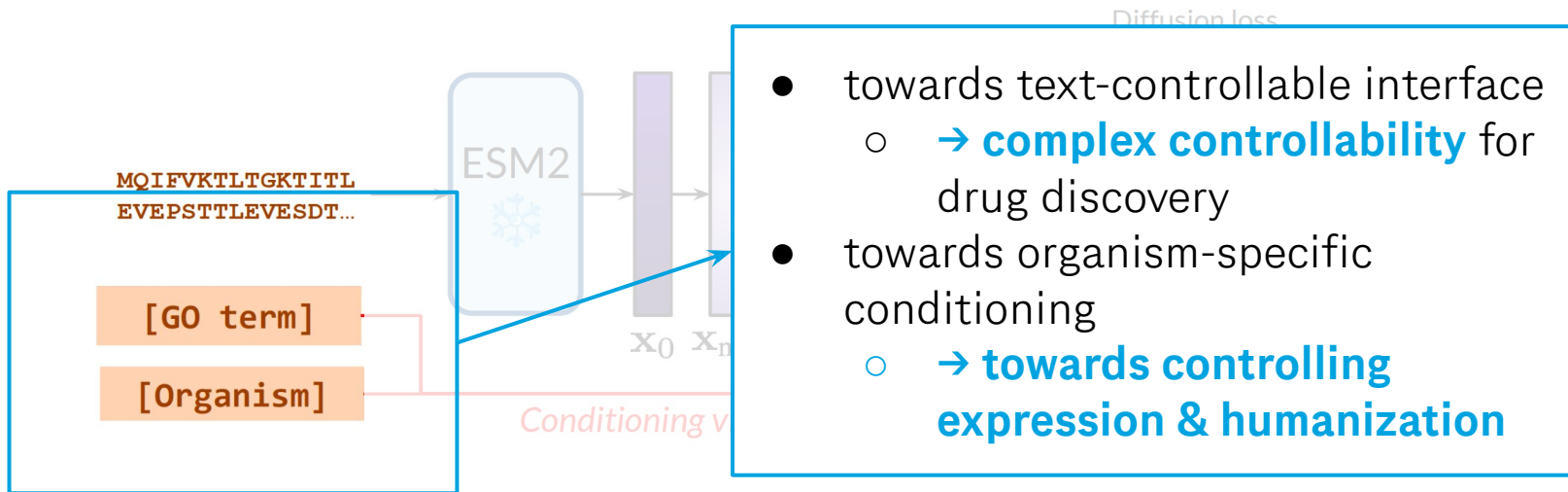# …but add embedding compression with CHEAP

# Adding compositional function + taxonomic conditioning



Sequence databases have more sample-annotation pairs!

# Adding compositional function + taxonomic conditioning



MQIFVKTLTGKTITL
EVEPSTTLEVESDT…

ESM2

[GO term]

[Organism]

$\mathbf{x}_0$  $\mathbf{x}_n$

Diffusion loss

*Conditioning v*

- towards text-controllable interface
  - **→ complex controllability** for drug discovery
- towards organism-specific conditioning
  - **→ towards controlling expression & humanization**
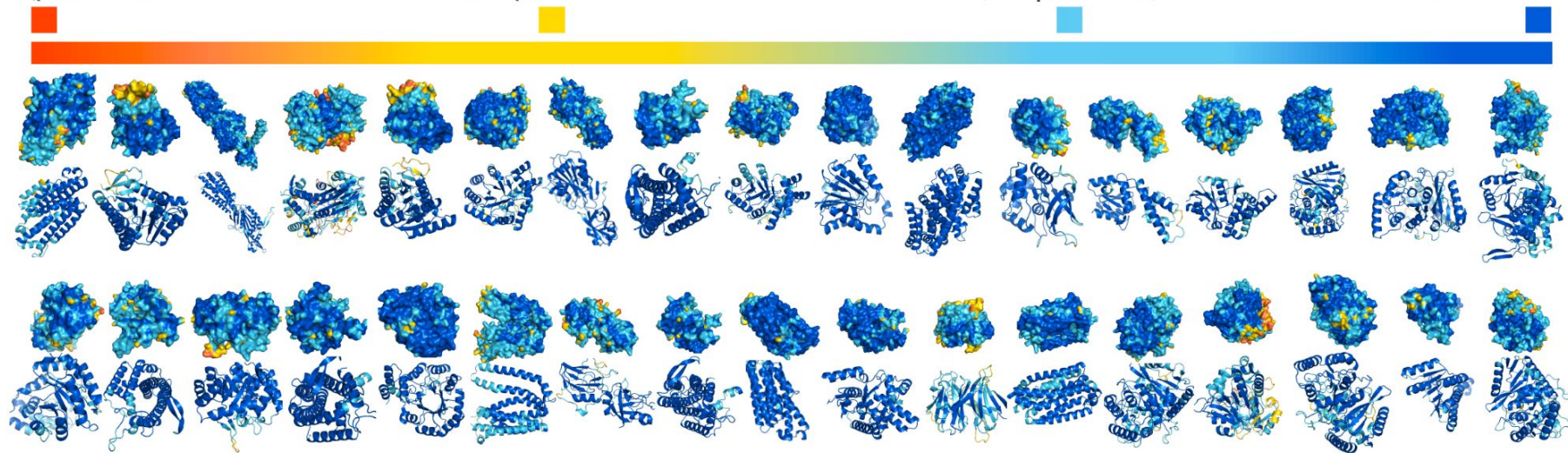
Sequence databases have more sample-annotation pairs!

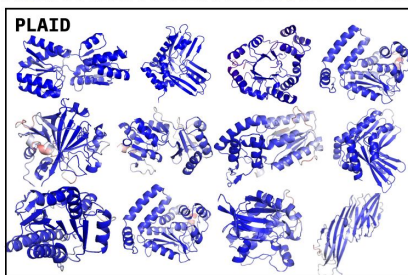# PLAID unconditionally generates diverse all-atom structures
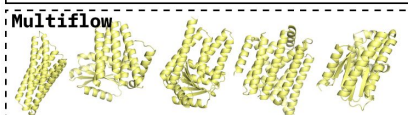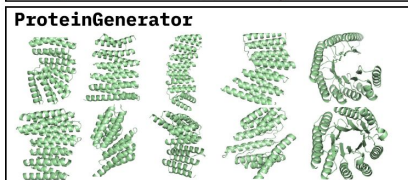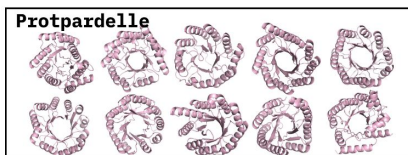
# PLAID unconditionally generates diverse, high-quality folds



**Protpardelle**
>len600_samp97
AGGGGGGGGGGGGGGGGGGGGGGGGGLGLGLLLPPAGL...
>len600_samp98
PPPPGGAGGGGAAAALAGGSPGGPPGGGGGGGGGGGG...
>len600_samp99
PPGPALPPSPGPGGVPPPPPLPPPPLPGGAPPAGGGLL...

**ProteinGenerator**
>len600_000097
GAAGLTAAAAVVGAAAAAGAAAAAALAAAAGAGAAAAA...
>len600_000098
AGAAGAAAAAAAAAAGAAAAGAGGGAGGAAAAAAAAAG...
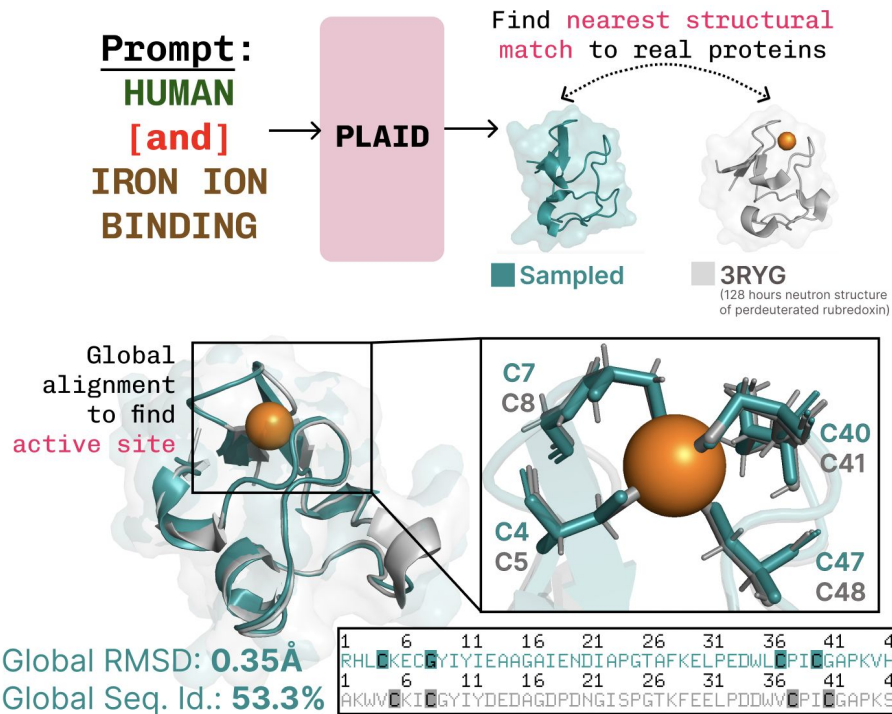>len600_000099
VAAAQAVQGAIAAAAALAATAALGLTAAGIAAPLLALV...

**Multiflow**
>len600_sample_97
LLGGLLGGLLGGAAGGAGAGAAAAGGGAVGVGVAGAVT...
>len600_sample_98
ADAATLTVGGGGTGGGGGAGGALGGAAAGGGGRVTLVV...
>len600_sample_99
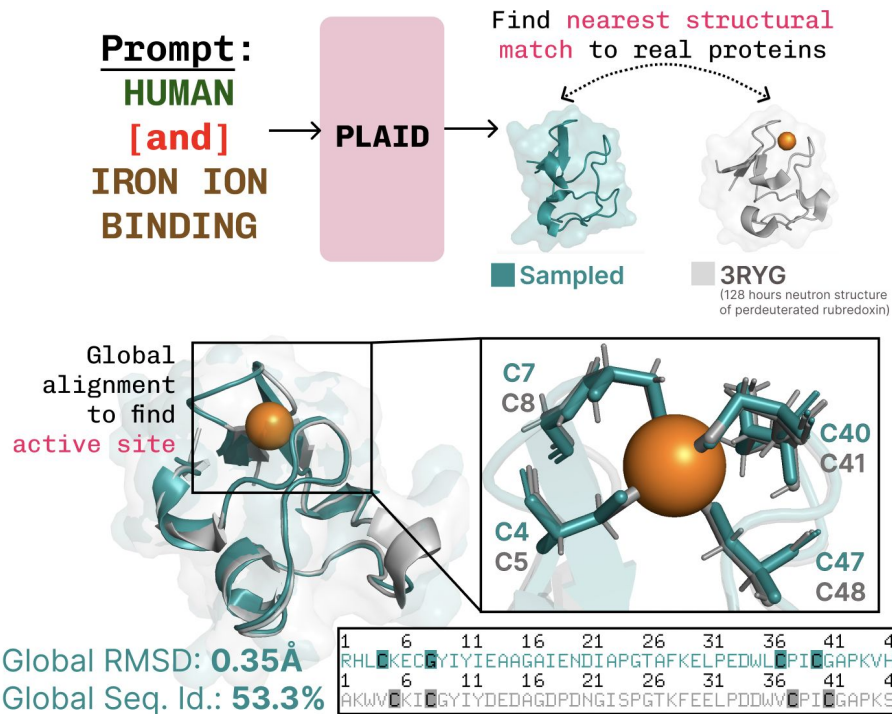AGGGAGLAGGAGGAGGAAAAAAAAAAAAAAGAGGGAAAA...

**PLAID**
>len600_sample97
PDMGTVLGLAHSVGHLDFKTPDLSVADLETNLALLAAH...
>len600_sample98
FEMFDDKGGDLWERAASSGQLLIDVAYLANNGLRDGAT...
>len600_sample99
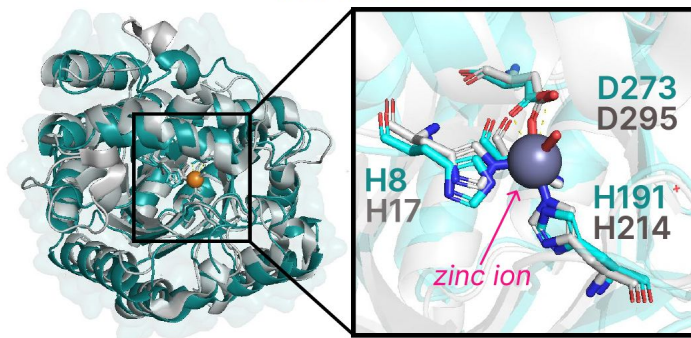GNGGQARGTDDPLTHALQTLFQSAALDQSLQGDPENAV...

# 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



Prompt:
HUMAN
[and]
IRON ION
BINDING

→ PLAID →

Find nearest structural match to real proteins

Sampled    3RYG
(128 hours neutron structure of perdeuterated rubredoxin)

Global alignment to find active site

C7
C8
C40
C41
C4
C5
C47
C48

Global RMSD: 0.35Å
Global Seq. Id.: 53.3%

```
1      6    11   16   21   26   31   36   41   46
RHLCKECGYIYIEAAGAIENDIAPGTAFKELPEDWLCPICGAPKVHFK
1      6    11   16   21   26   31   36   41   46   51
AKWVCKICGYIYDEDAGDPDNGISPGTKFEELPDDWVCPICGAPKSEFEKL FE
```

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

# Function-prompted generations learn active site sidechains



**Prompt:**
HUMAN [and] DEAMINASE ACTIVITY

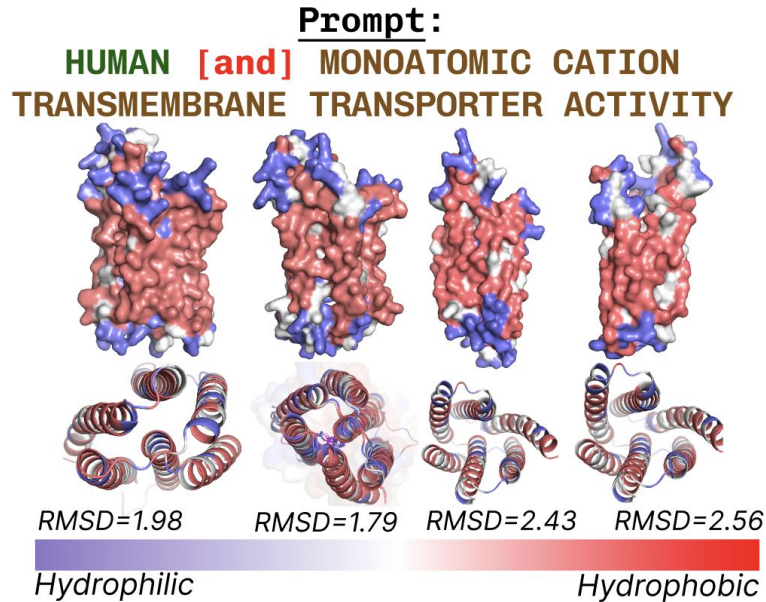D273
D295

H8
H17

H191
H214

zinc ion

RMSD: **2.25Å**
Seq. Id.: **24.3%**

■ Sampled
■ 7RTG (Crystal Structure of the Human Adenosine Deaminase 1)
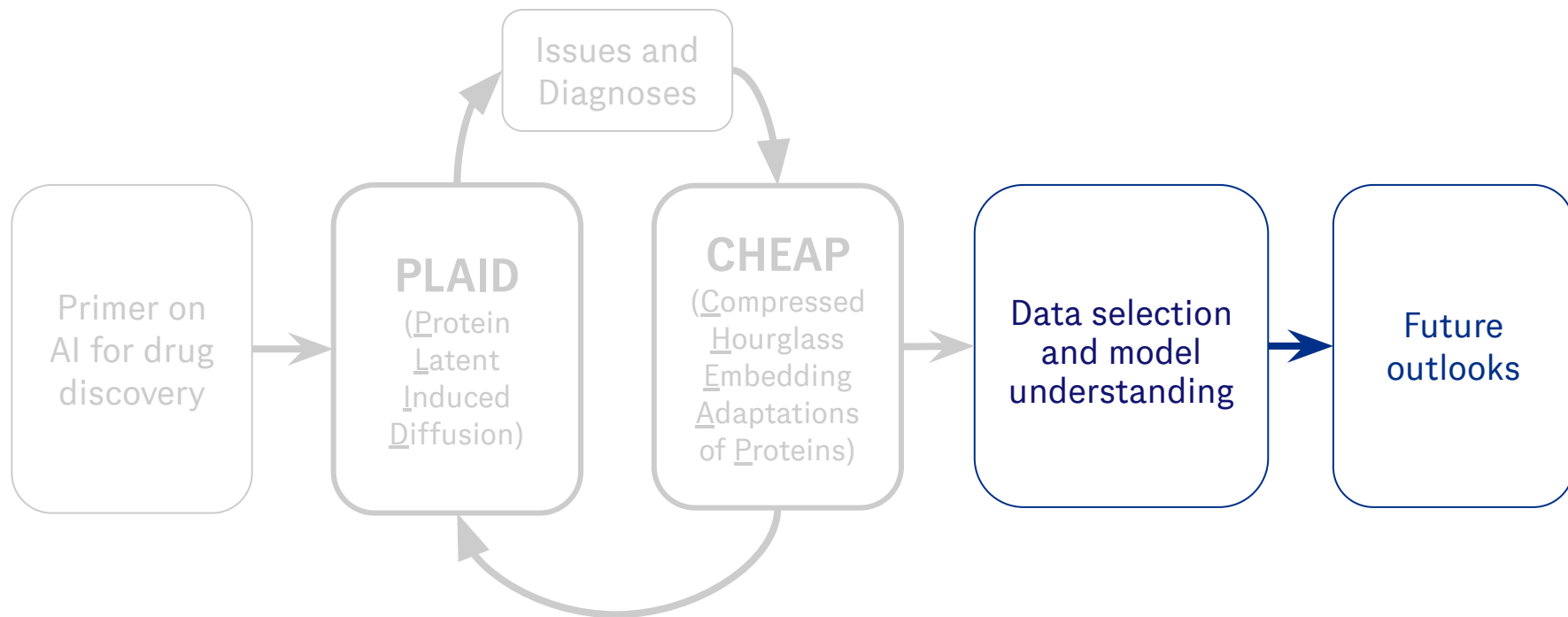
...despite these key residues not being adjacent in the sequence.

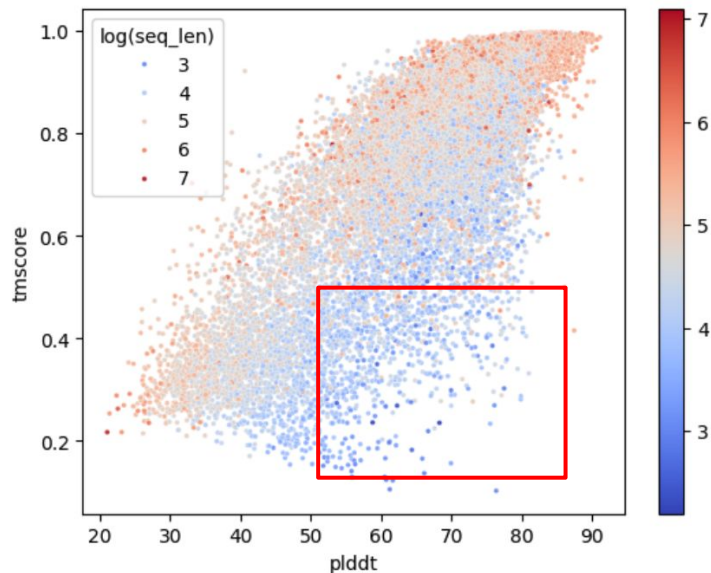# Transmembrane proteins exhibit expected hydrophobicity patterns



Hydrophobic residues are found at the core, as expected.

# Agenda

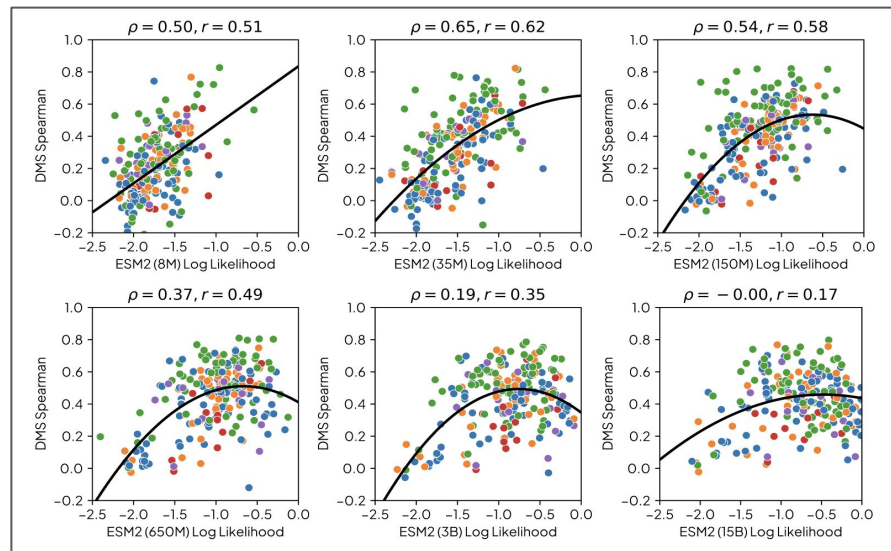# From proof-of-concept to deployment in AI for drug discovery

- Is the data learning a "biological world model", or artifacts of the training data?



Length determines overconfident predictions, but we often use pLDDT for generative model evaluation.

# From proof-of-concept to deployment in AI for drug discovery

- Is the data learning a "biological world model", or artifacts of the training data?

# From proof-of-concept to deployment in AI for drug discovery

## Genome modeling and design across all domains of life with Evo 2

Garyk Brixi*,1,2,3, Matthew G. Durrant*,1,2, Jerome Ku*,1,2, Michael Poli*,2,3,5,
Greg Brockman**,2,6,§, Daniel Chang**,1,2,3, Gabriel A. Gonzalez**,1,2, Samuel H. King**,1,2,3,
David B. Li**,1,2,3, Aditi T. Merchant**,1,2,3, Mohsen Naghipourfar**,1,2,7, Eric Nguyen**,2,3,
Chiara Ricci-Tam**,1,2, David W. Romero**,2,4, Gwanggyu Sun**,1,2, Ali Taghibakshi**,2,4,
Anton Vorontsov**,2,4, Brandon Yang**,2,6, Myra Deng8, Liv Gorton8, Nam Nguyen8,
Nicholas K. Wang8, Etowah Adams9, Stephen A. Baccus3, Steven Dillmann3,
Stefano Ermon3, Daniel Guo1,3, Rajesh Ilango1, Ken Janik4, Amy X. Lu7, Reshma Mehta6,
Mohammad R.K. Mofrad7, Madelena Y. Ng3, Jaspreet Pannu3, Christopher Ré3,
Jonathan C. Schmok1, John St. John4, Jeremy Sullivan1, Kevin Zhu7, Greg Zynda4,
Daniel Balsam8,10, Patrick Collison1,10, Anthony B. Costa4,10, Tina Hernandez-
Boussard3,10, Eric Ho8,10, Ming-Yu Liu4,10, Thomas McGrath8,10,
Kimberly Powell4,10, Dave P. Burke‡,1,2,10, Hani Goodarzi‡,1,2,10,11,
Patrick D. Hsu‡,†,1,2,7,10, Brian L. Hie‡,†,1,2,3,10

1Arc Institute; 2Core Contributor, Evo 2 Team; 3Stanford University; 4NVIDIA;
5Liquid AI; 6Independent Researcher; 7University of California, Berkeley;
8Goodfire; 9Columbia University; 10Senior Contributor, Evo 2 Team;
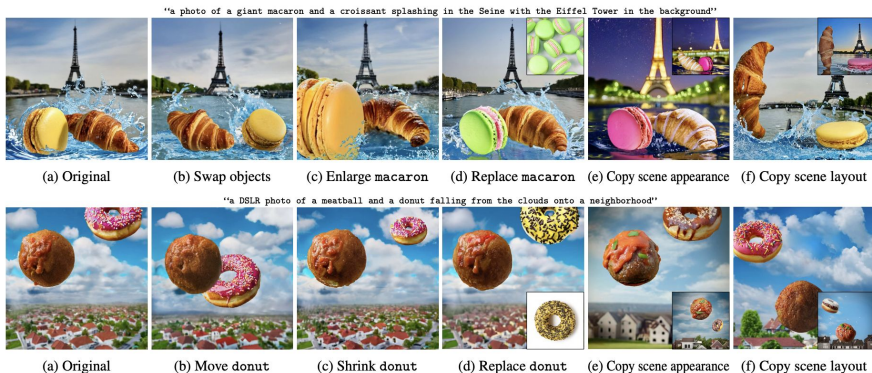11University of California, San Francisco

# Medium-term directions…

Latent diffusion for drug design

- Leveraging **"fuzziness"** in # of atoms and binding position
- Alleviating computational challenges for large complexes with **compression**
- Semantic control in latent space



"a photo of a giant macaron and a croissant splashing in the Seine with the Eiffel Tower in the background"

(a) Original    (b) Swap objects    (c) Enlarge macaron    (d) Replace macaron    (e) Copy scene appearance    (f) Copy scene layout

"a DSLR photo of a meatball and a donut falling from the clouds onto a neighborhood"

(a) Original    (b) Move donut    (c) Shrink donut    (d) Replace donut    (e) Copy scene appearance    (f) Copy scene layout

Diffusion self-guidance for controllable image generation. Epstein et al., 2023

# Medium-term directions…

Multimodal biophysical reasoning / chain-of-thought "scratchpad"



April 16, 2025   Release

# Thinking with images

OpenAI o3 and o4-mini represent a significant breakthrough in visual perception by reasoning with images in their chain of thought.



Example of how biochemical implausibilities can be reasoned through language.

Image source: PoseBusters documentation

# Long term goals…

- How can we move to a "target-agnostic" paradigm in drug discovery using advances in task-agnostic AI systems?
  - AI for biology as reasoning about the molecular-level world

- How can we extrapolate / "reason" beyond human intelligence?
  - How does data availability & simulation fidelity affect how this is done?

- How can we work *with* rather than *against* Moravec's Paradox, using scientific applications as a testbed?

# Acknowledgements

# Acknowledgements: Pieter!

# Acknowledgements: Prescient Design / Genentech



*and very many more Prescient team members!*

# Acknowledgements: MSR & Google Brain

# Acknowledgements: collaborators



*and very many more!*

# Acknowledgements: committee & other faculty



Arc Institute

Innovative Genomics Institute

# Acknowledgements: RLL labmates

# Acknowledgements: admin

Acknowledgements:
friends!! berkeley!!

Acknowledgements:
friends!! berkeley!!

Acknowledgements:
friends!! toronto/sf!!

# Acknowledgements: family <3

# Acknowledgements: family <3

# Acknowledgements: family <3

# Acknowledgements

the end!!!! 🫶

# Linear interpolation in the latent space

# Linear interpolation in the latent space



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

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

# ...noising the compressed latent space <u>does</u> map to corrupted structures

# …noising the compressed latent space <u>does</u> map to corrupted structures

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

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

→ a pervasive issue across LLMs, ViTs, etc.

**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
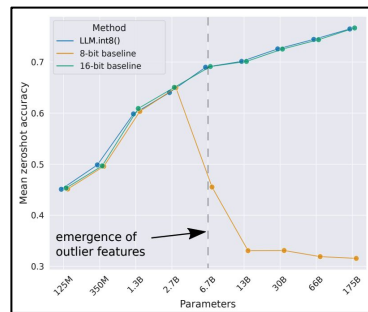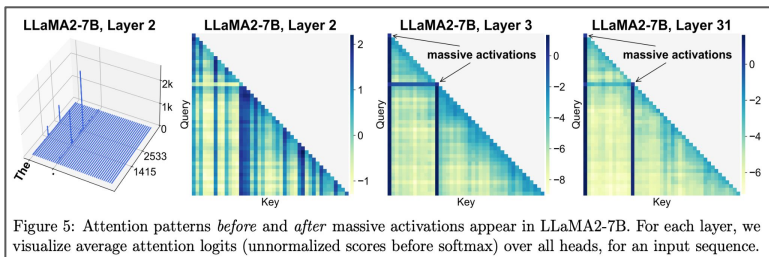


Figure 5: Attention patterns *before* and *after* massive activations appear in LLaMA2-7B. For each layer, we visualize average attention logits (unnormalized scores before softmax) over all heads, for an input sequence.

`LLM.int8()`: **8-bit Matrix Multiplication for Transformers at Scale**
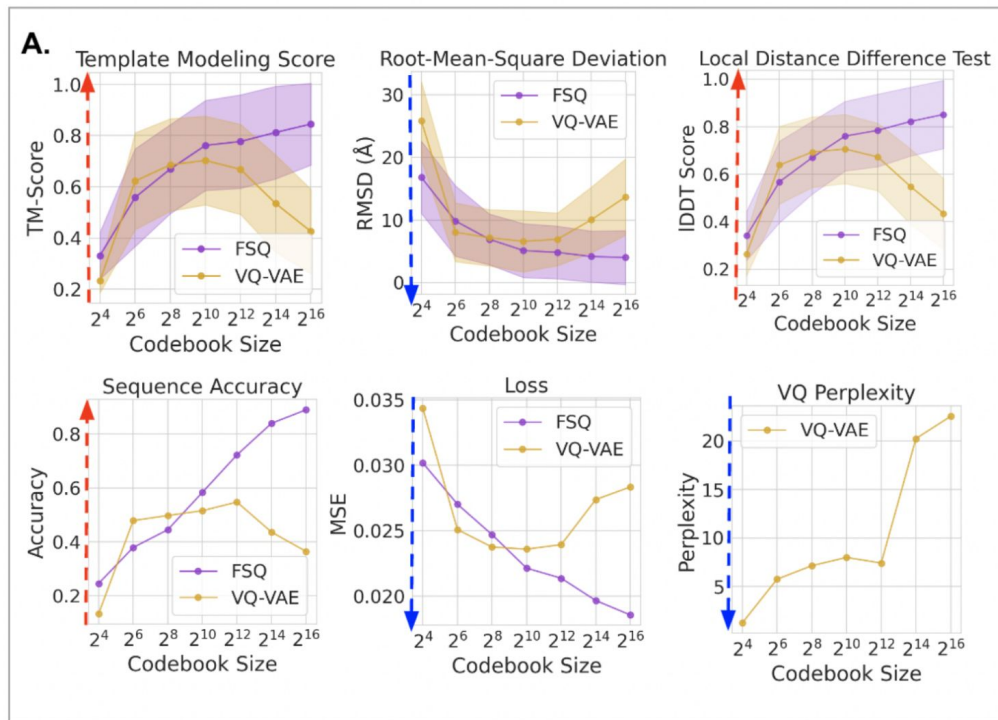
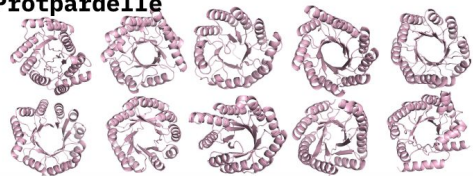Tim Dettmers[λ*]     Mike Lewis[†]     Younes Belkada[§∓]     Luke Zettlemoyer[†λ]

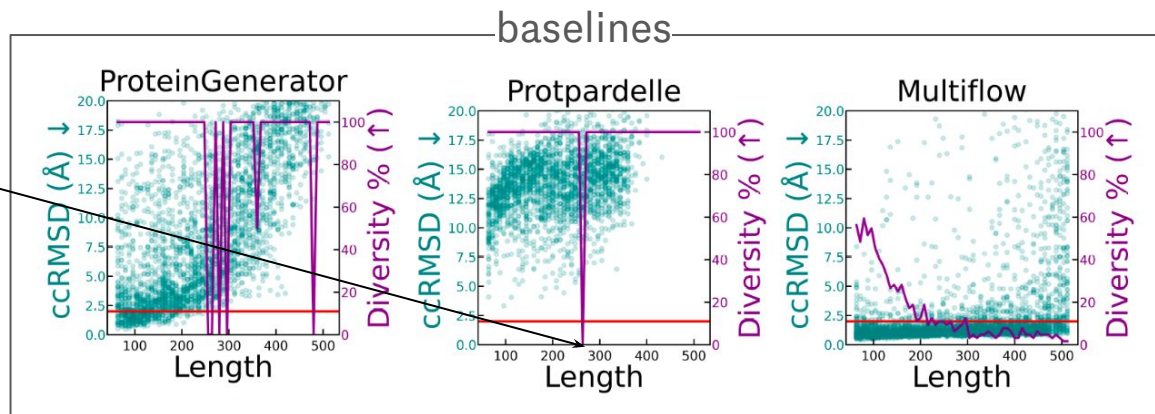# All-atom structural tokenizer, obtained from sequence alone

# PLAID unconditionally generates diverse, high-quality folds
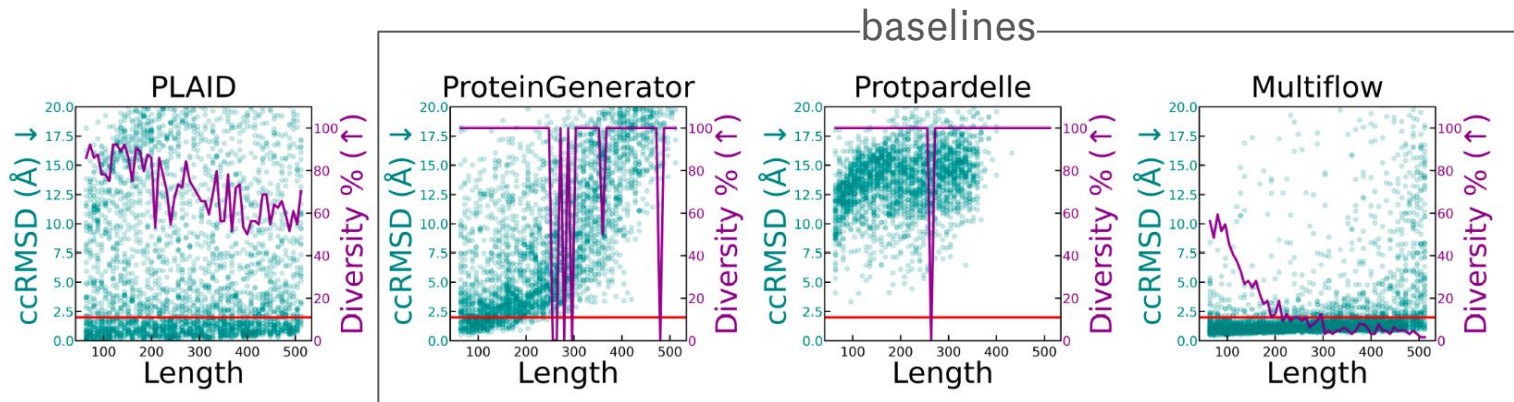


baselines

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

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

Protpardelle

Protpardelle
>len600_samp97
AGGGGGGGGGGGGGGGGGGGGGGGGGLGLGLLLPPAGL...
>len600_samp98
PPPPGGAGGGGAAAALAGGSPGGPPGGGGGGGGGGGGG...
>len600_samp99
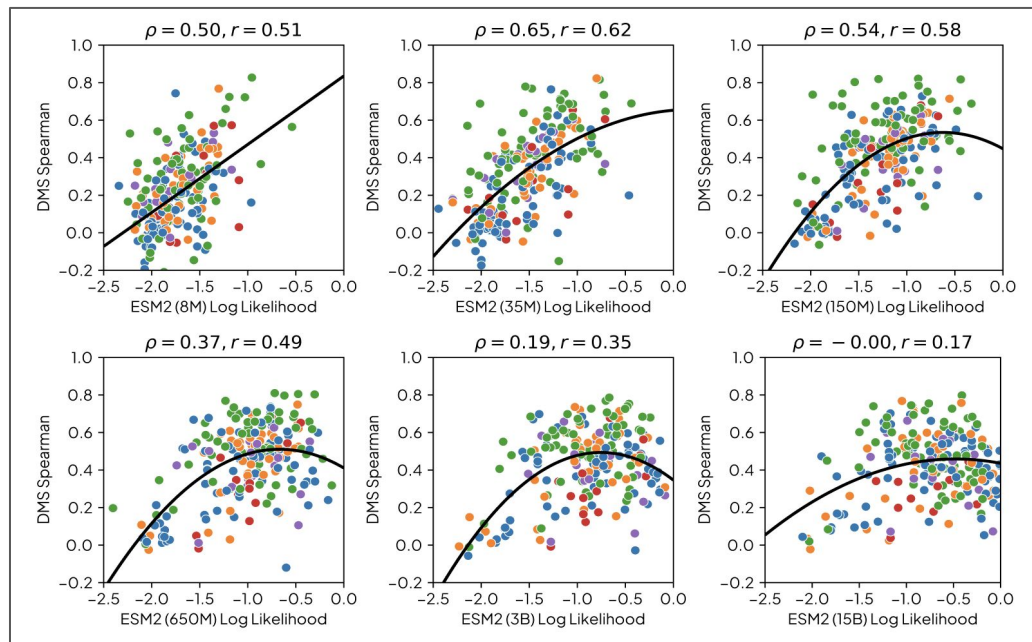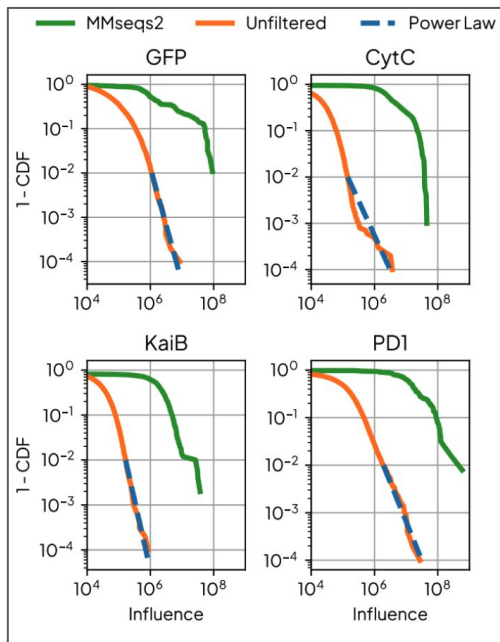PPGPALPPSPGPGGVPPPPPLPPPPLPGGAPPAGGGLL...

# PLAID unconditionally generates diverse, high-quality folds

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

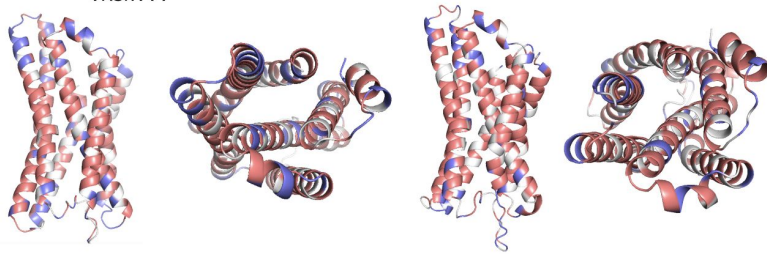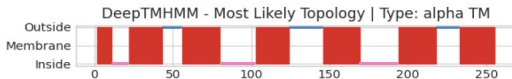# From proof-of-concept to deployment in AI for drug discovery

# Transmembrane proteins exhibit expected numbers of helices



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