



COS 597N: Machine Learning for Structural Biology

Lecture 4

Fall 2023

Course Logistics

- Optional student-only “precept”, Tuesdays at 4:30p in CS 401.
- Today:
 - Protein structure determination and cryo-EM reconstruction
- Next week 9/28: Protein language modeling — modified format!
 - Flash talks (groups of 1-2) + guest instructor (Adam Lerer)
 - (Very) short writing assignment
 - More details and paper sign up by the end of this week (end of day Friday): <https://docs.google.com/spreadsheets/d/1WznSeVYRaCFk8cLzGpxKRhe5JM65TZLd-Ge29byuze4/edit#gid=0>
- Oct 12: Protein design
- Oct 19 (fall break): No class + Project proposal due
 - Guidelines: <https://docs.google.com/document/d/1bKyklL9v-N-Yac1tBQCNi8CGQHsN5wZ5BQM7ZDo4WN4/edit>

This lecture

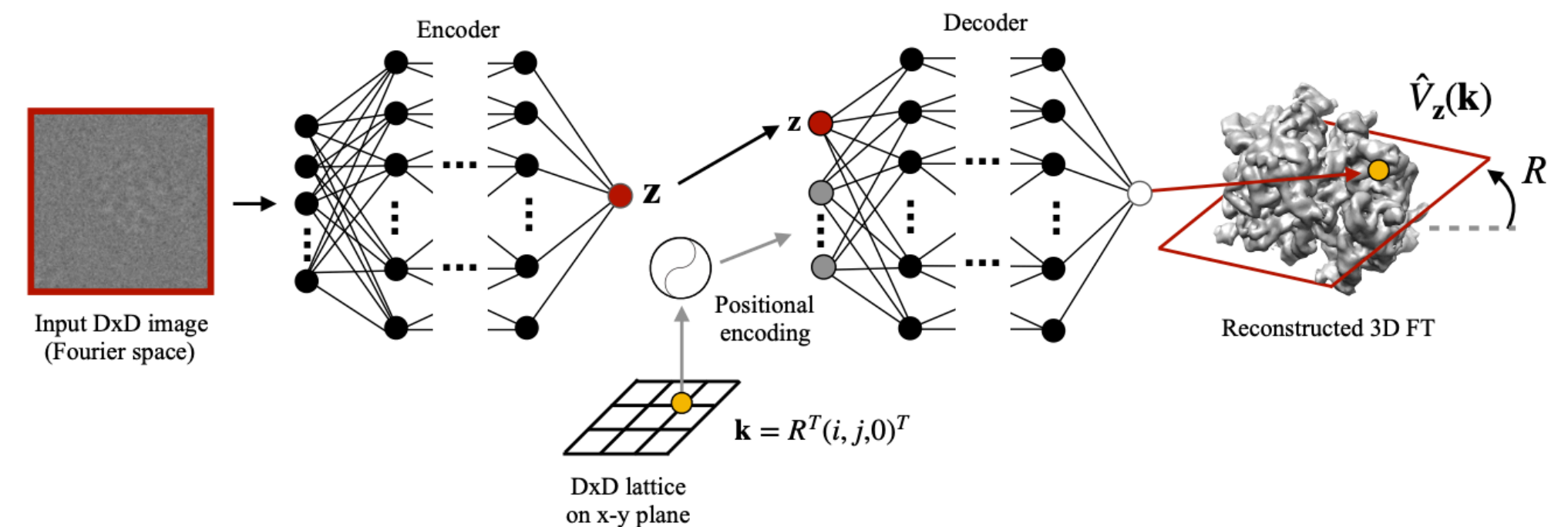
CryoDRGN

- (Recap): Who went to John Jumper's talk?
- CryoDRGN: Deep Reconstructing Generative Networks
 - Seminar
 - Figure by figure
- Questions:
 - What did you think of the papers?
 - What are the differences between conference vs. journal paper?
 - Who is familiar with NeRFs and implicit neural representations?
 - Any other thoughts/reflections?

[Submitted on 11 Sep 2019 (v1), last revised 15 Feb 2020 (this version, v3)]

Reconstructing continuous distributions of 3D protein structure from cryo-EM images

Ellen D. Zhong, Tristan Bepler, Joseph H. Davis, Bonnie Berger



Article | [Published: 04 February 2021](#)

CryoDRGN: reconstruction of heterogeneous cryo-EM structures using neural networks

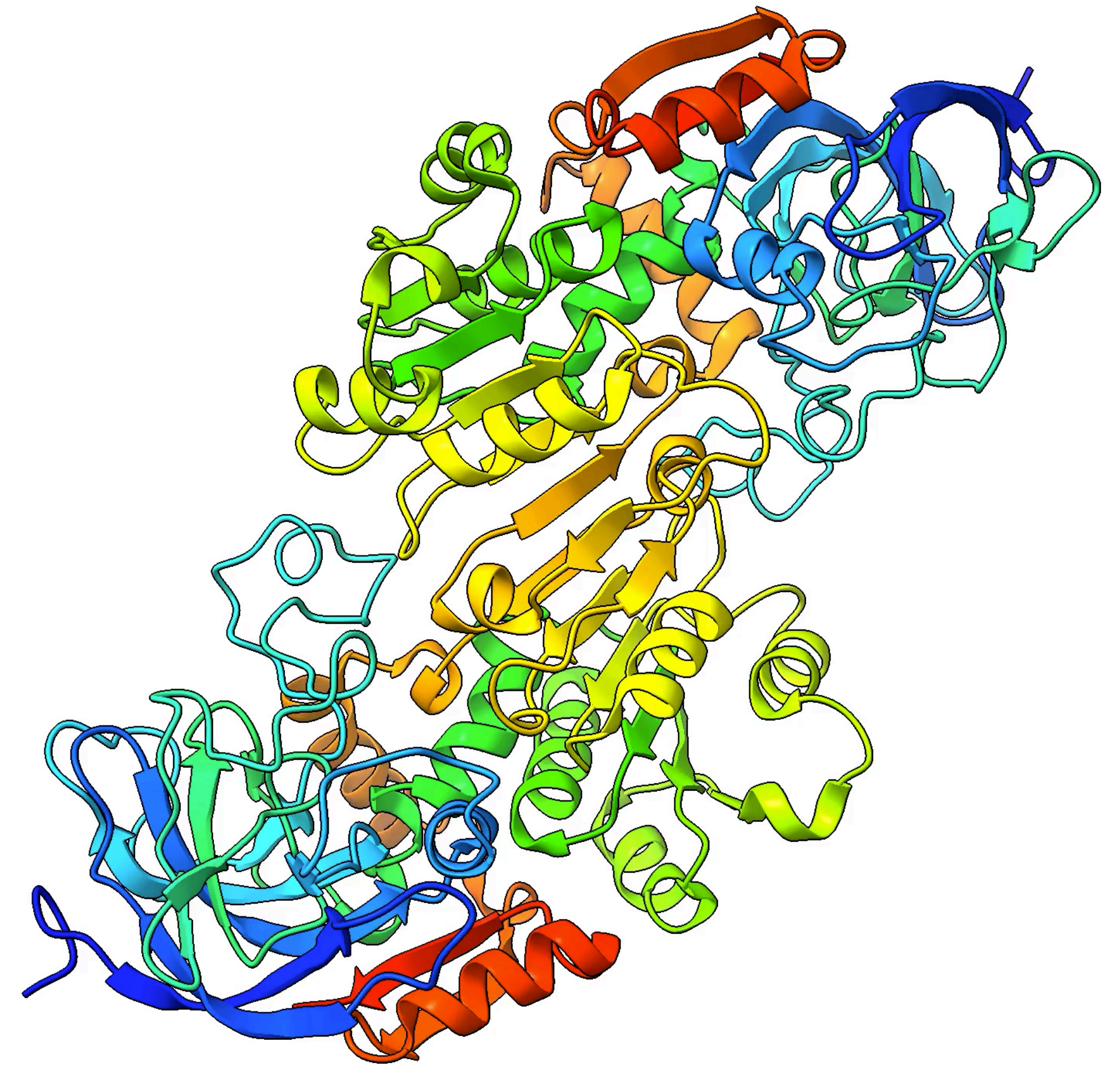
[Ellen D. Zhong](#), [Tristan Bepler](#), [Bonnie Berger](#)  & [Joseph H. Davis](#) 

[Nature Methods](#) **18**, 176–185 (2021) | [Cite this article](#)

31k Accesses | **171** Citations | **177** Altmetric | [Metrics](#)

Outline

- **Motivation:** Why do we care about protein structure?



Outline

- **Motivation:** Why do we care about protein structure?
- **Background:** Cryo-EM reconstruction & the heterogeneity problem

**Continuously heterogeneous
hyper-objects in cryo-EM and 3-D
movies of many temporal dimensions**

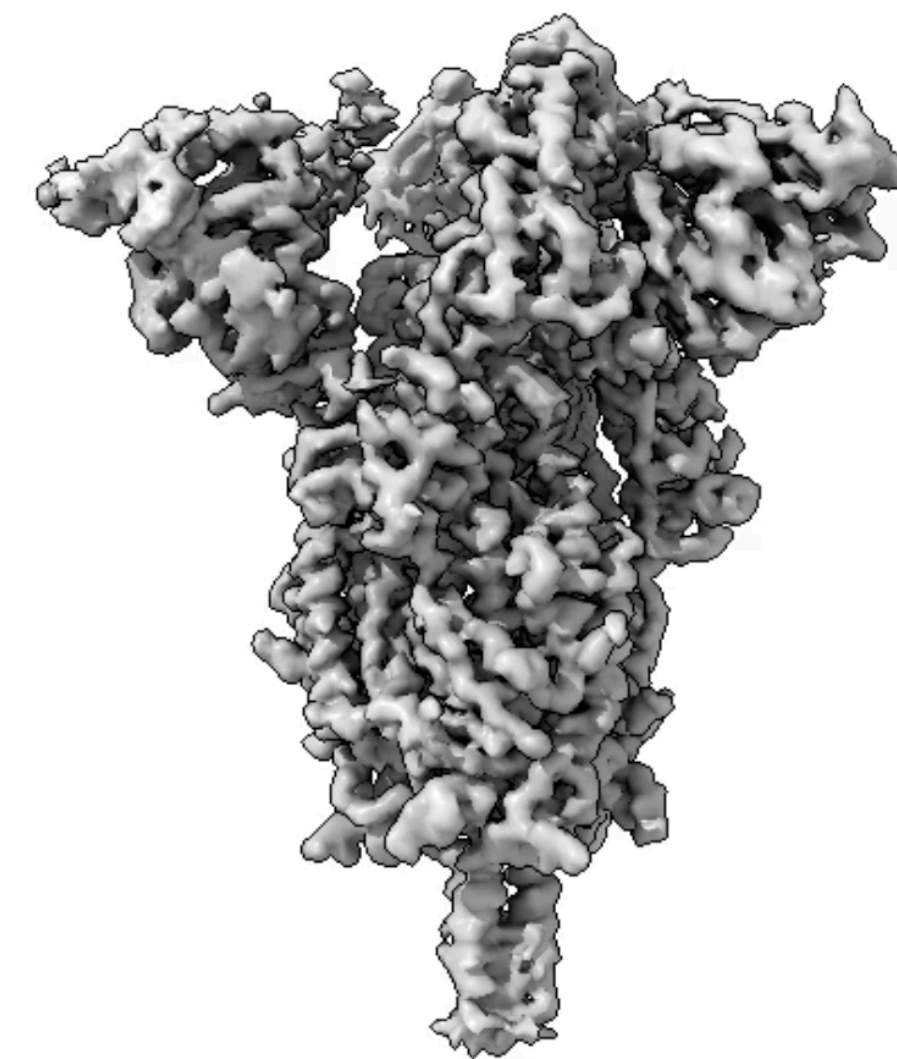
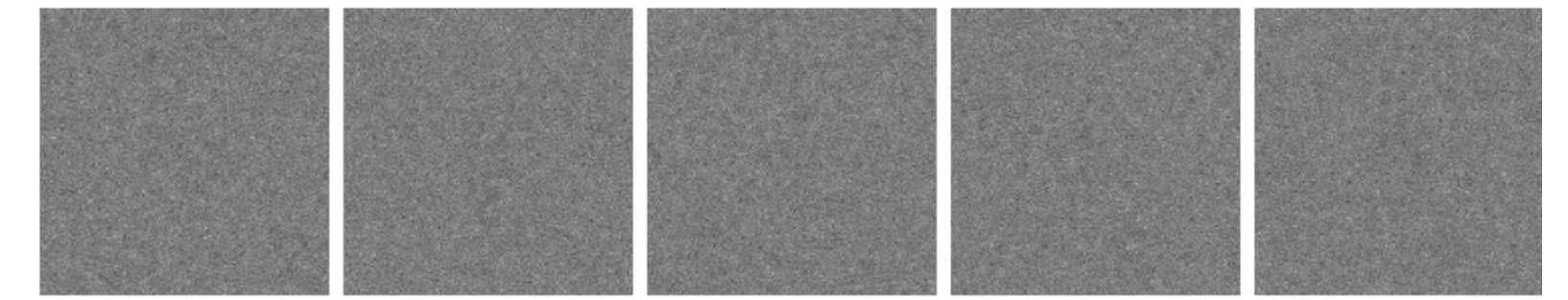
Roy R. Lederman* and Amit Singer†

April 11, 2017

Abstract

Single particle cryo-electron microscopy (EM) is an increasingly popular method for determining the 3-D structure of macromolecules from noisy 2-D images of single macromolecules whose orientations and positions are random and unknown. One of the great opportunities in cryo-EM

The cryo-EM reconstruction task



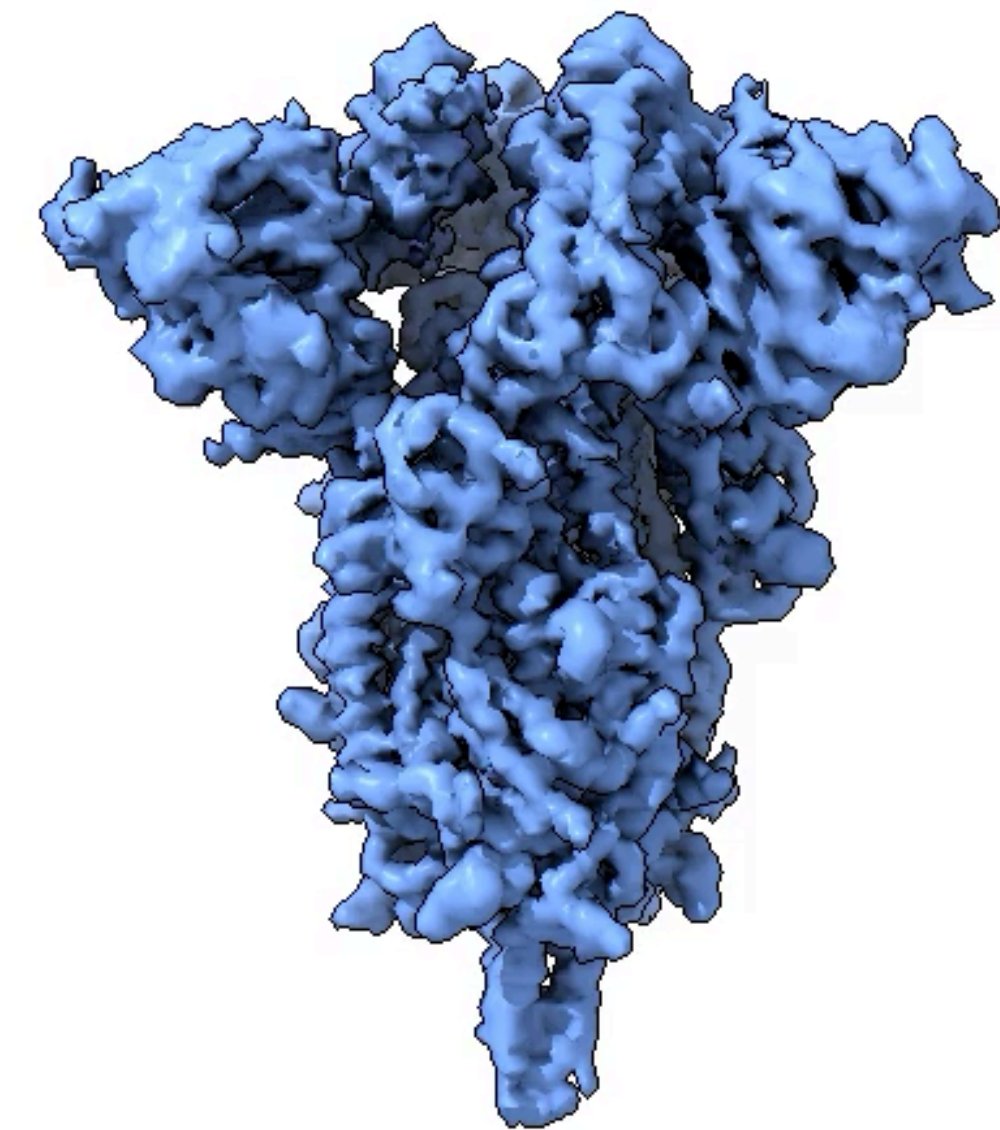
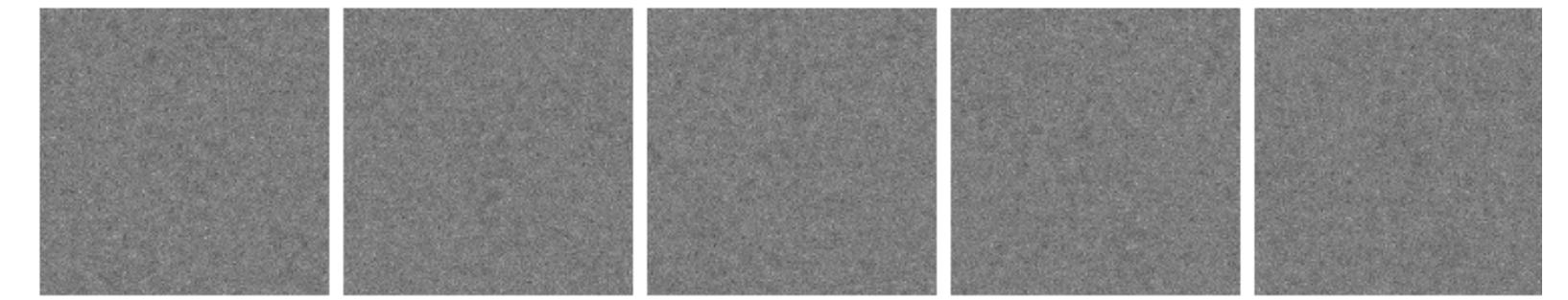
Walls et al, 2020

Cryo-EM structure of the SARS CoV-2 Spike protein

Outline

- **Motivation:** Why do we care about protein structure?
- **Background:** Cryo-EM reconstruction & the heterogeneity problem
- **CryoDRGN:** Neural 3D reconstruction of dynamic protein structure with cryoDRGN ❄️🐉

The cryo-EM reconstruction task

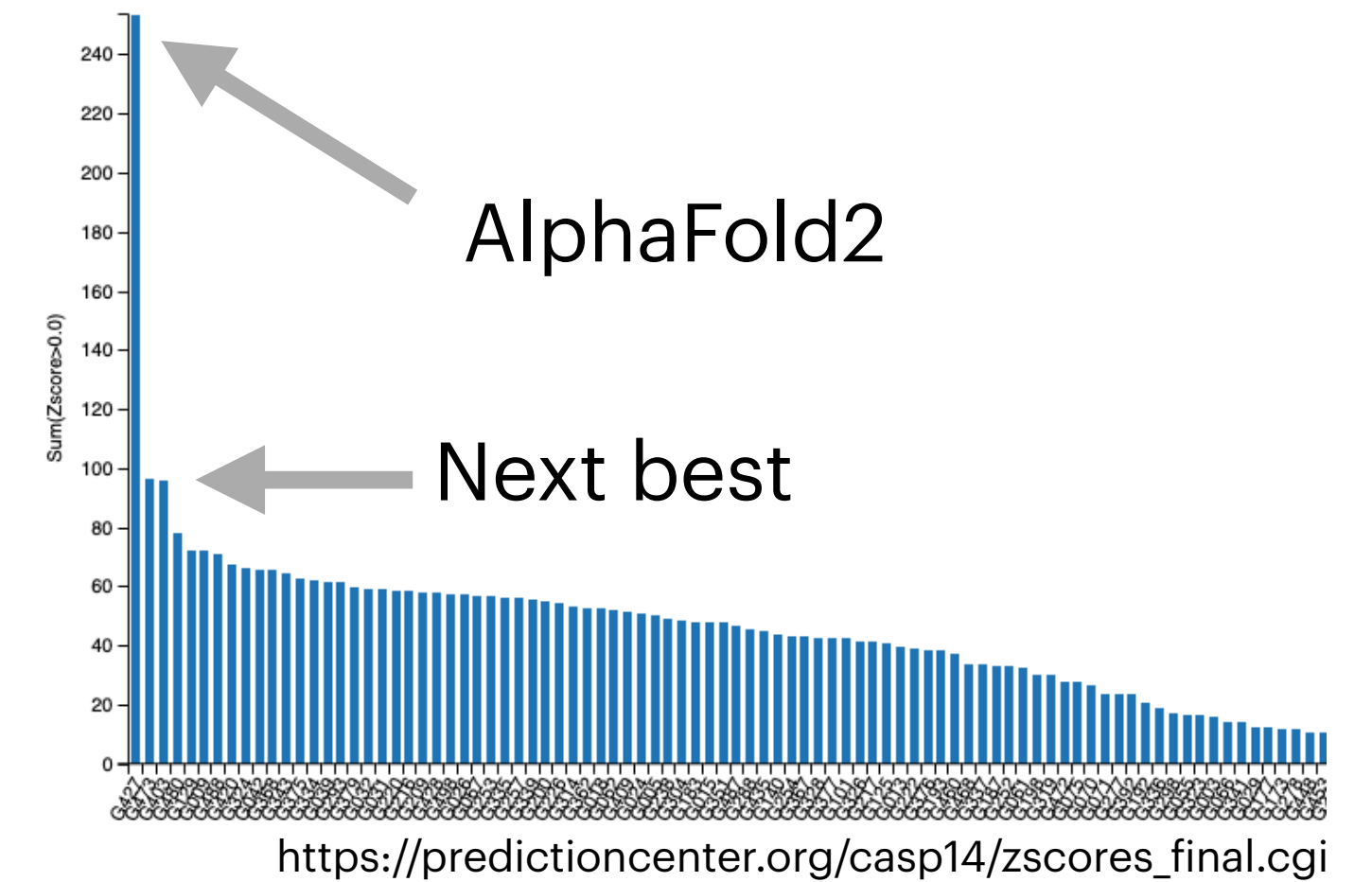


cryoDRGN trajectory of the SARS CoV-2 Spike protein

Outline

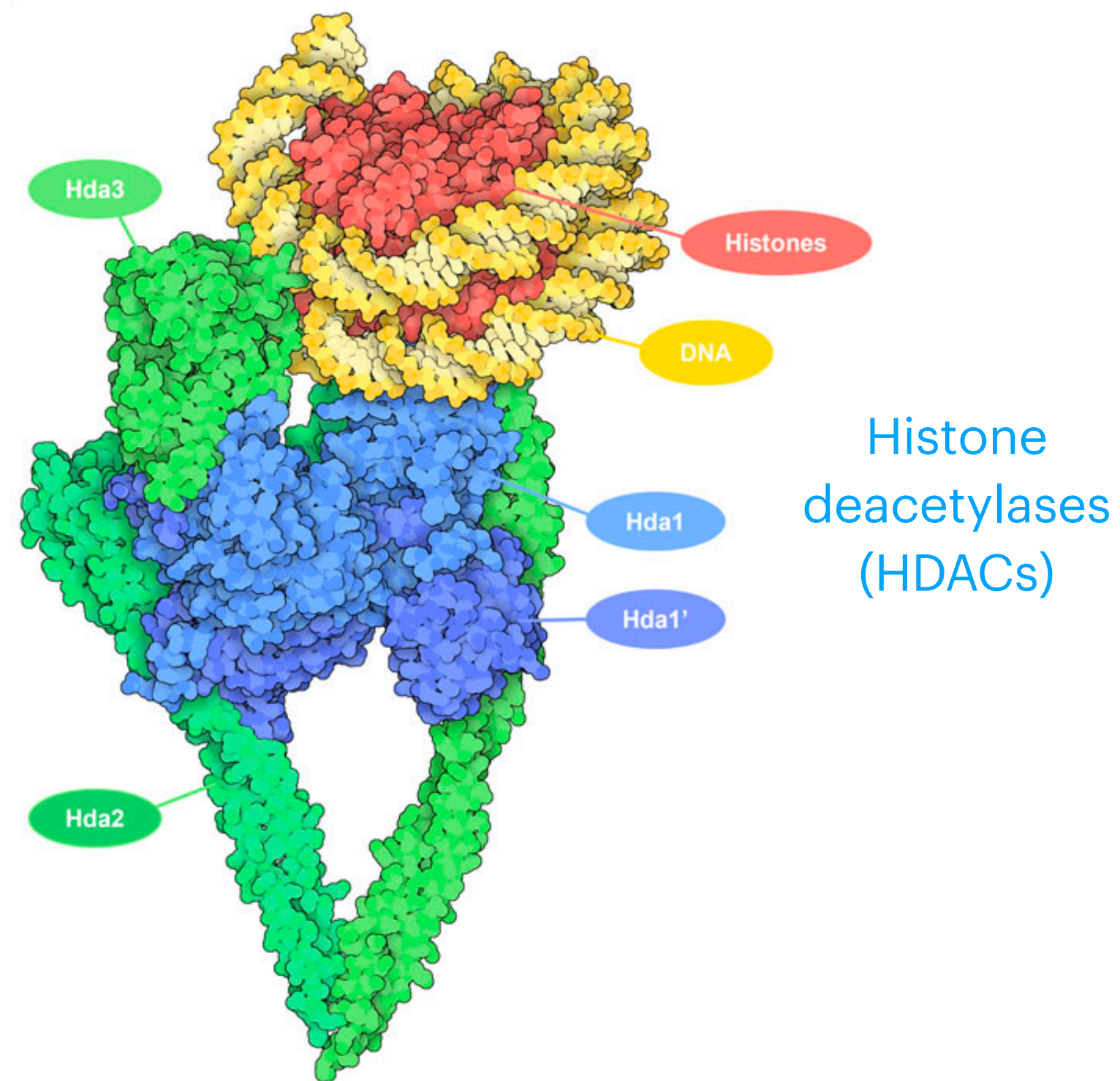
- **Motivation:** Why do we care about protein structure?
- **Background:** Cryo-EM reconstruction & the heterogeneity problem
- **CryoDRGN:** Neural 3D reconstruction of dynamic protein structure with cryoDRGN ❄️🐉
- **Future directions:** Machine learning for structure determination at the proteome scale

CASP14 Results, Dec 2020

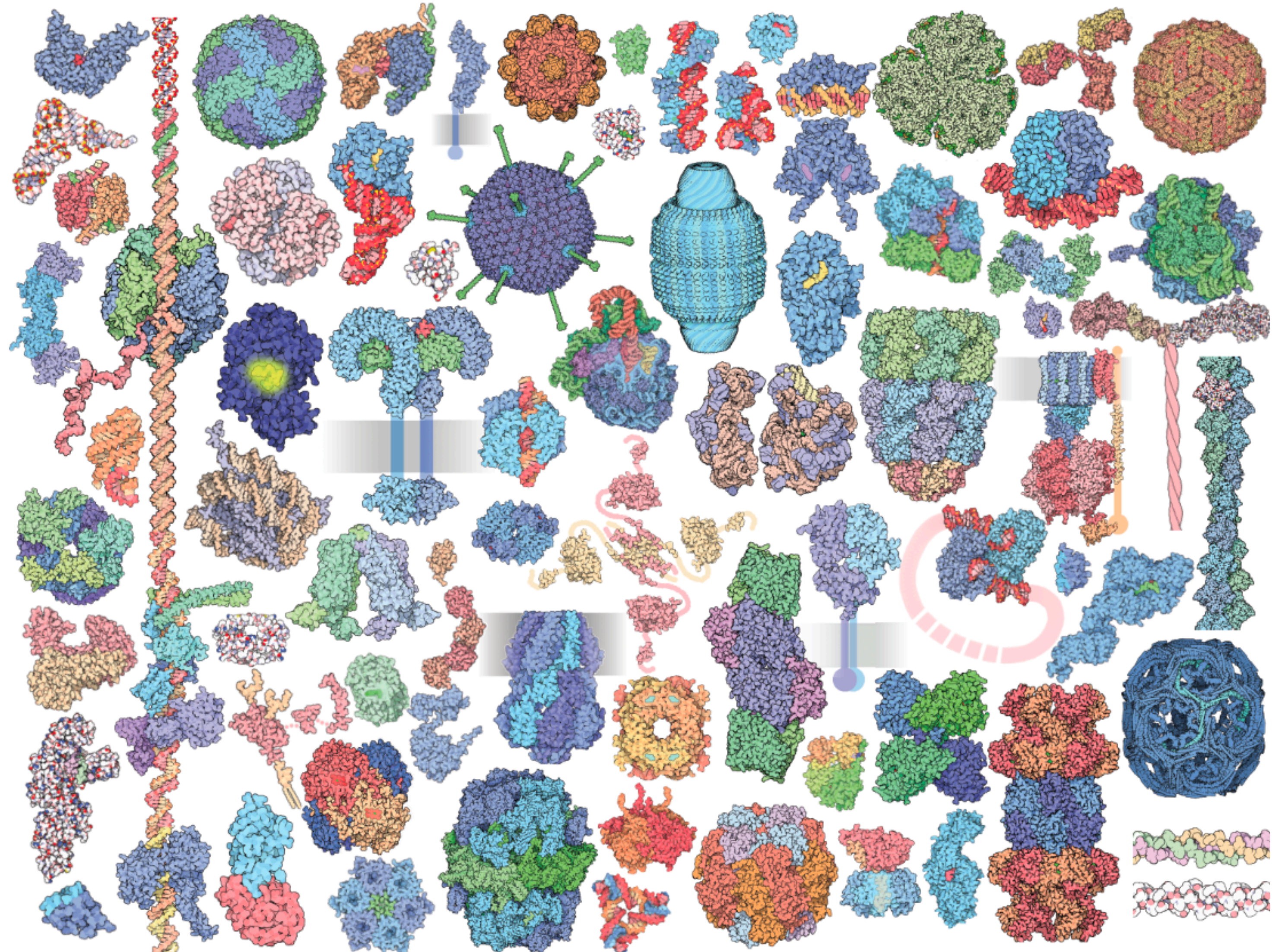


All essential biological processes are carried out by proteins and protein complexes

- Fundamental molecules of life
- Medicine and health
- Nanotech and biotech



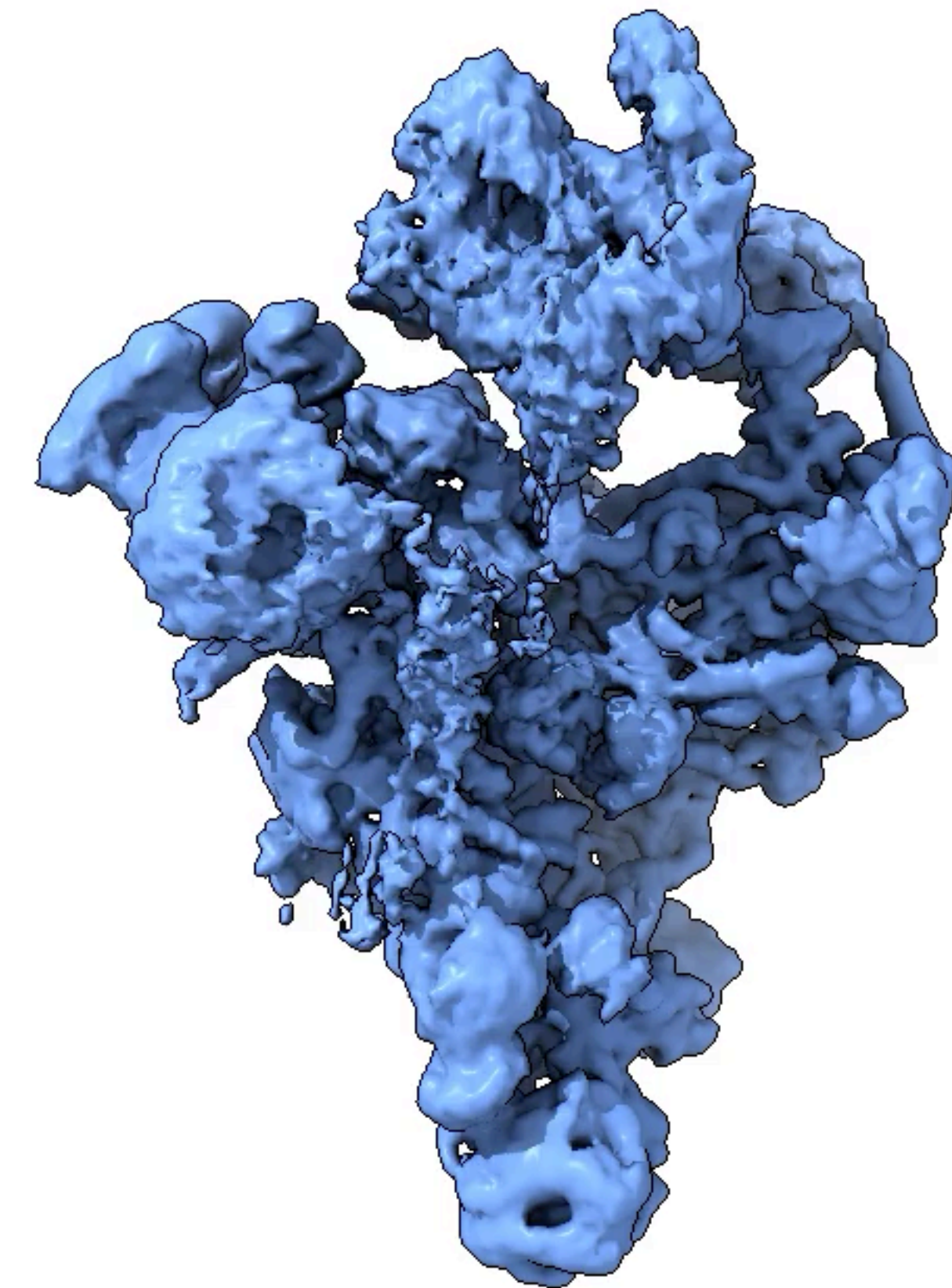
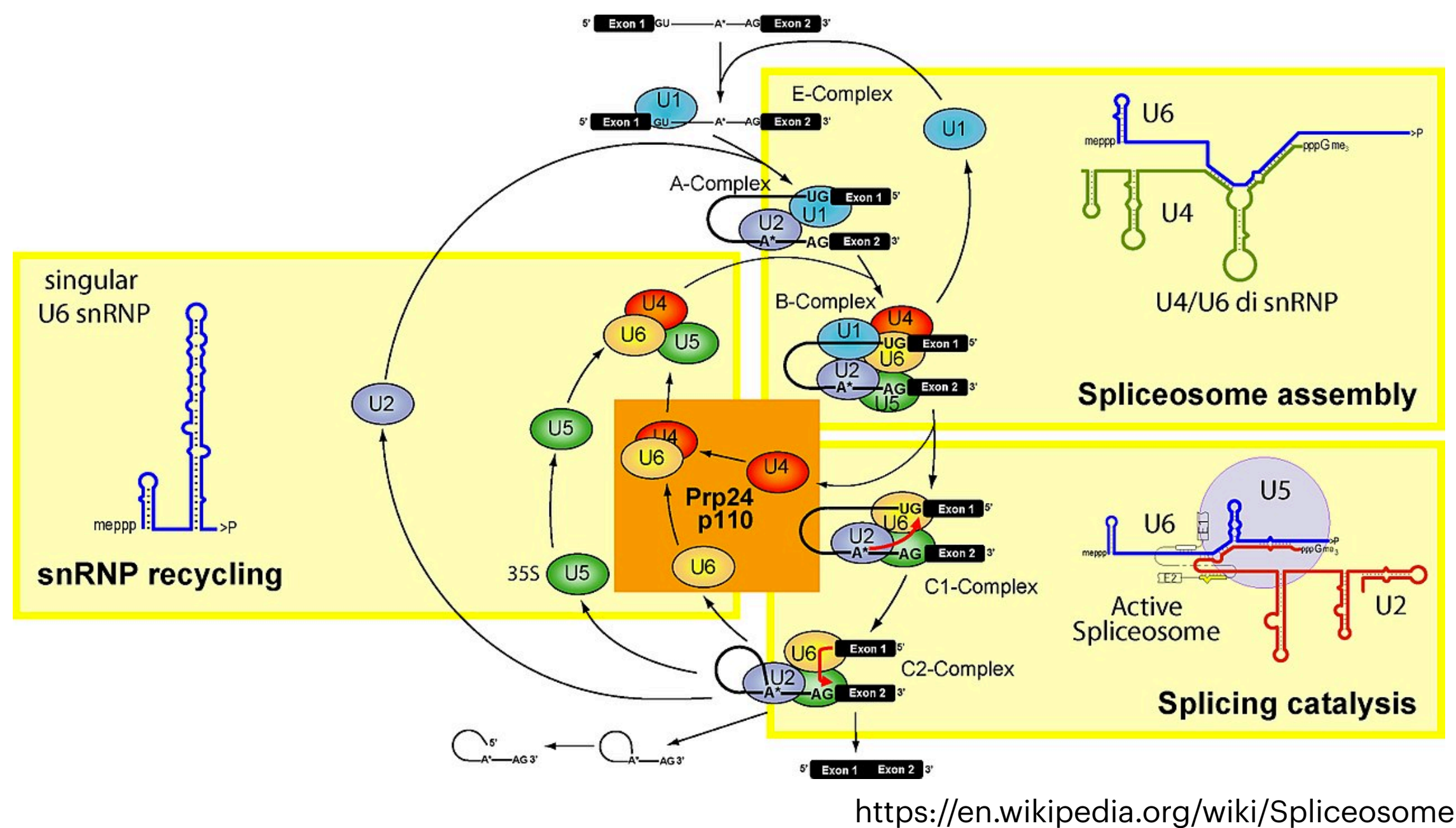
PDB-101 Molecule of the Month



All essential biological processes are carried out by proteins and protein complexes

... which are dynamic macromolecular machines

Spliceosome splicing cycle

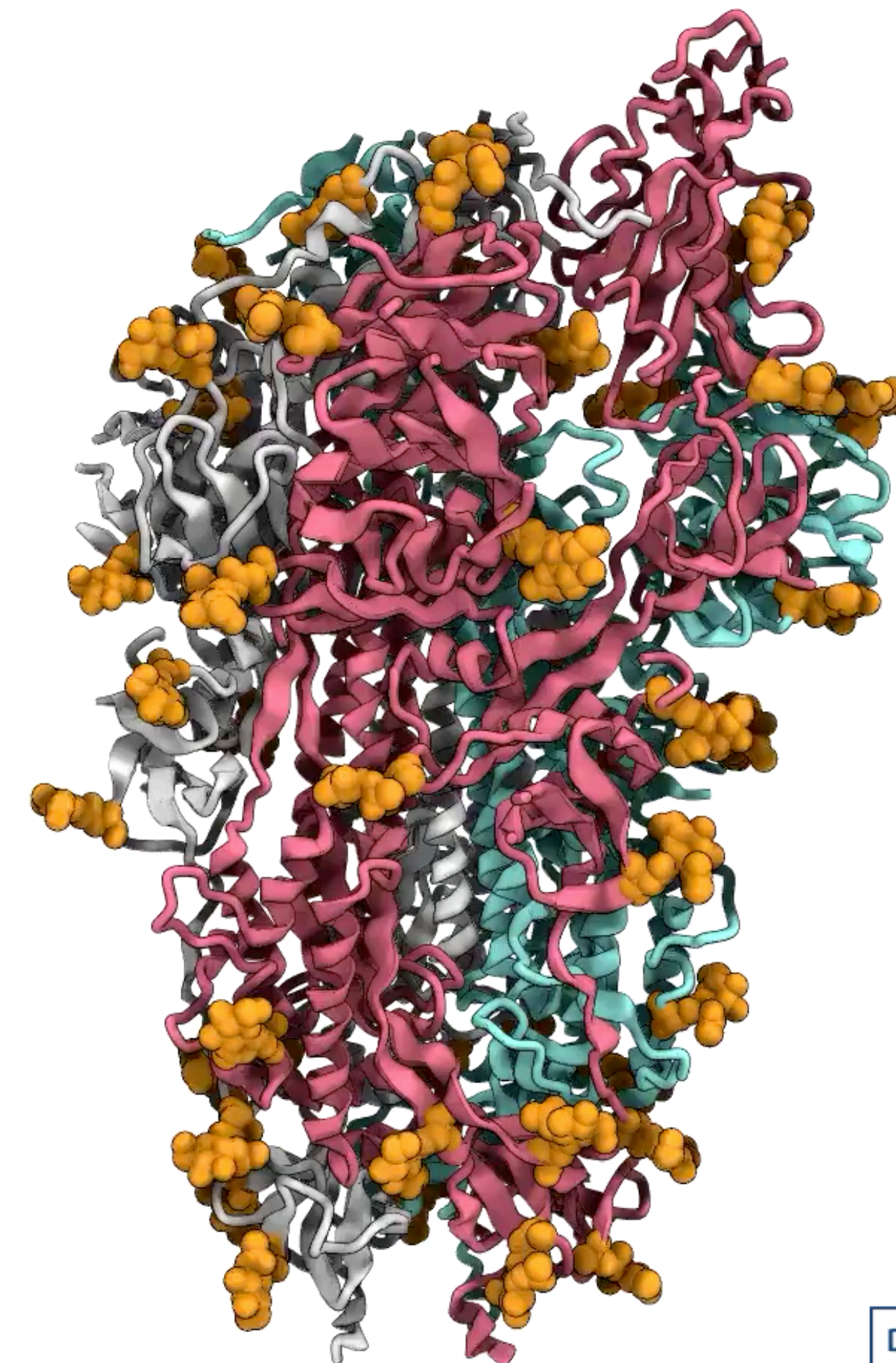


cryoDRGN trajectory of the pre-catalytic spliceosome

Techniques to study molecular motions are limited

- Nuclear magnetic resonance (NMR) spectroscopy
 - Small proteins (<100 AA in length)
- Electric field crystallography (EF-X), multi-temperature and XFEL crystallography
 - Requires sample crystallization
- Computational modeling
 - Molecular dynamics simulations
 - Hacking AlphaFold?
- **Cryo-electron microscopy (cryo-EM)**

0.0 μ s

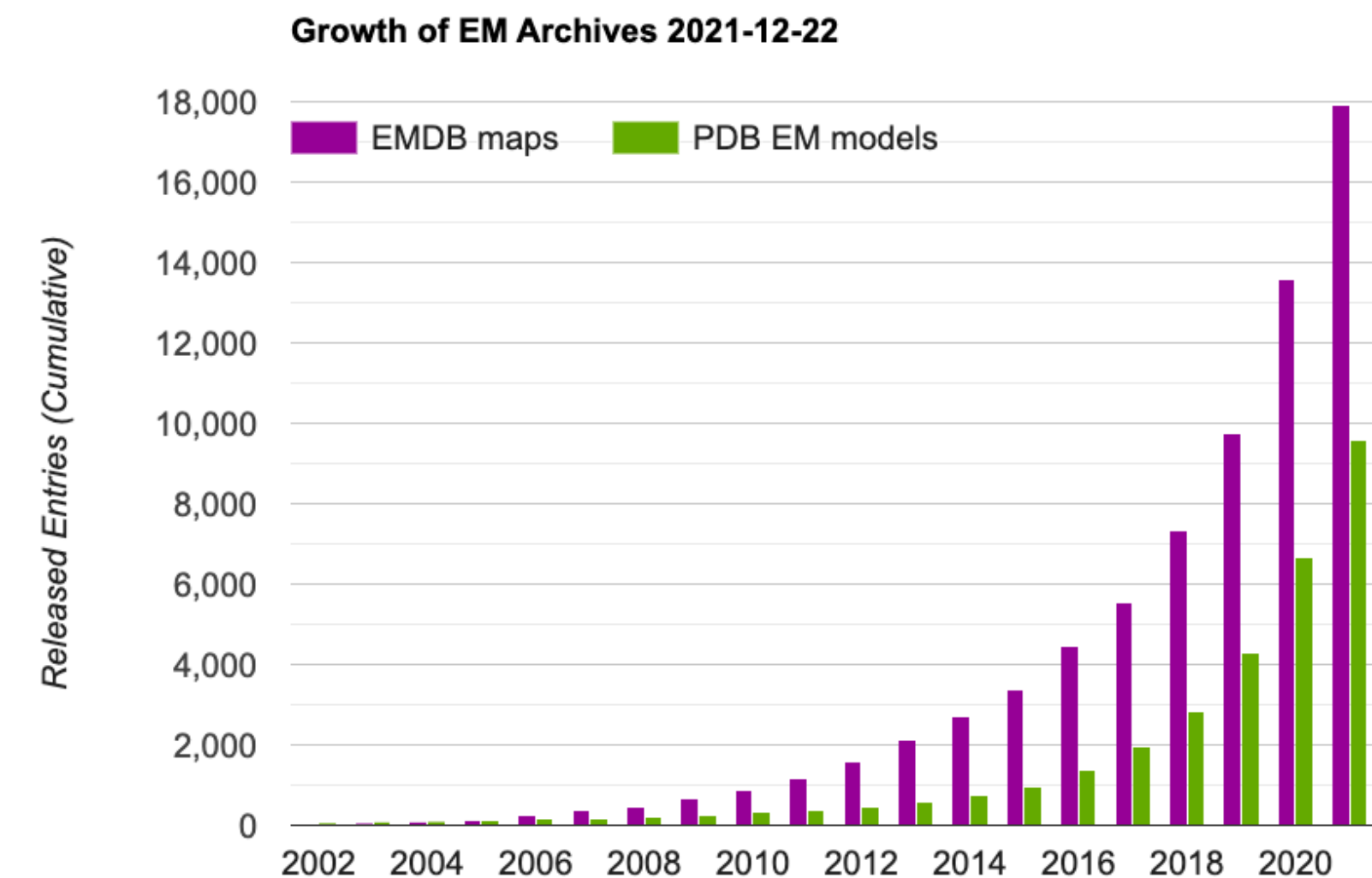
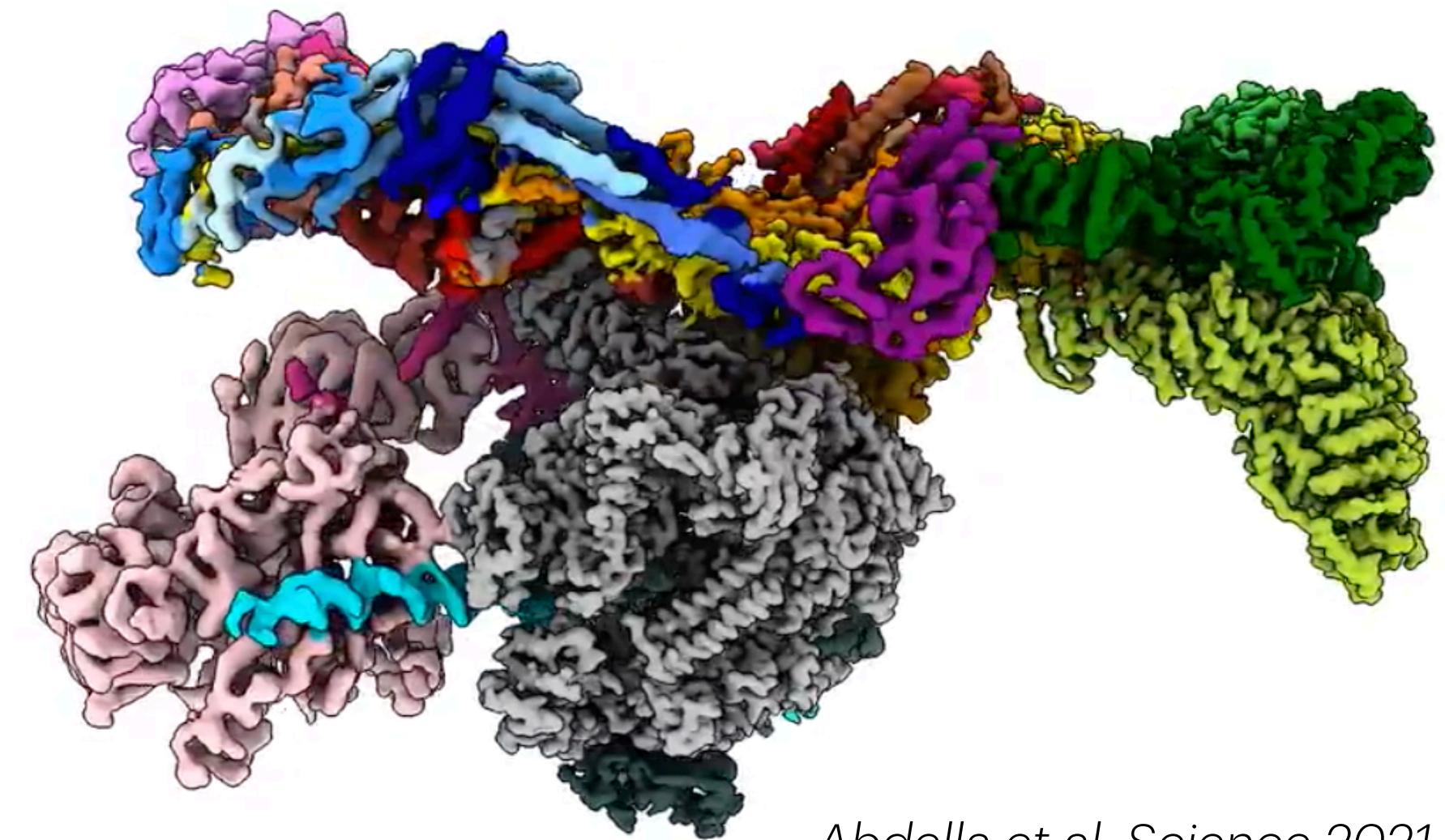


MD simulation of
SARS CoV-2
Spike



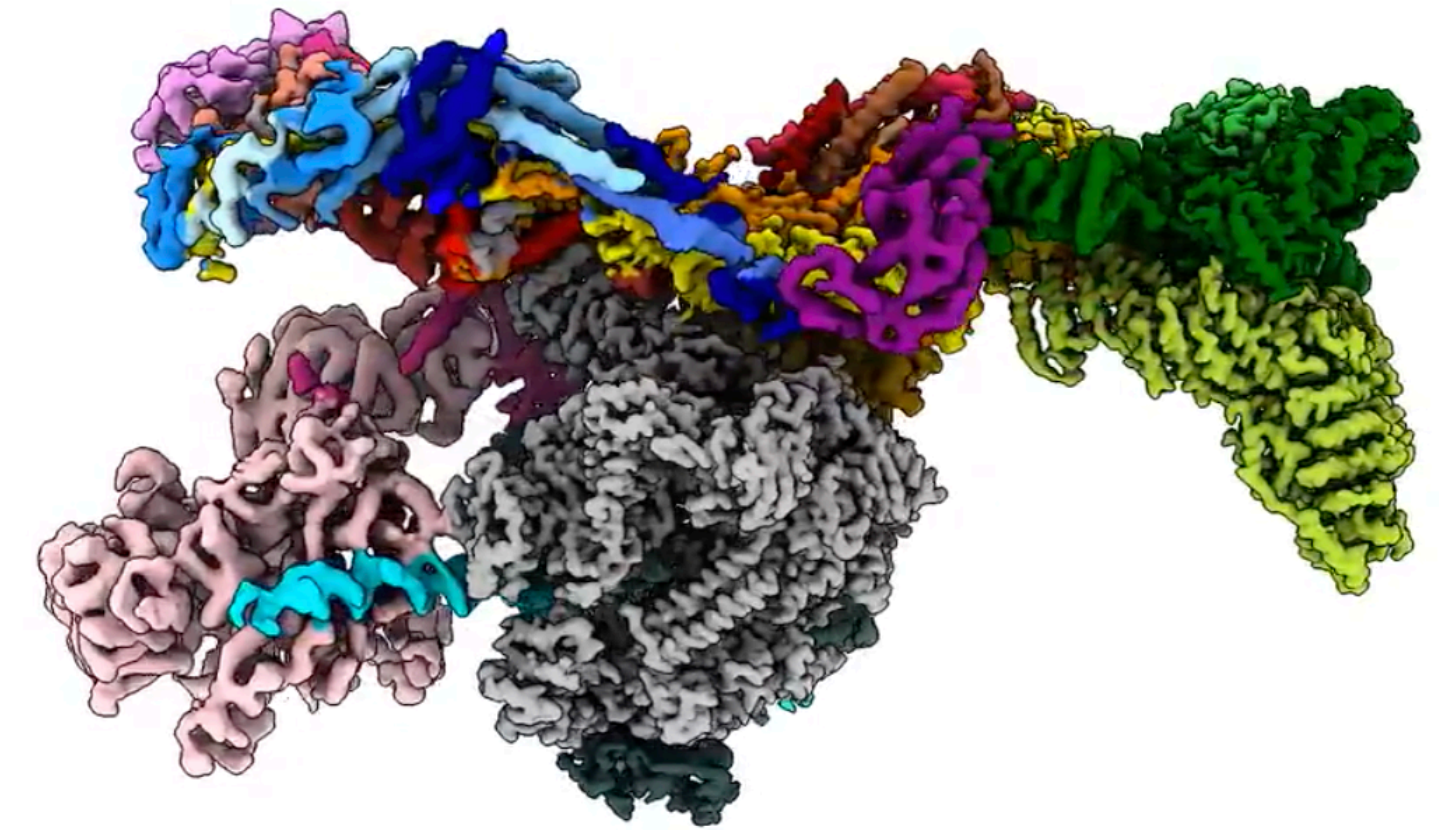
The ongoing cryo-EM “resolution revolution”

- 2017 Nobel Prize in Chemistry
- Cryo-EM has opened up new areas of structural biology
- Recent hardware and software breakthroughs:
 - **Hardware:** direct electron detectors
 - **Software:** New reconstruction algorithms, GPU compute
 - **Faster:** Automation and democratization of cryo-EM imaging
- New computational challenges and opportunities



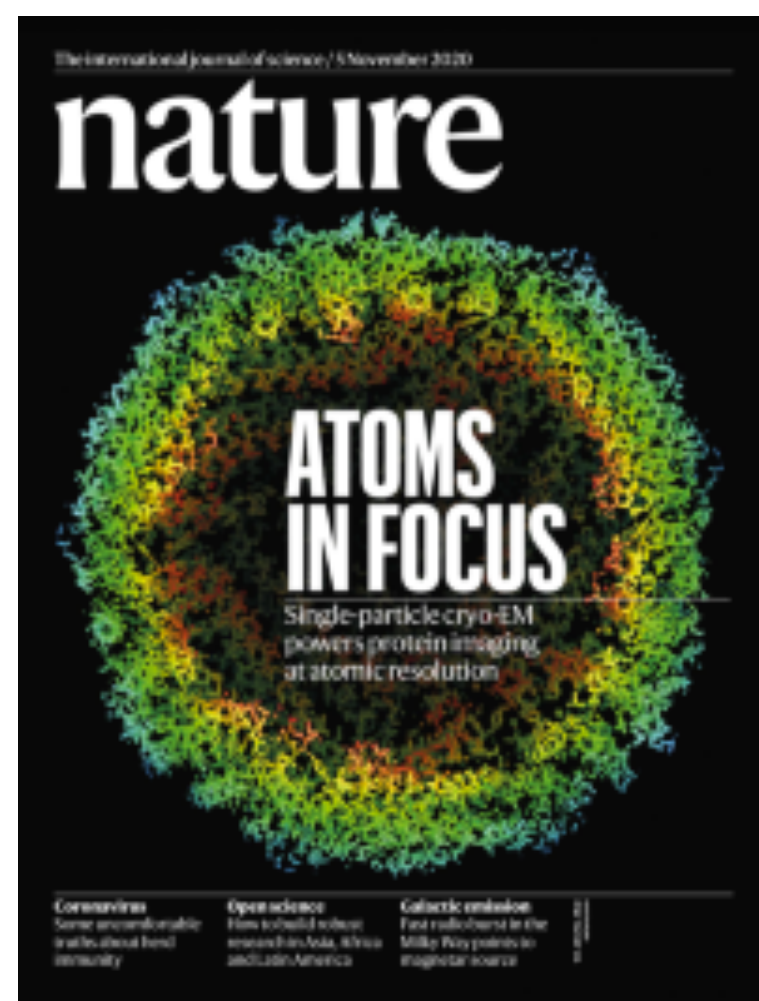
Frontiers of single particle cryo-EM

- Higher resolution structures
- Small proteins
- Time-resolved cryo-EM
- **Large, dynamic complexes**
 - (MDa scale, 10s-100s of proteins)

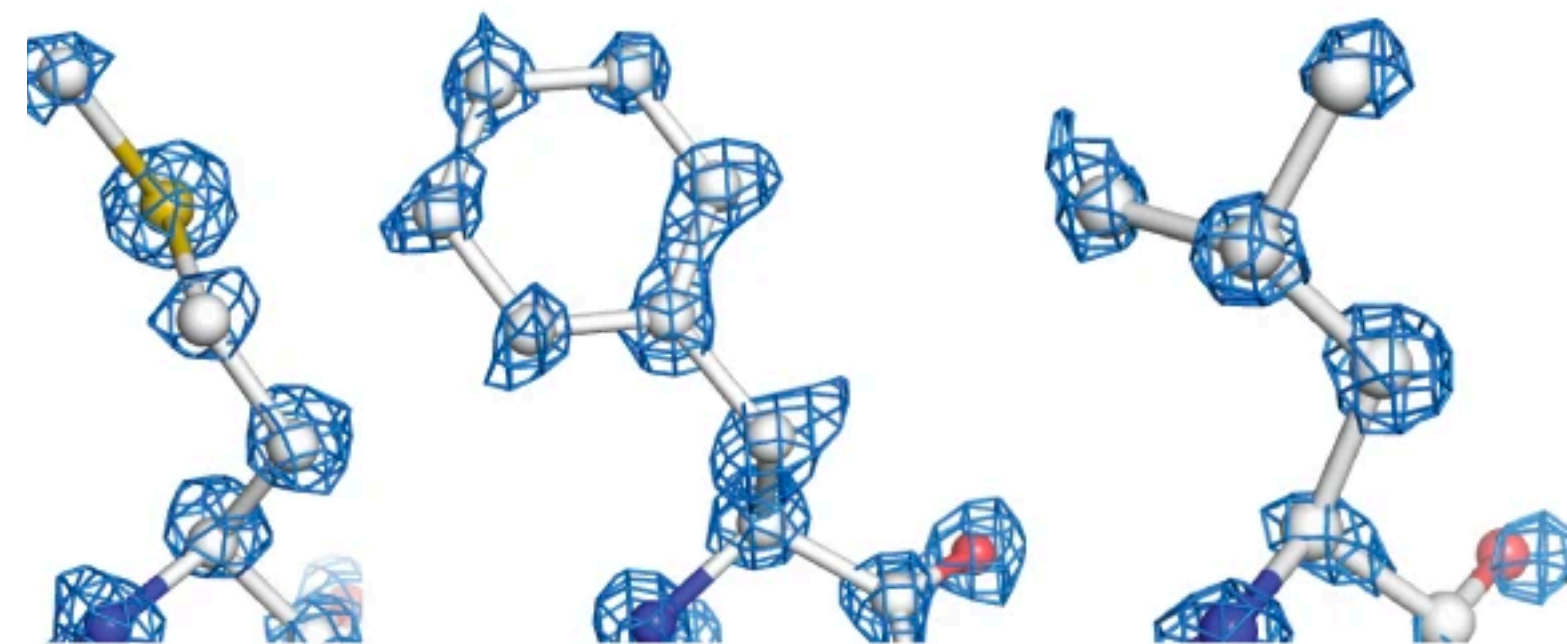


cPIC: RNA Pol II GTFs DNA
TFIIH: cTFIIH Mat1 cyclin-H CDK7
MedHead: Med6 Med8 Med11 Med17 Med18 Med20 Med22 Med27 Med28 Med29 Med30
MedMiddle: Med1 Med4 Med7 Med9 Med10 Med19 Med21 Med26 Med31

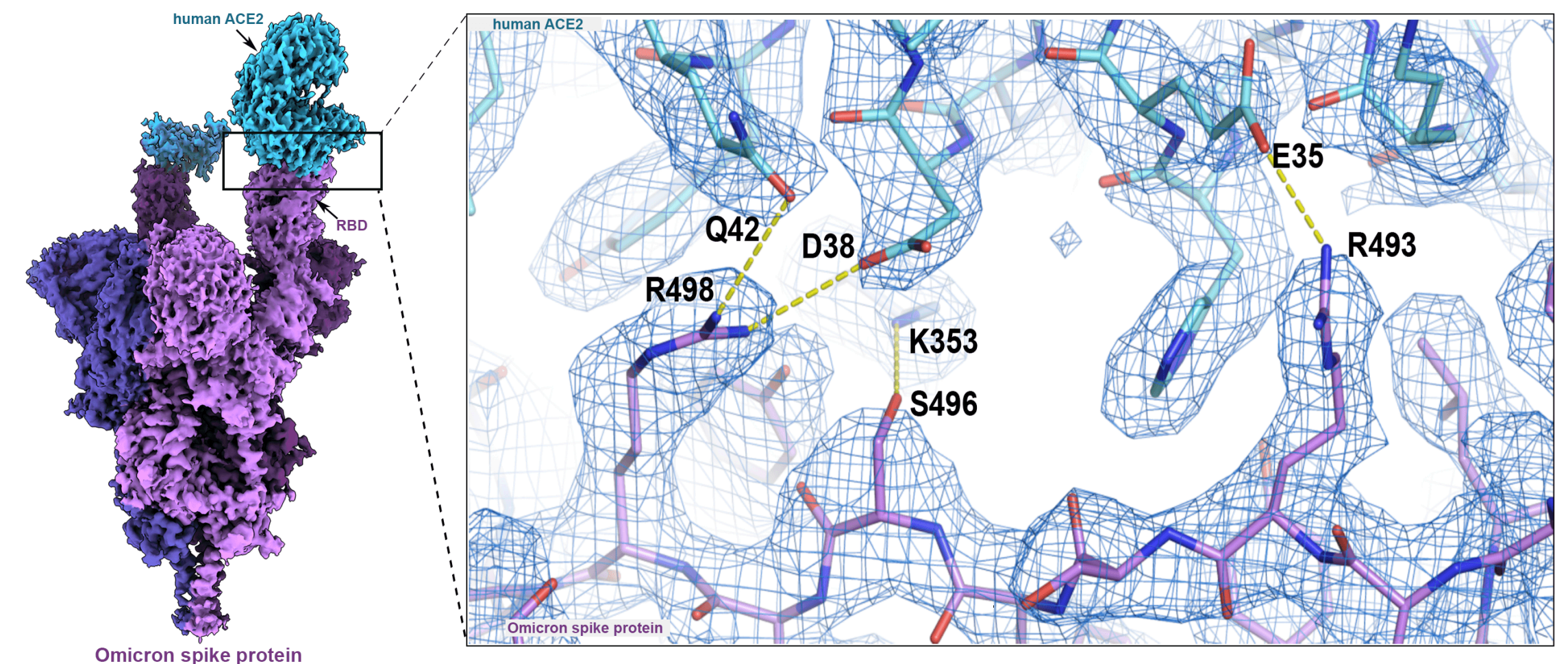
<https://twitter.com/DanielHurdiss/status/1372659832780623872>



November 2020



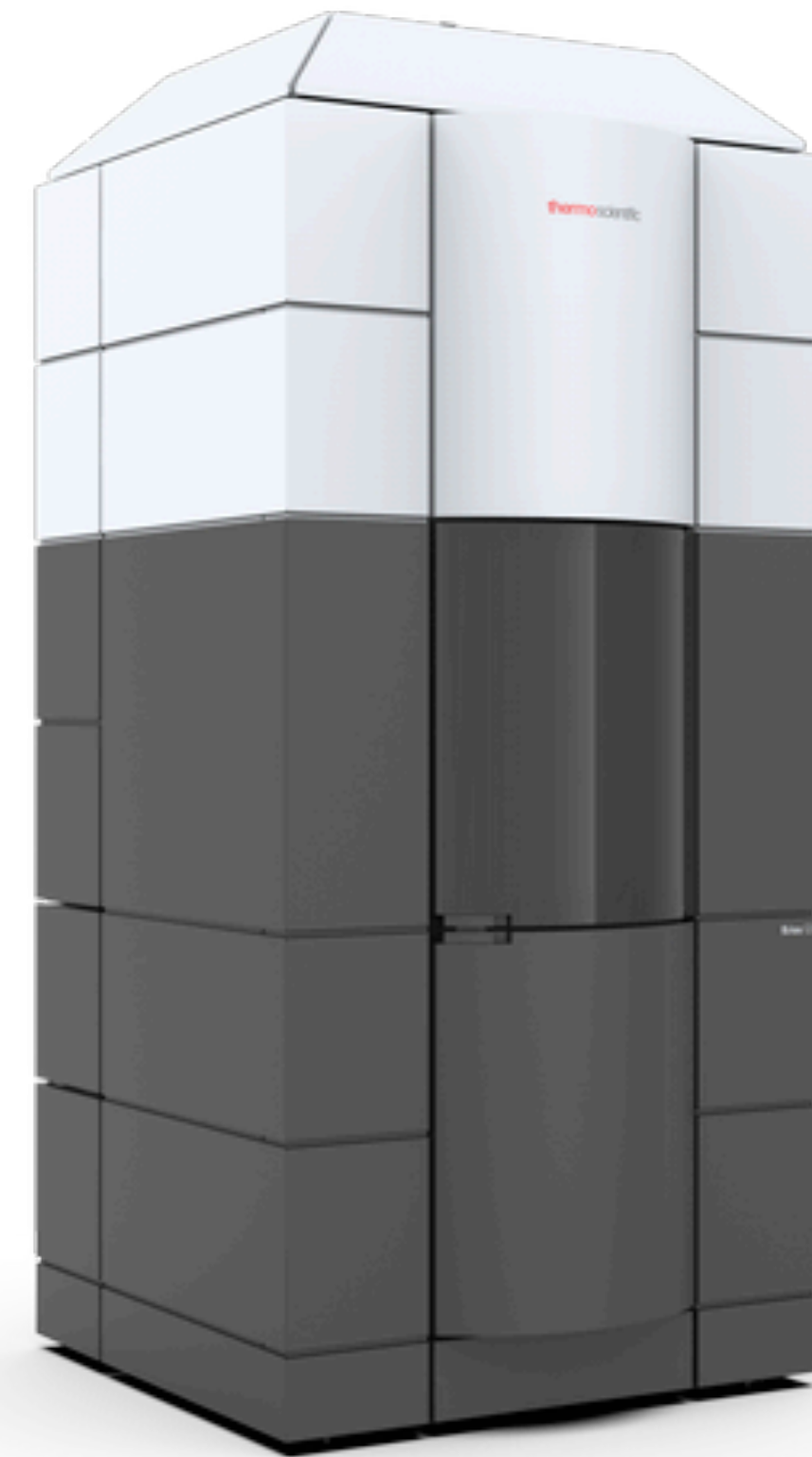
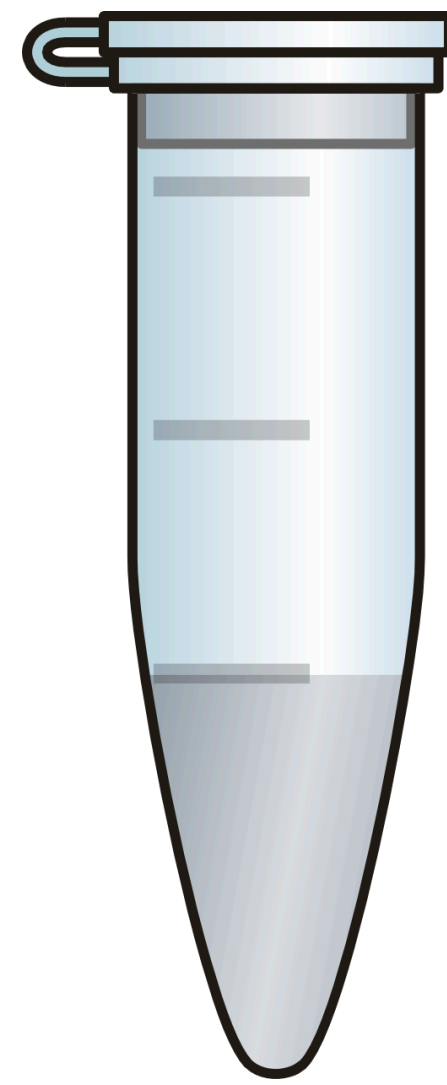
Nakane et al Nature 2020



Omicron spike protein structure. Mannar et al bioRxiv, 2022

The cryo-EM image processing pipeline: From micrograph to atomic coordinates

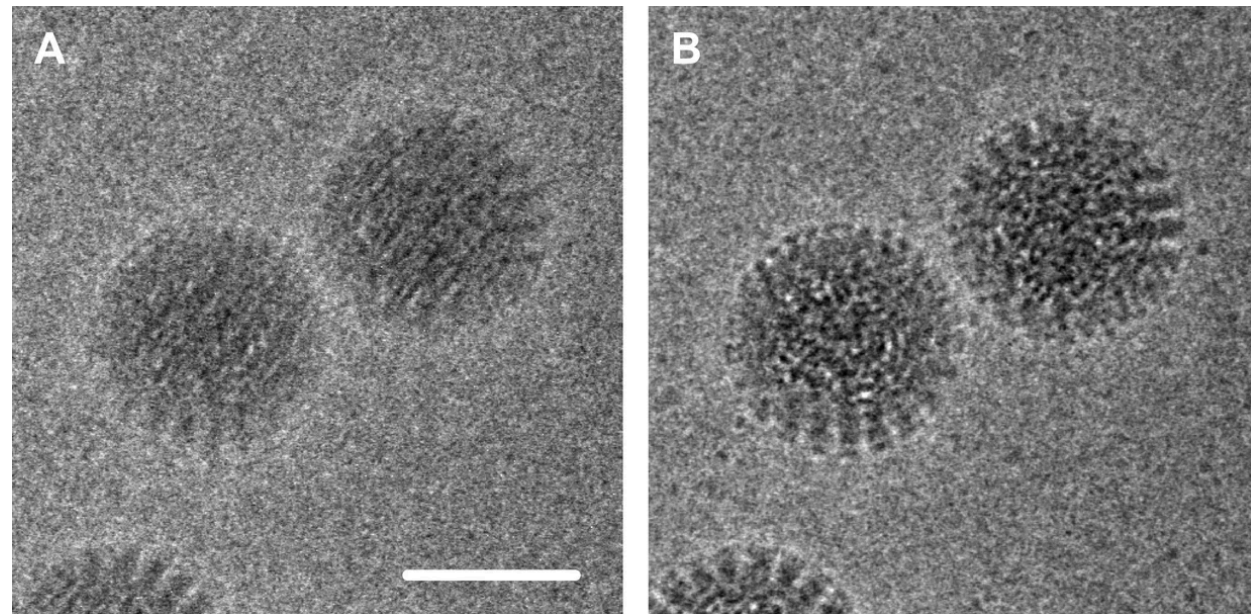
[Step 0) Sample preparation and imaging]



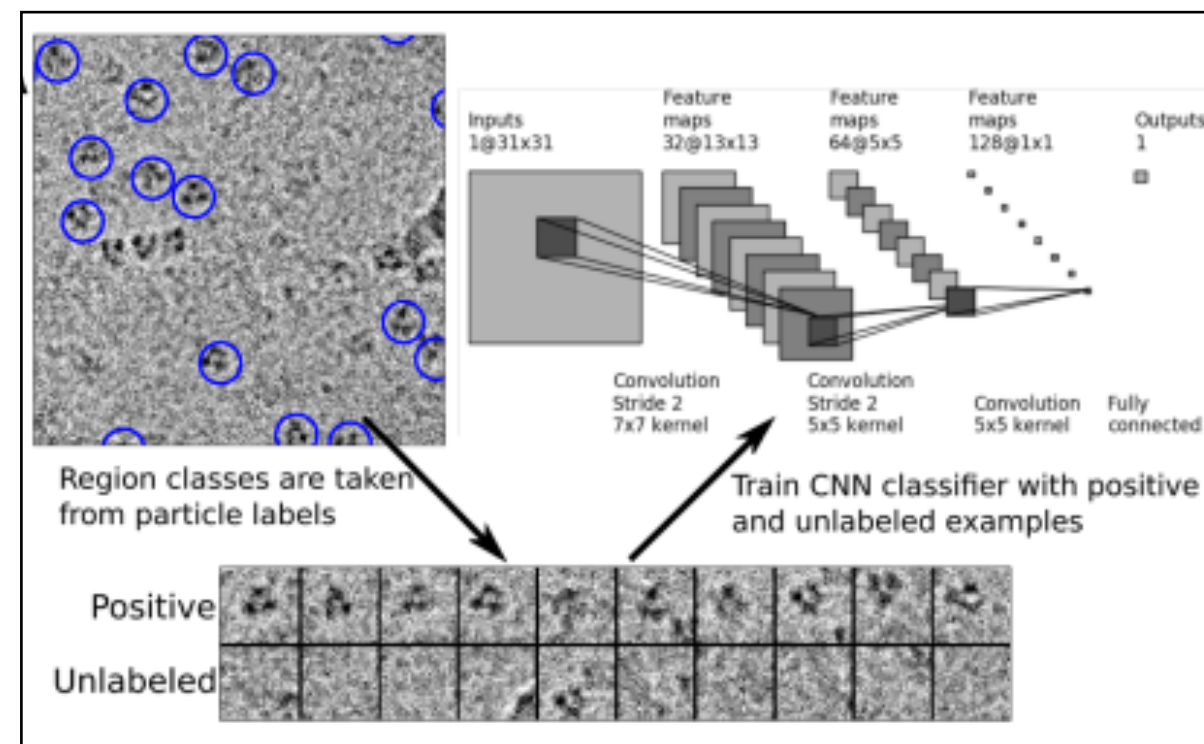
The cryo-EM image processing pipeline: From micrograph to atomic coordinates

[Step 0) Sample preparation and imaging]

1) Micrograph pre-processing



Grigorieff, 2013

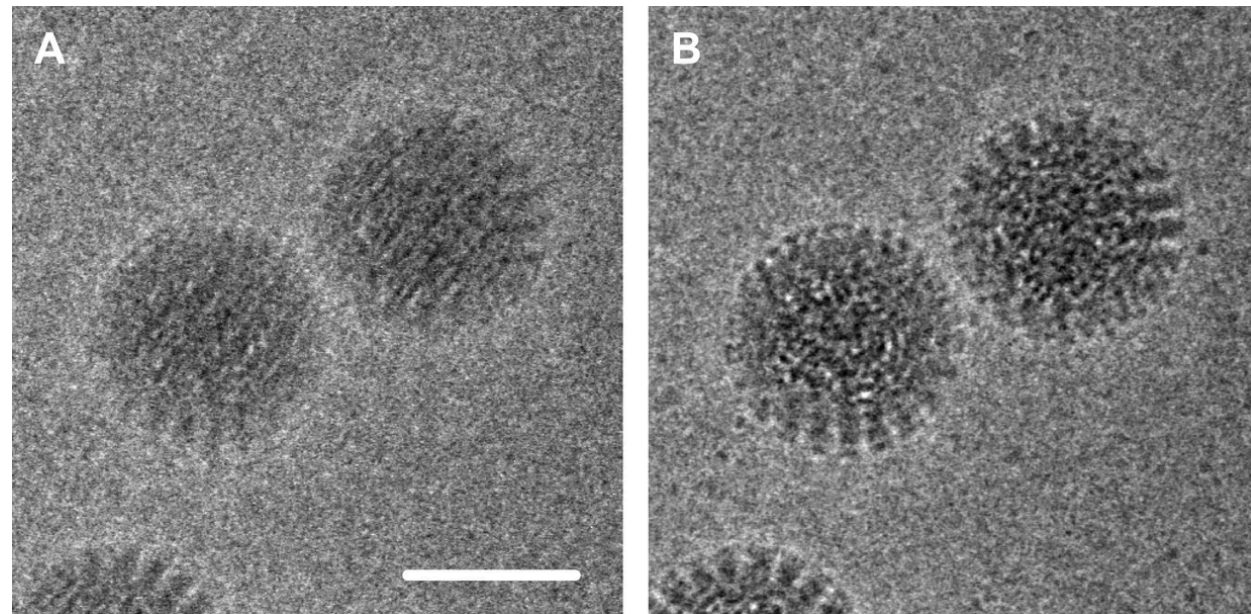


Topaz, Bepler et al, 2019

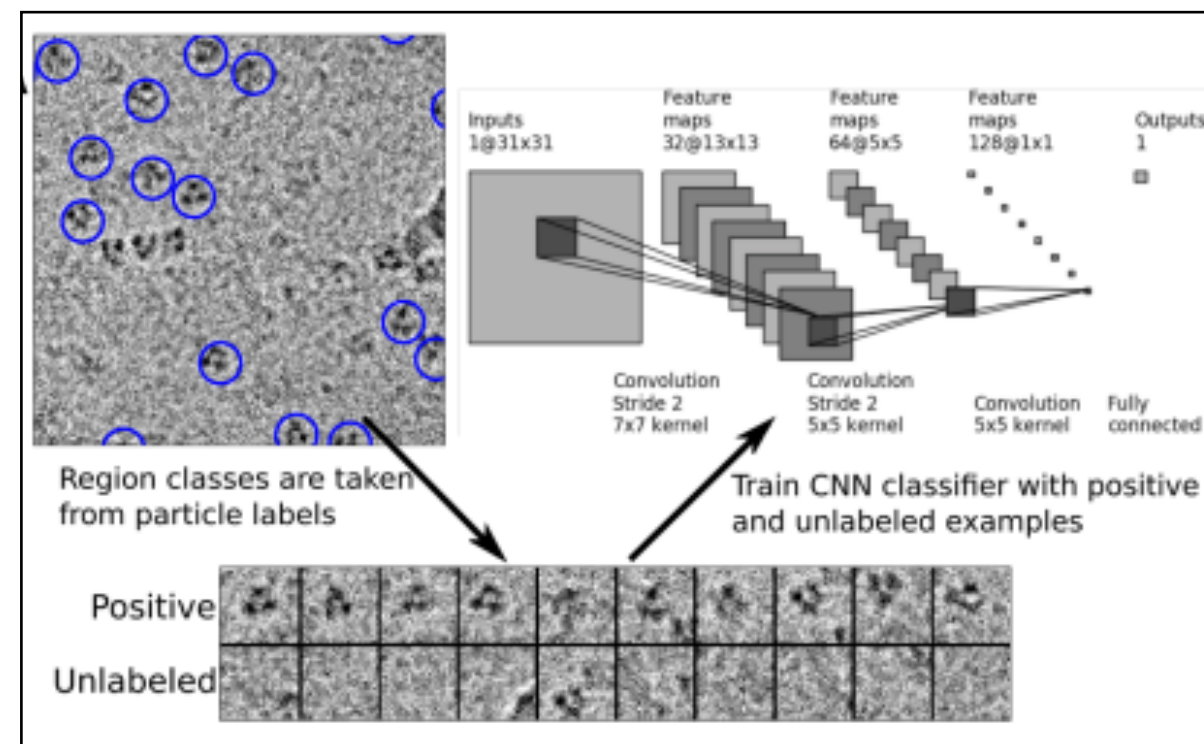
The cryo-EM image processing pipeline: From micrograph to atomic coordinates

[Step 0) Sample preparation and imaging]

1) Micrograph pre-processing



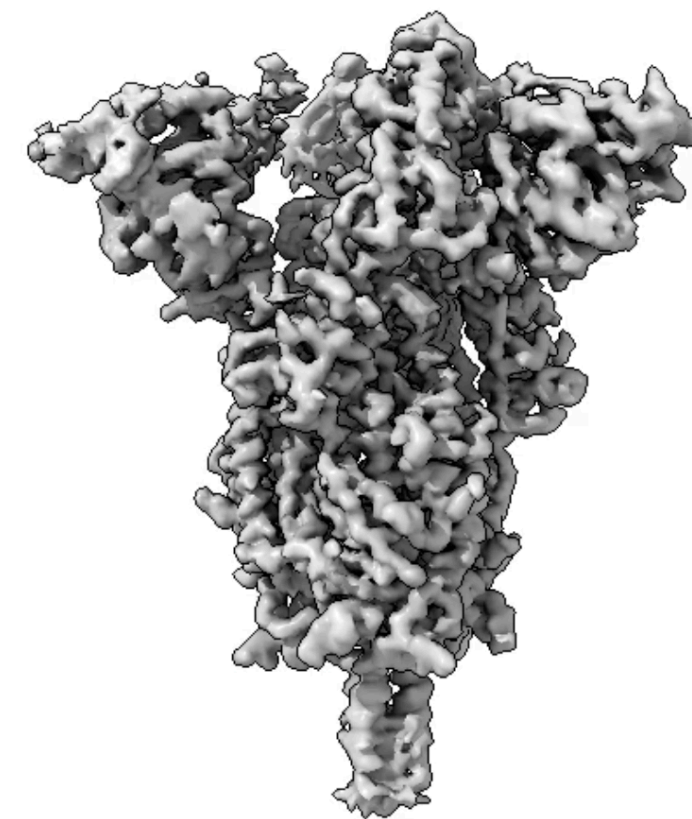
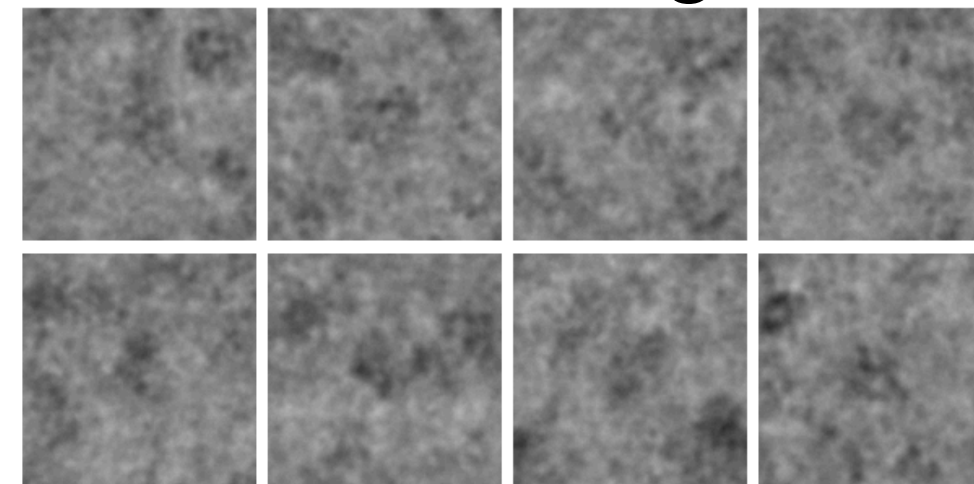
Grigorieff, 2013



Topaz, Bepler et al, 2019

2) 2D to 3D reconstruction

10⁴-10⁷ images

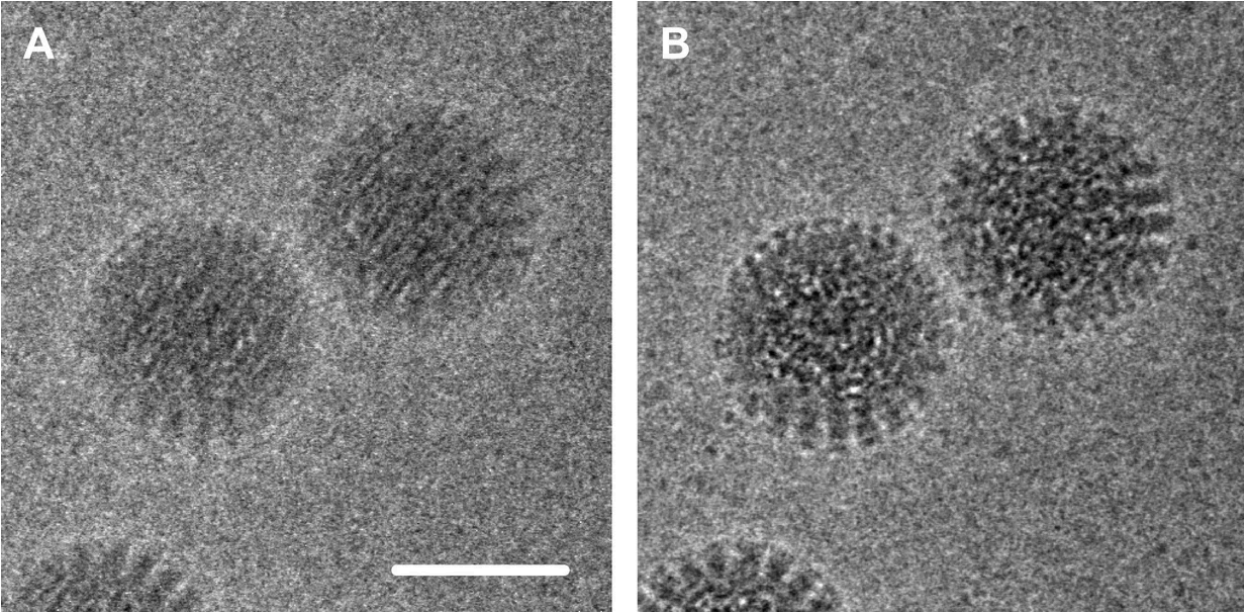


Walls et al, 2020

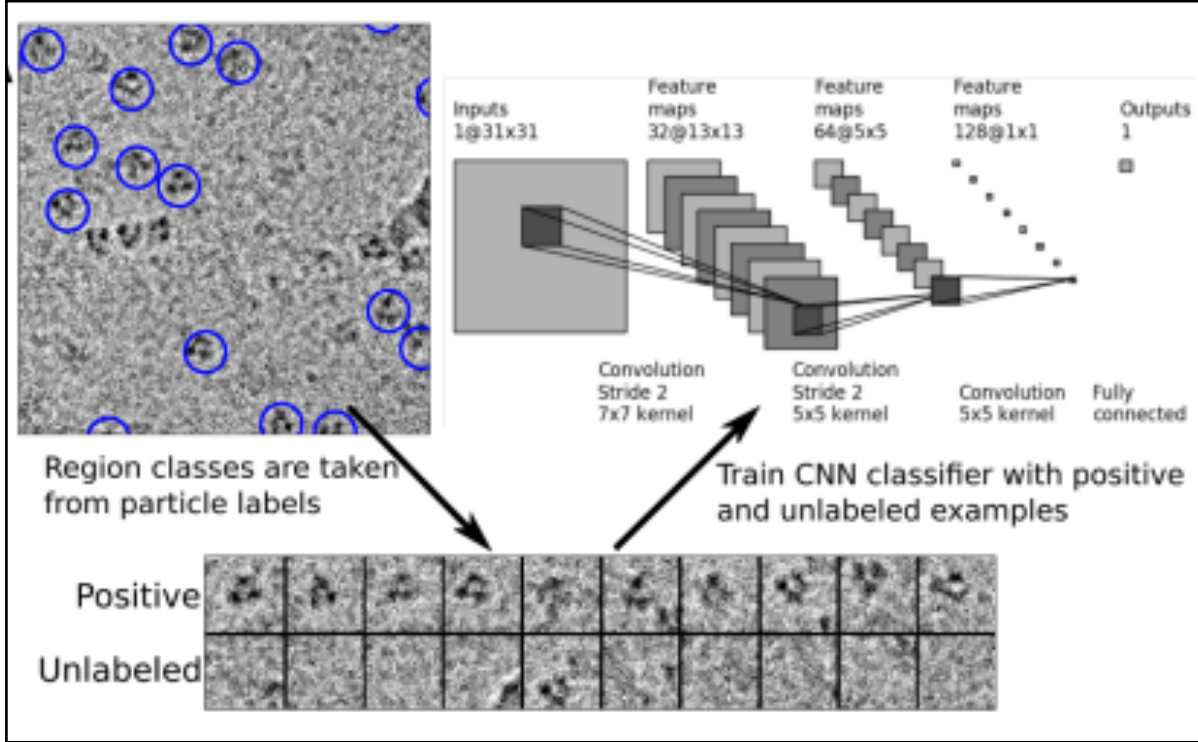
The cryo-EM image processing pipeline: From micrograph to atomic coordinates

[Step 0) Sample preparation and imaging]

1) Micrograph pre-processing



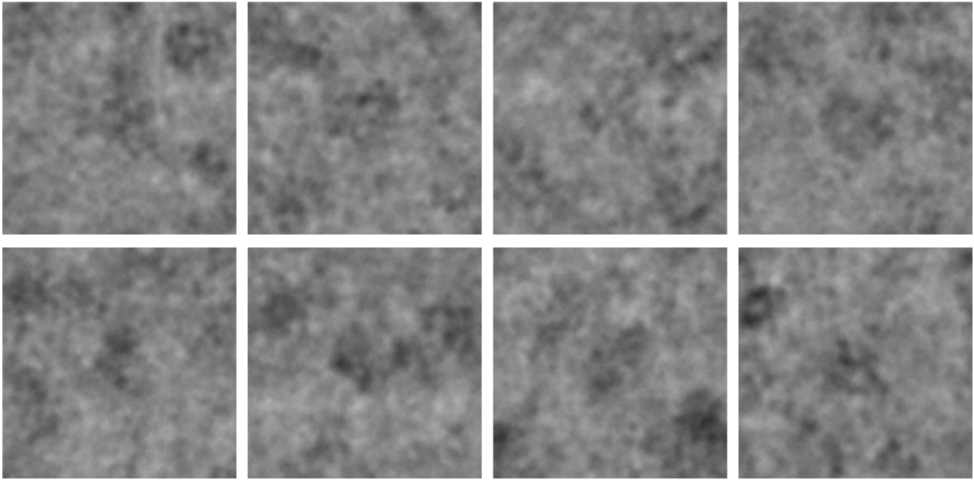
Grigorieff, 2013



Topaz, Bepler et al, 2019

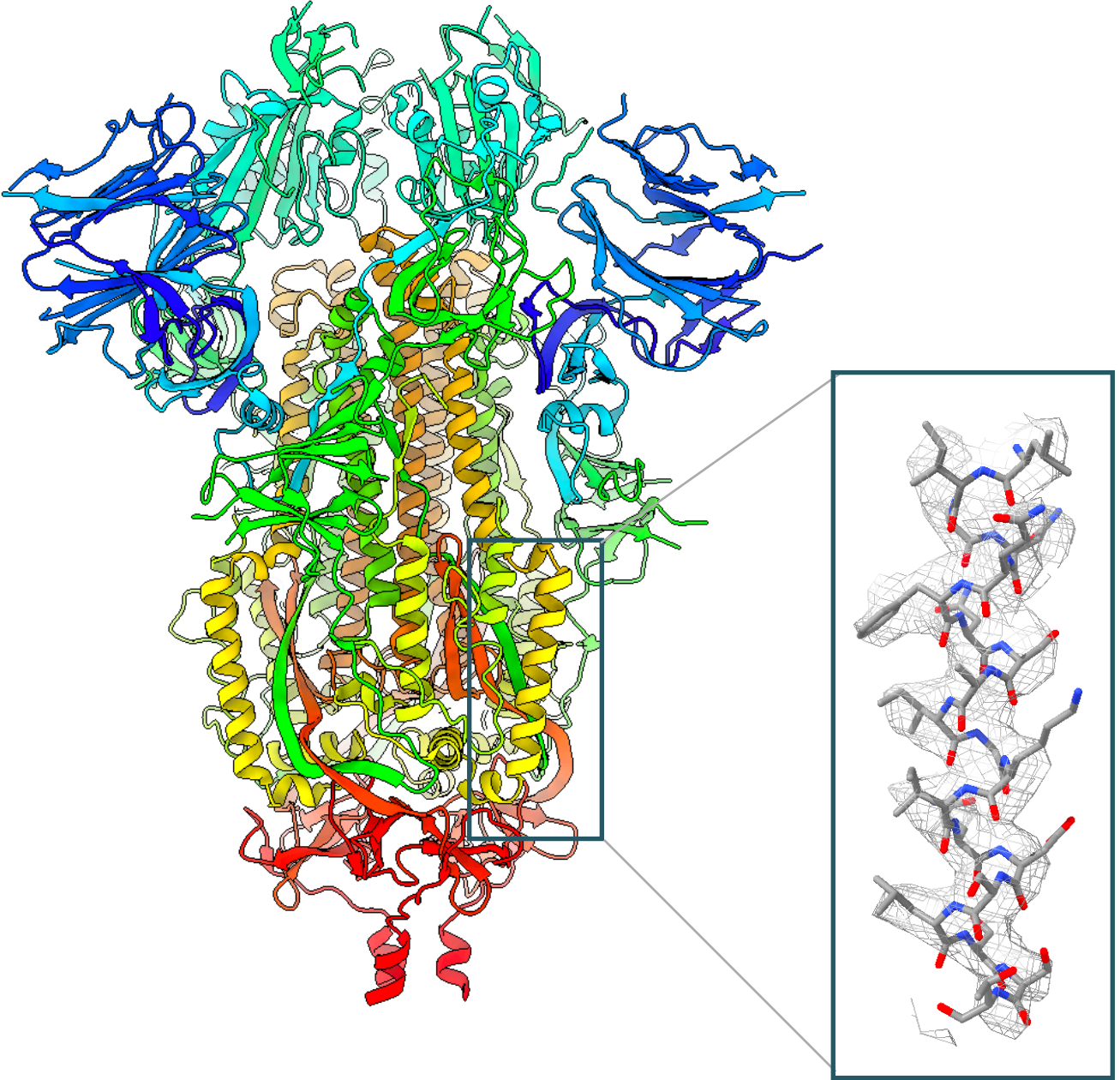
2) 2D to 3D reconstruction

10⁴-10⁷ images



Walls et al, 2020

3) Atomic model fitting

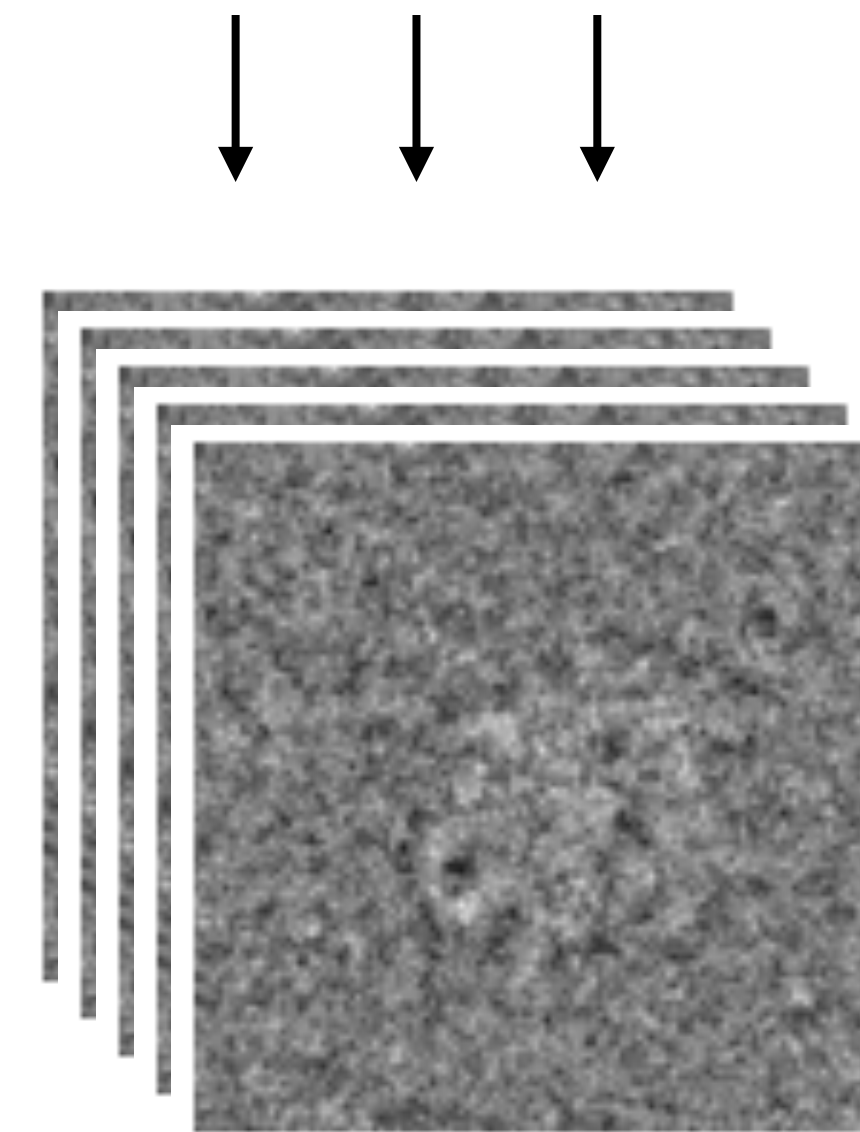
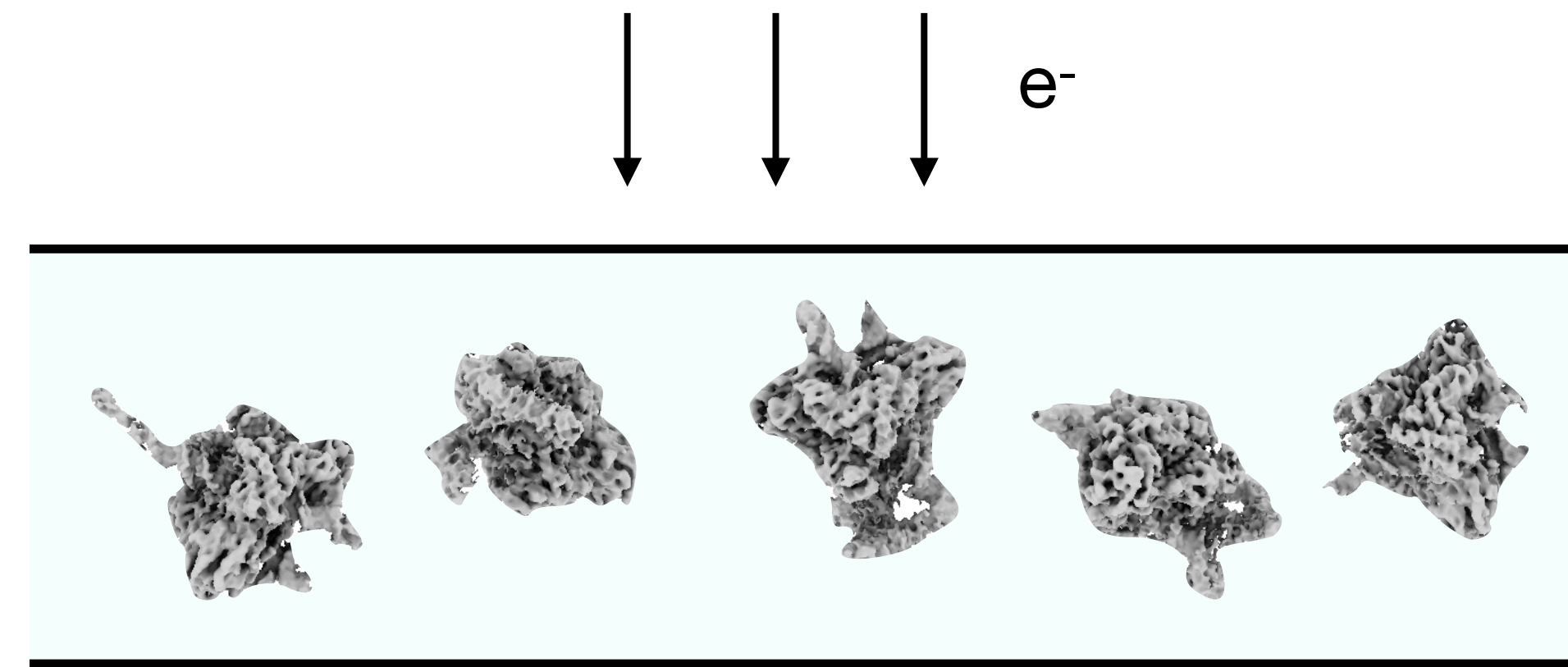


Single particle cryo-EM image formation

- A purified solution of the molecule is fixed in a thin layer of vitreous ice
- Each cryo-EM image $X : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a tomographic projection of a volume $V : \mathbb{R}^3 \rightarrow \mathbb{R}$

$$X(x, y) = PSF * T_t * \int V(R^T(x, y, z)^T) dz + noise$$

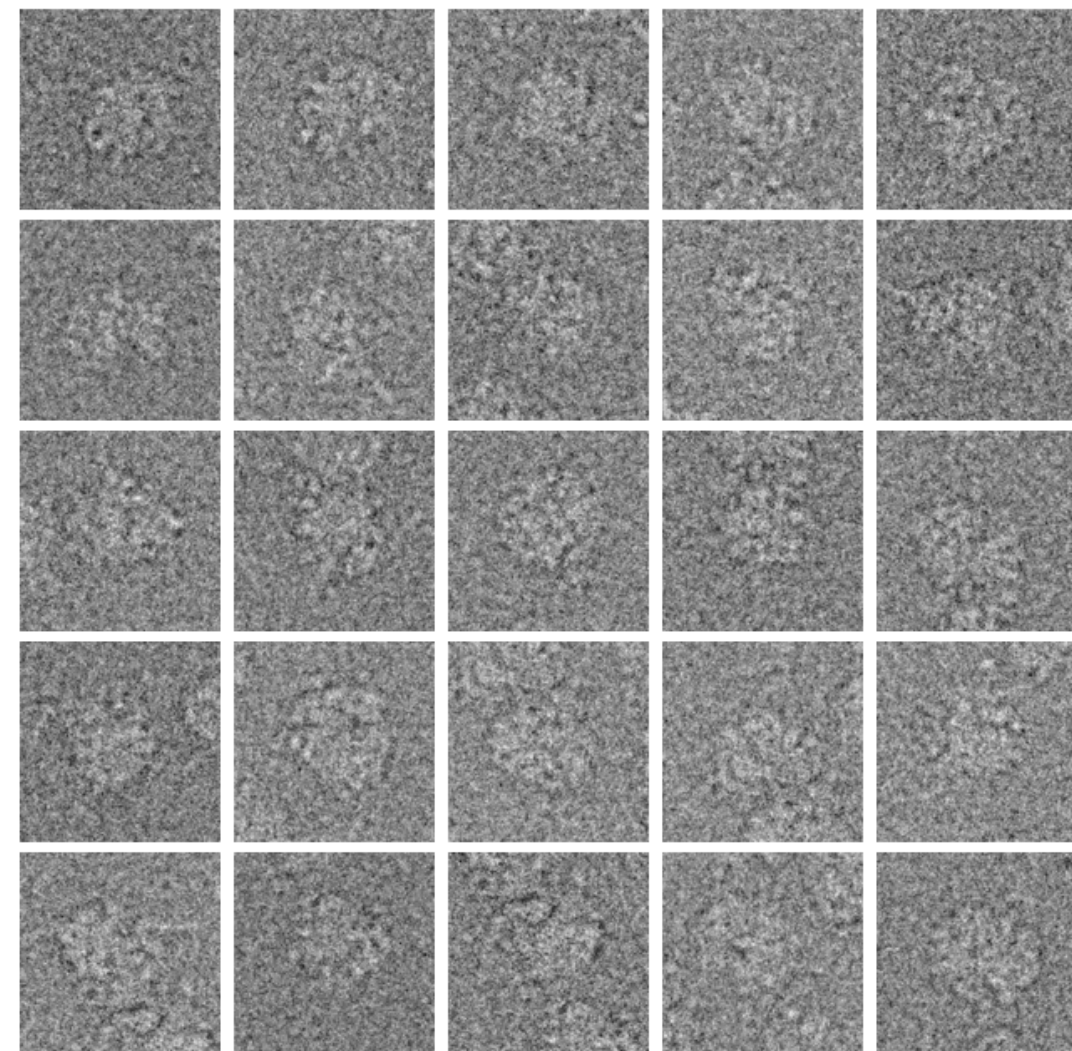
↑ ↑ ↑
 Microscope point spread function In-plane shift by $t \in \mathbb{R}^2$ 3D rotation by $R \in SO(3)$



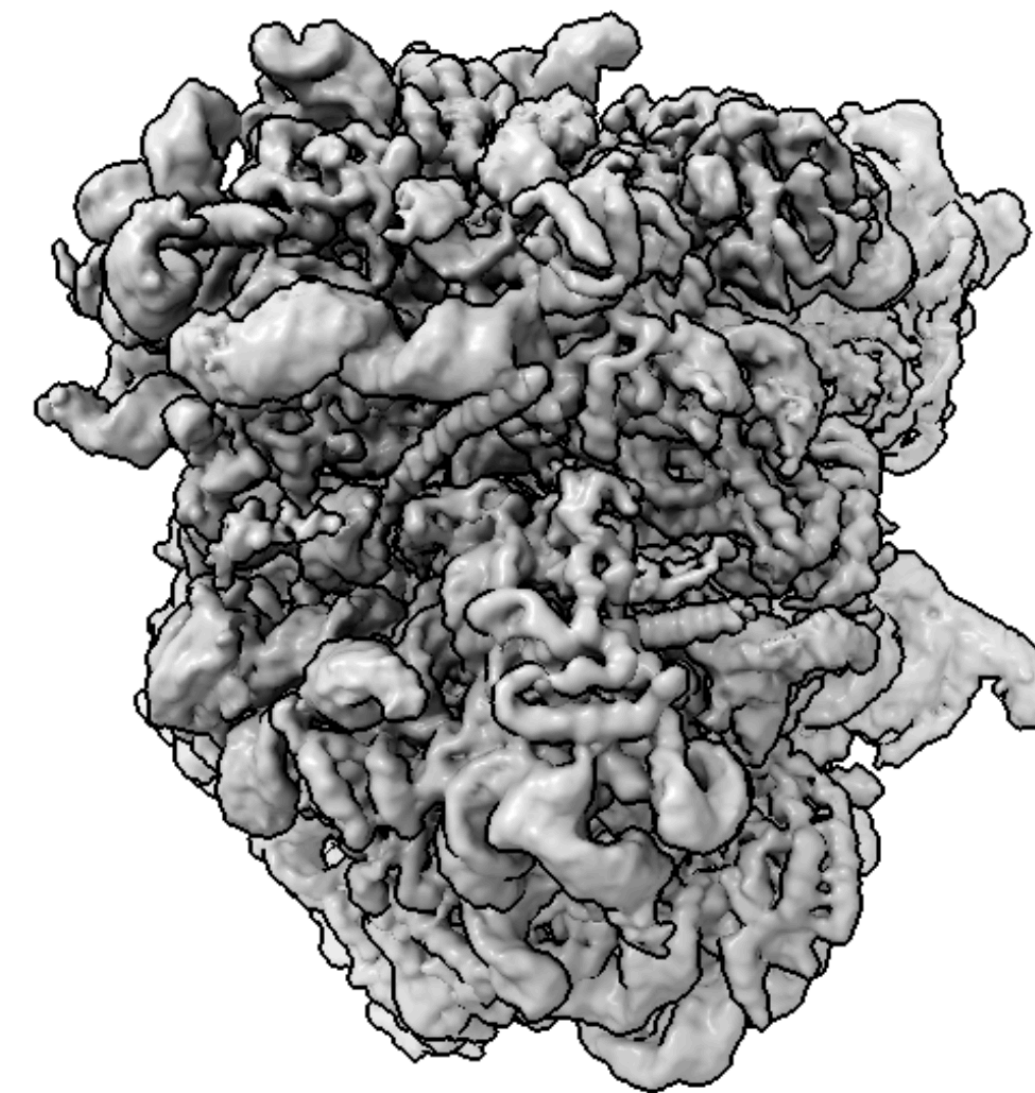
~10⁴⁻⁷ projection images

The cryo-EM reconstruction task

- **Goal:** Reconstruct a volume $V : \mathbb{R}^3 \rightarrow \mathbb{R}$ describing a molecule's 3D structure from a set of noisy projection images X_1, \dots, X_N each containing a copy of V captured from an unknown pose $\phi_i \in (SO(3) \times \mathbb{R}^2)$

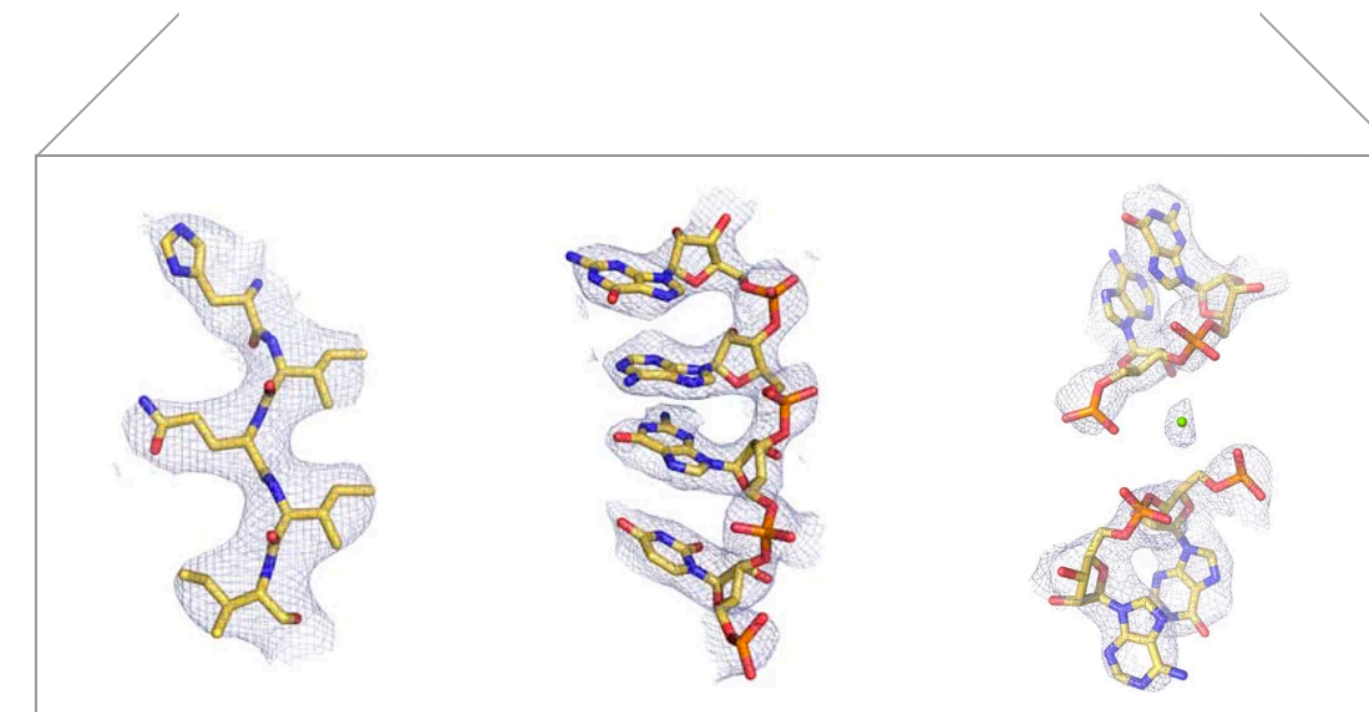


[EMPIAR-10028]



Challenges

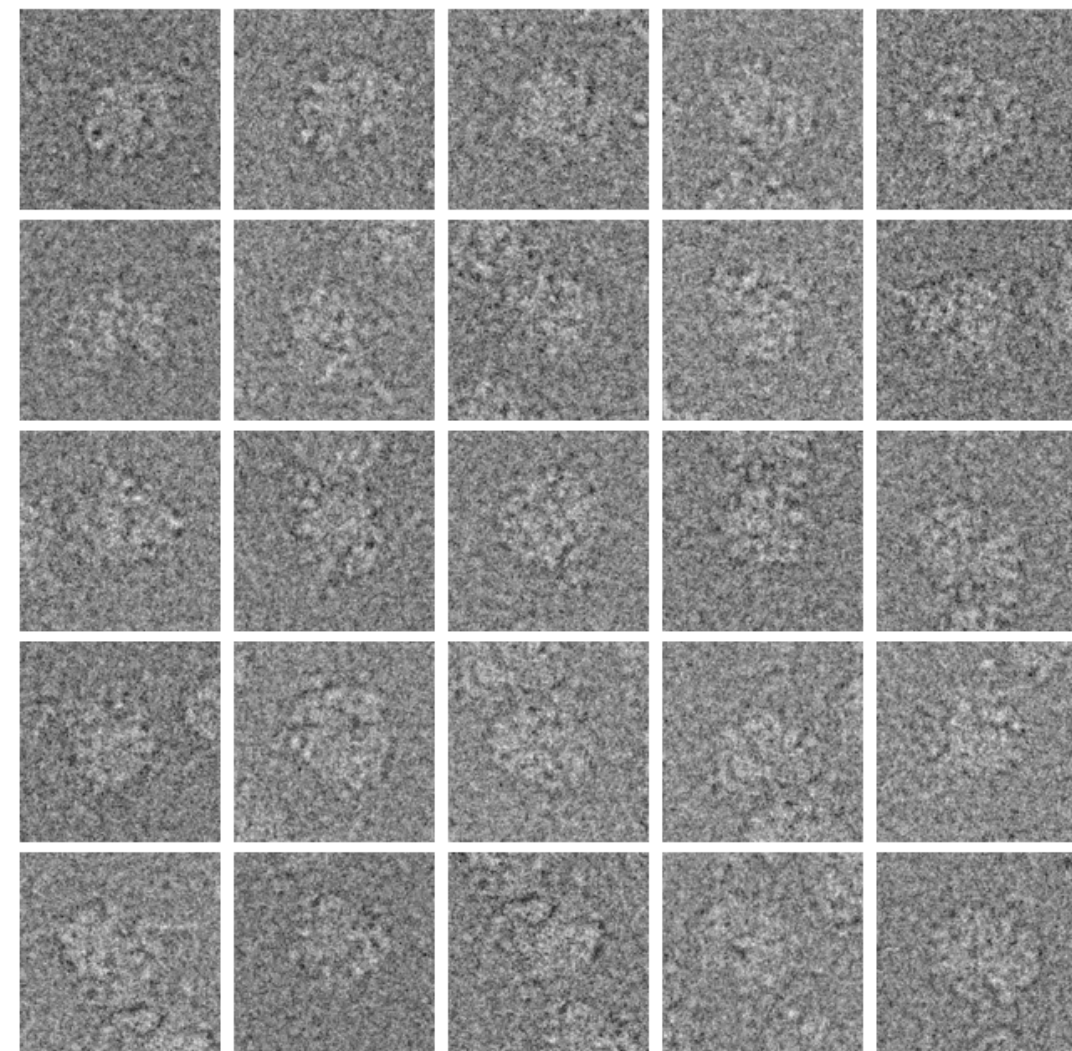
- Unknown particle poses
- Low signal to noise ratio
- Image degrading filters in microscopy
- Discretization of the measurements



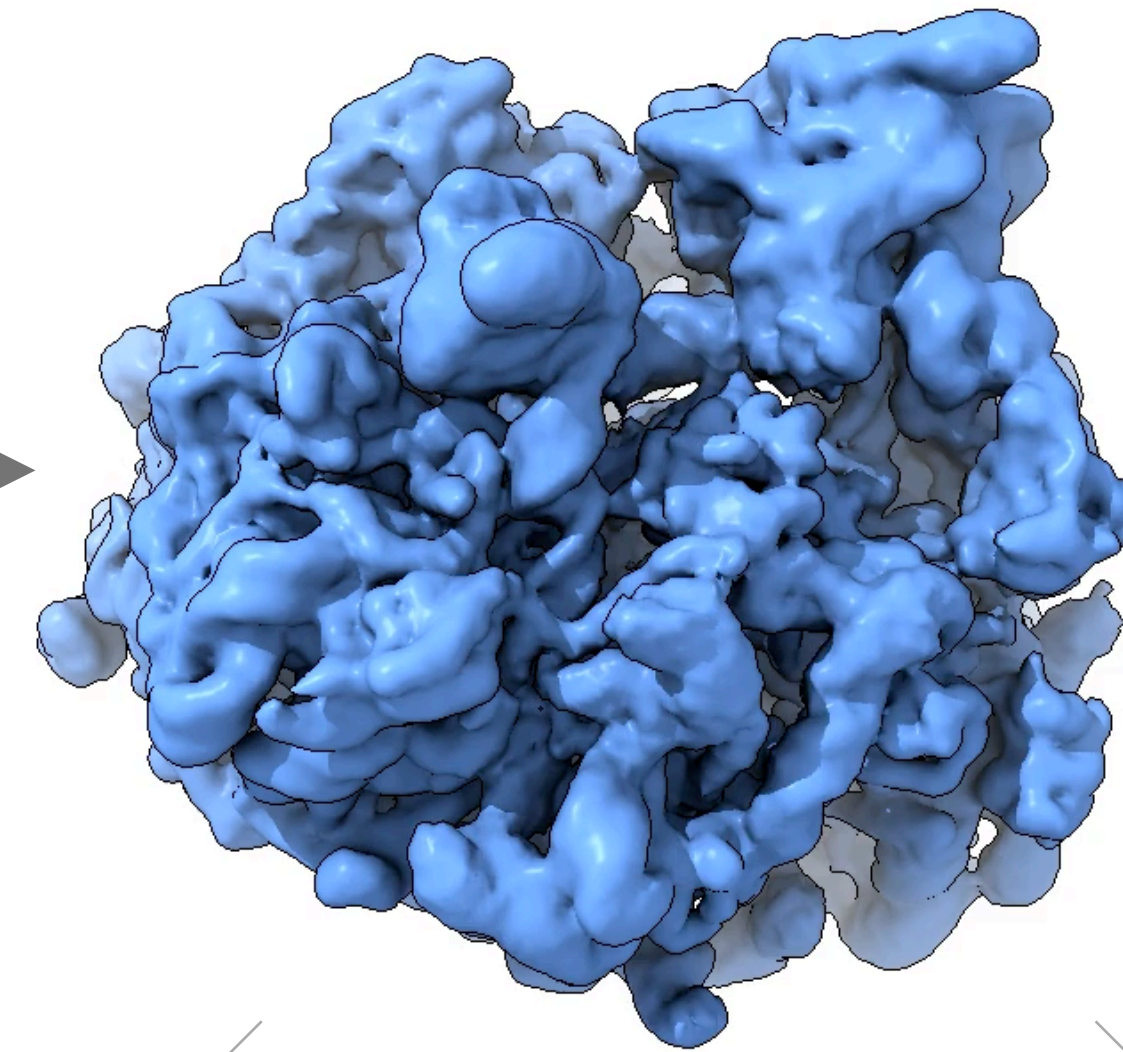
Wong et al. 2014

The cryo-EM reconstruction task

- **Goal:** Reconstruct a volume $V : \mathbb{R}^3 \rightarrow \mathbb{R}$ describing a molecule's 3D structure from a set of noisy projection images X_1, \dots, X_N each containing a copy of V captured from an unknown pose $\phi_i \in (SO(3) \times \mathbb{R}^2)$

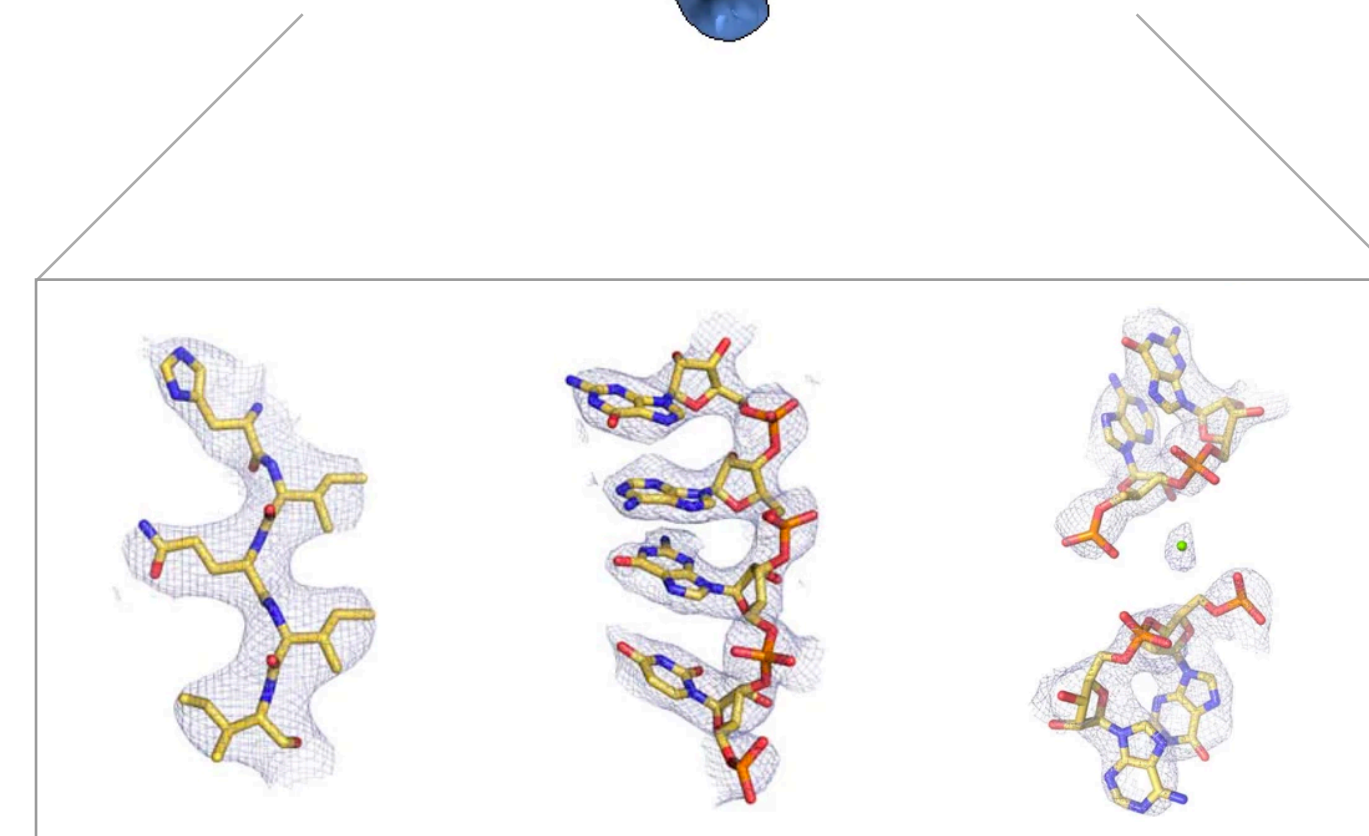


[EMPIAR-10028]



Challenges

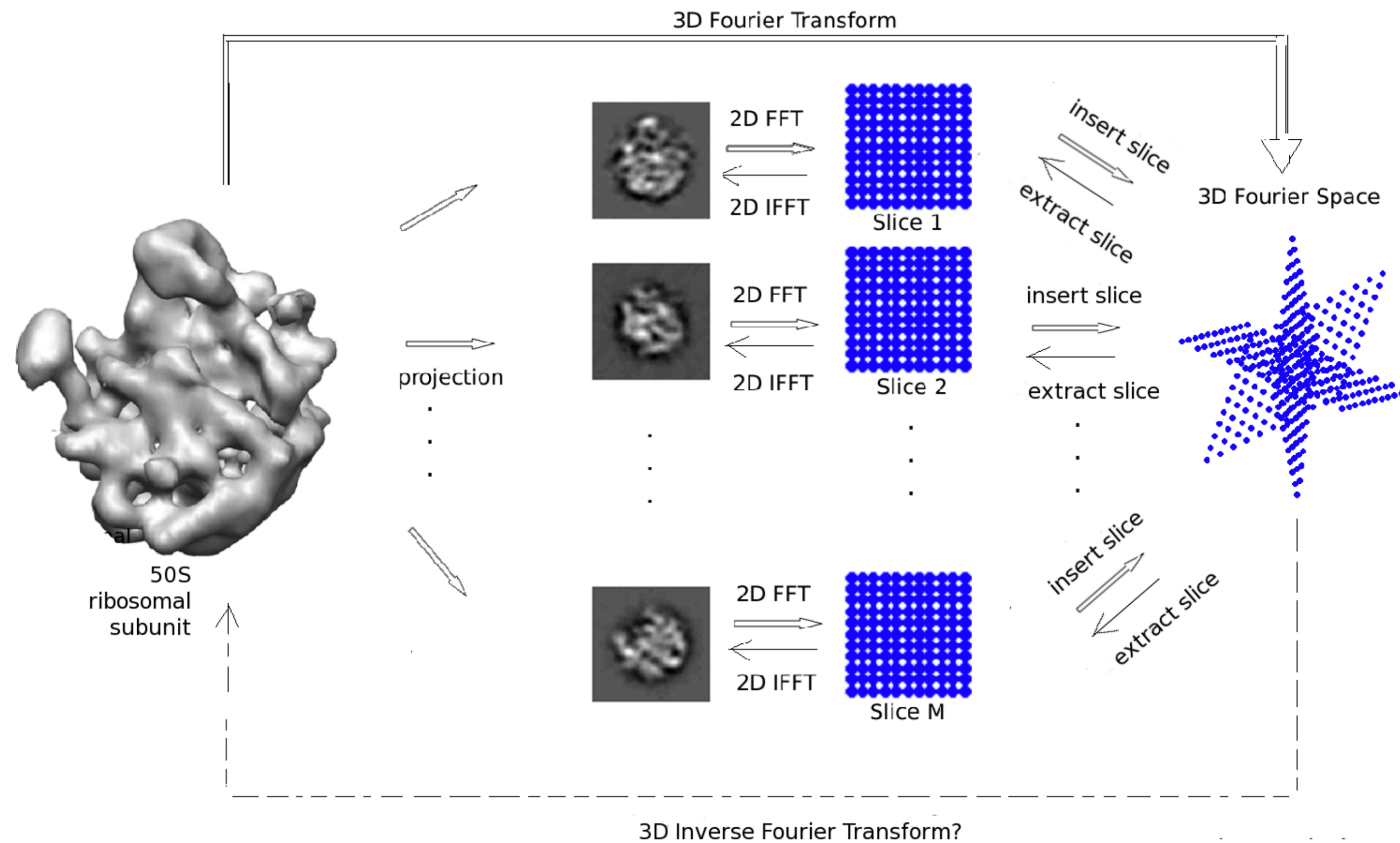
- Unknown particle poses
- Low signal to noise ratio
- Image degrading filters in microscopy
- Discretization of the measurements
- *The heterogeneity problem*



Wong et al. 2014

The Fourier slice theorem

“The Fourier transform of a 2D projection of a volume is a central slice out of the 3D Fourier transform of the volume, perpendicular to the projection direction.”

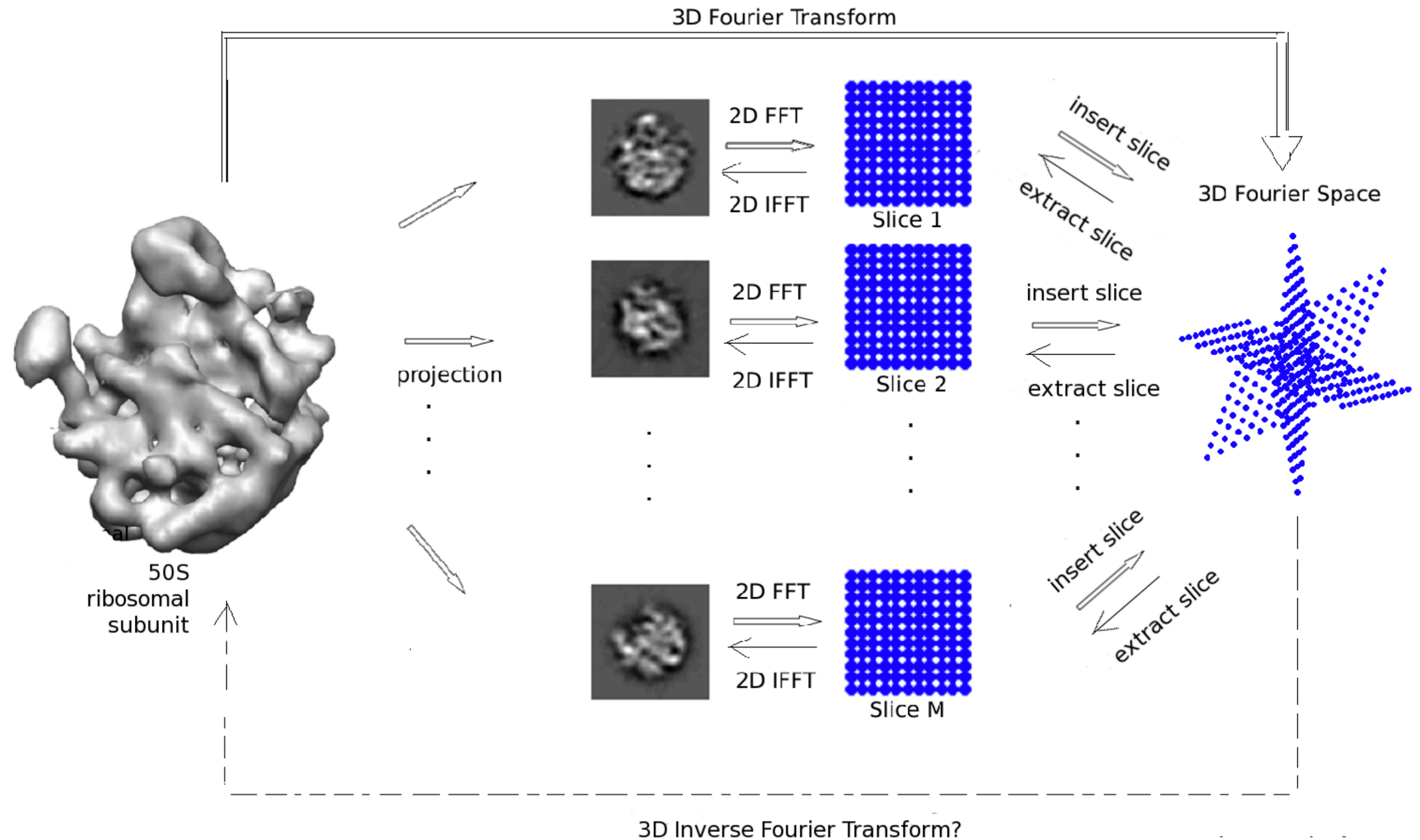


Traditional homogeneous reconstruction algorithms

Goal: Find the 3D structure V_θ that maximizes the likelihood of data $\mathbf{x} = \{x_1, \dots, x_N\}$, marginalizing over unknown poses $\{\phi_i\}$

$$p(\mathbf{x} | V_\theta) = \prod_i \int_{SO(3) \times \mathbb{R}^2} p(x_i, \phi_i | V_\theta) d\phi_i$$

- E-step: Estimate $\{\phi_i\}$ with fixed V_θ
- M-step: Estimate V_θ with fixed $\{\phi_i\}$



Traditional homogeneous reconstruction algorithms

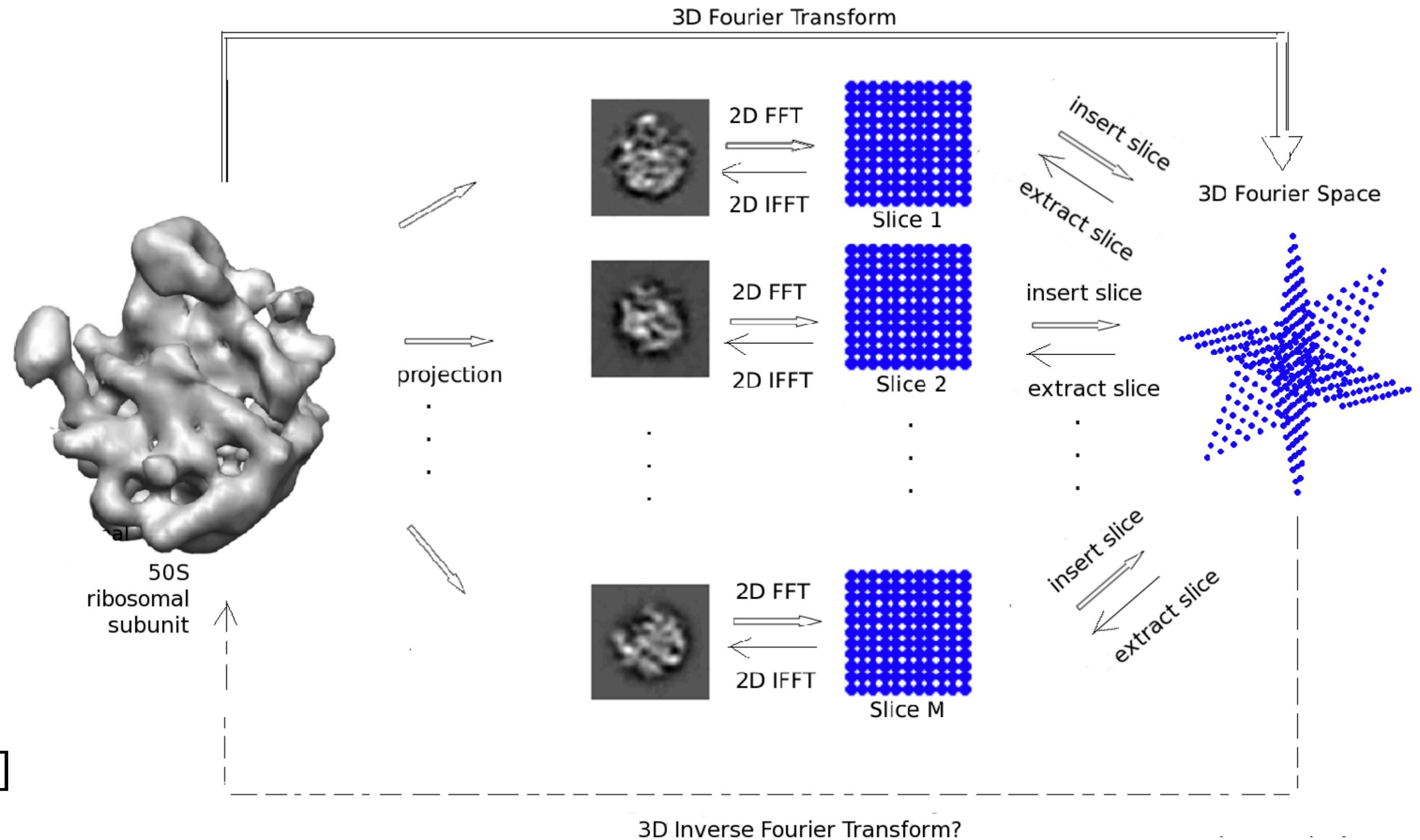
Goal: Find the 3D structure V_θ that maximizes the likelihood of data $\mathbf{x} = \{x_1, \dots, x_N\}$, marginalizing over unknown poses $\{\phi_i\}$

$$p(\mathbf{x} | V_\theta) = \prod_i \int_{SO(3) \times \mathbb{R}^2} p(x_i, \phi_i | V_\theta) d\phi_i$$

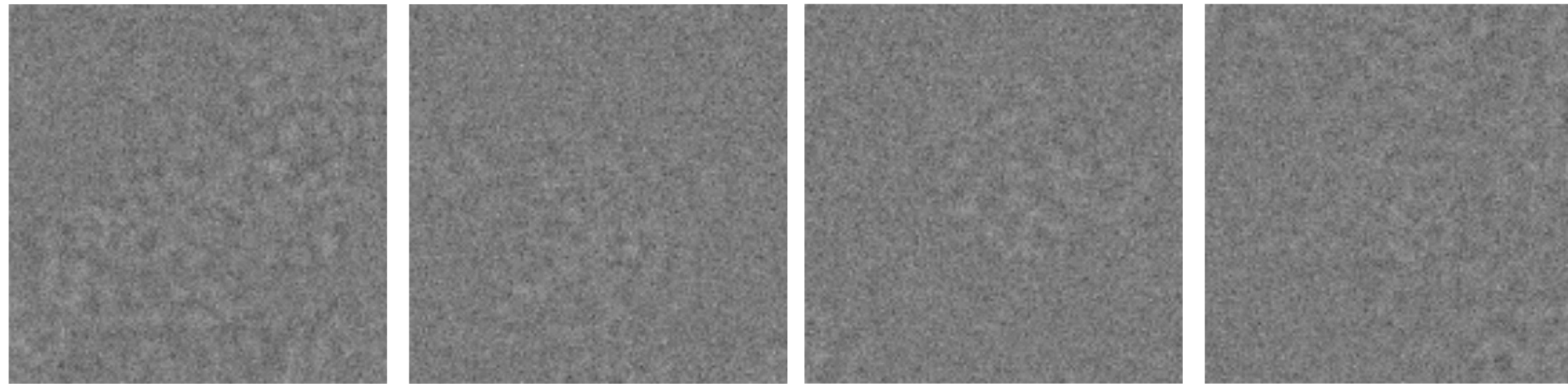
- E-step: Estimate $\{\phi_i\}$ with fixed V_θ
- M-step: Estimate V_θ with fixed $\{\phi_i\}$

Many state of the art software packages:

- RELION: Bayesian formulation for MAP estimation proposed by Sjors Scheres [JSB 2013]
- CryoSPARC: Stochastic optimization techniques proposed by Punjani, Rubinstein, Fleet, Brubaker [CVPR 2016, Nat Methods 2017]



“The heterogeneity problem”

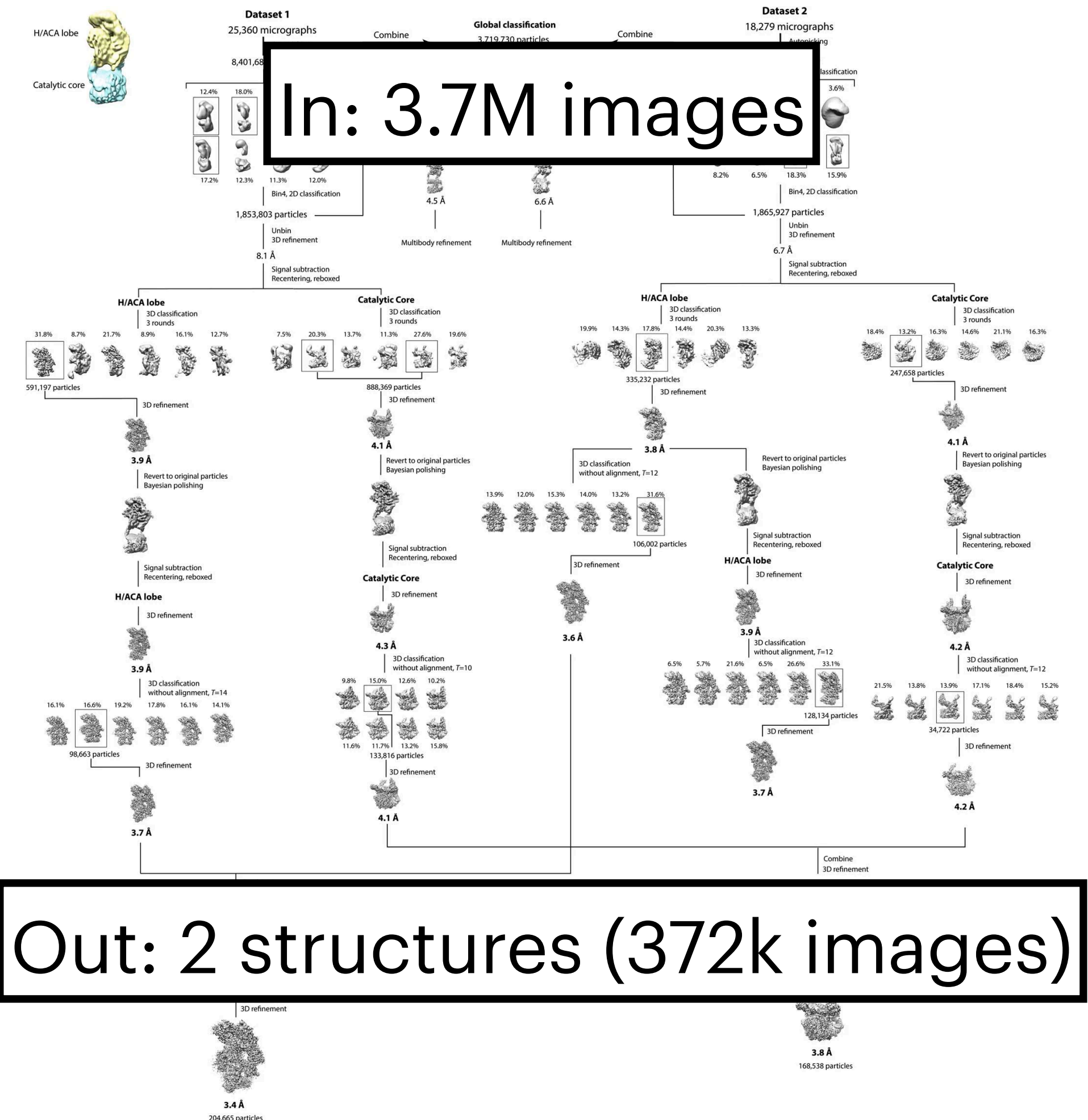


Each image contains a unique molecule

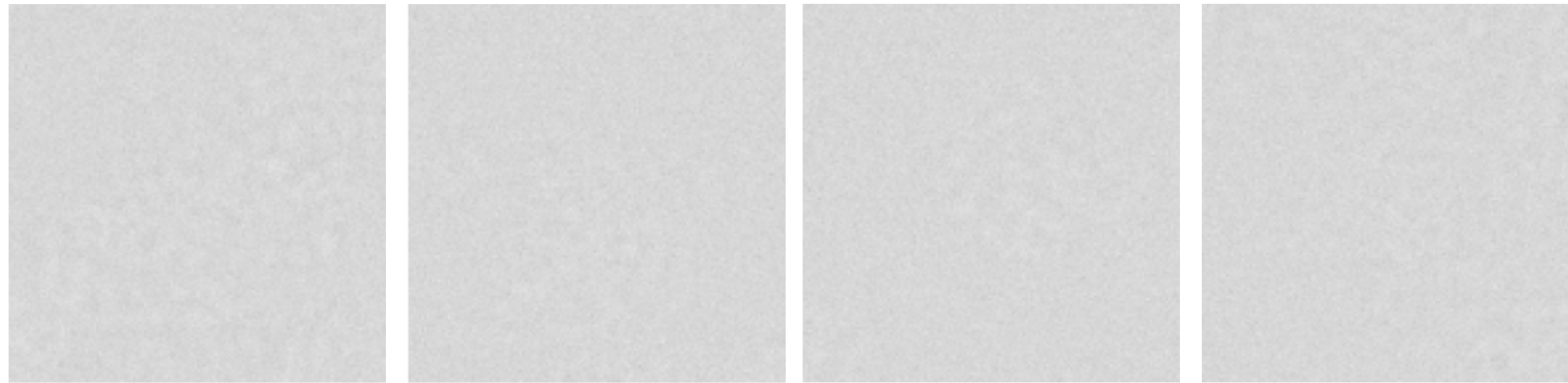
- The ability to image heterogeneous structures is a major opportunity in cryo-EM
- Standard approaches for heterogeneous structures include:
 - Discarding heterogeneous data
 - Multiclass reconstruction — a *discrete* mixture model of K independent structures
- The identification and analysis of heterogeneity — especially continuous forms — is an **open problem** in cryo-EM reconstruction.

Continuous heterogeneity: See Lederman & Singer 2017

A typical cryo-EM processing workflow



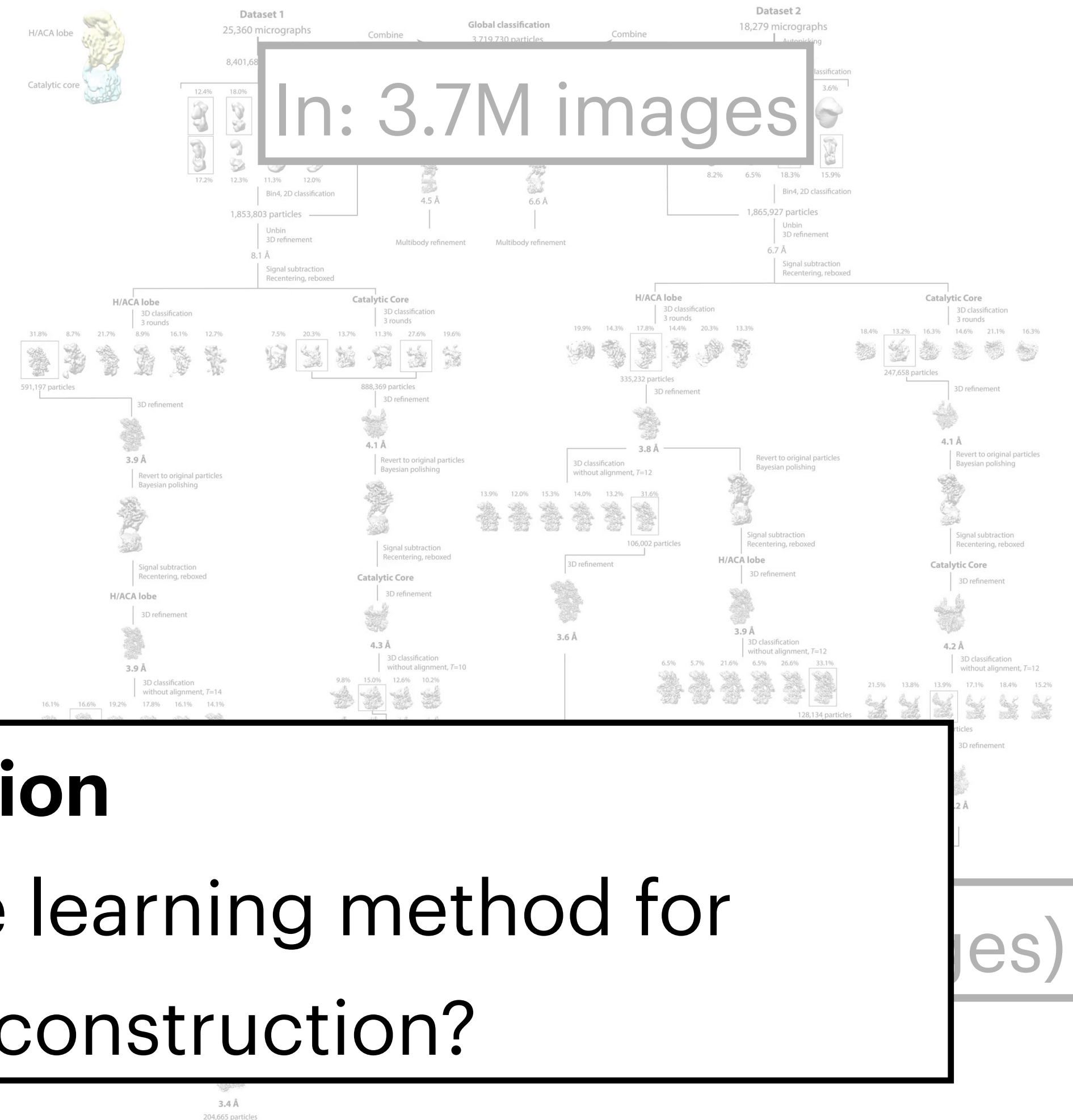
“The heterogeneity problem”



Each image contains a unique molecule

- The ability to image heterogeneous structures is a major opportunity in cryo-EM
- Standard approaches for heterogeneous structures include:
 - Discarding heterogeneous data
 - M
 - in
- The con reconstruction.

A typical cryo-EM processing workflow

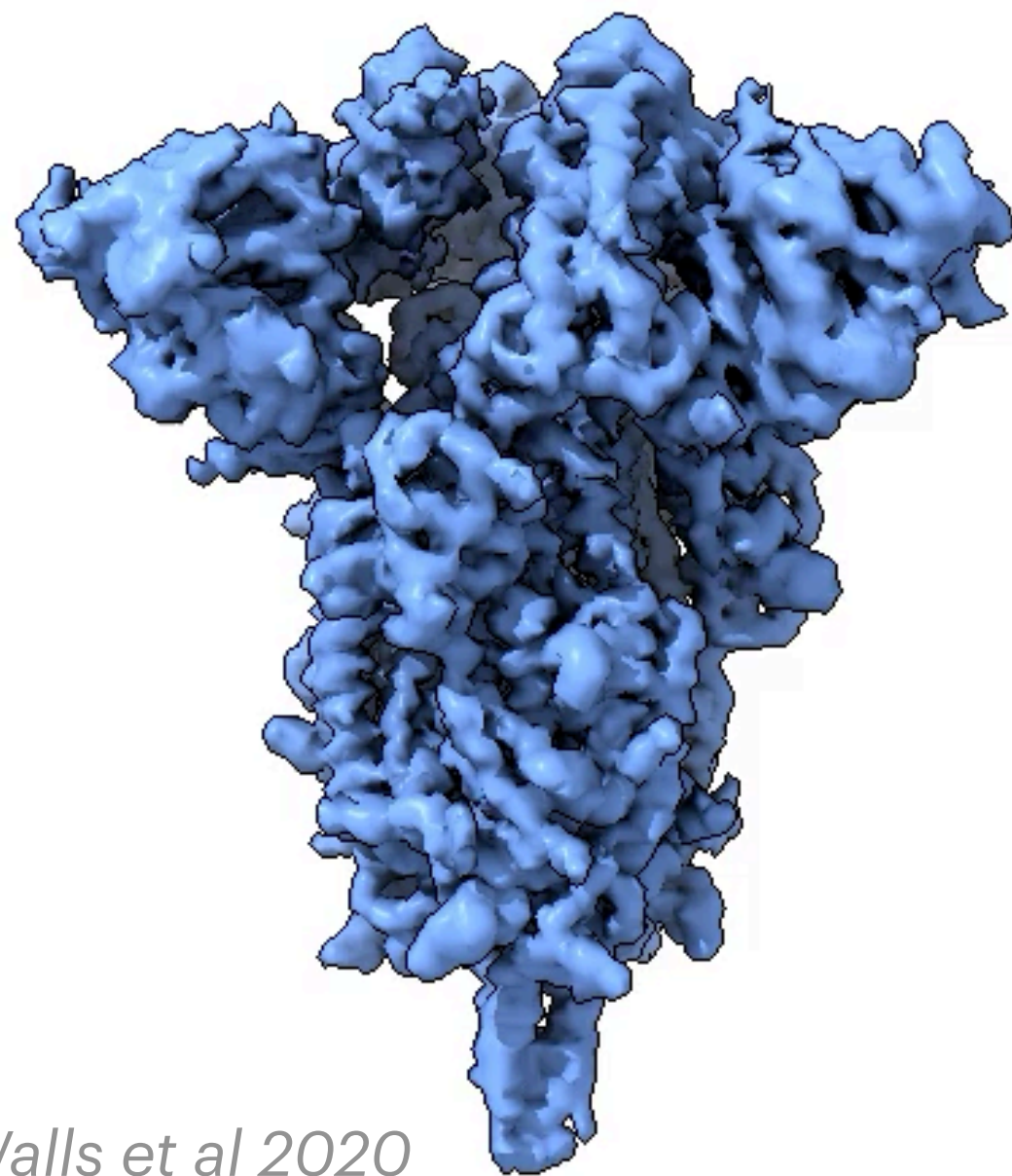
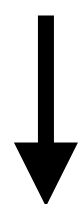
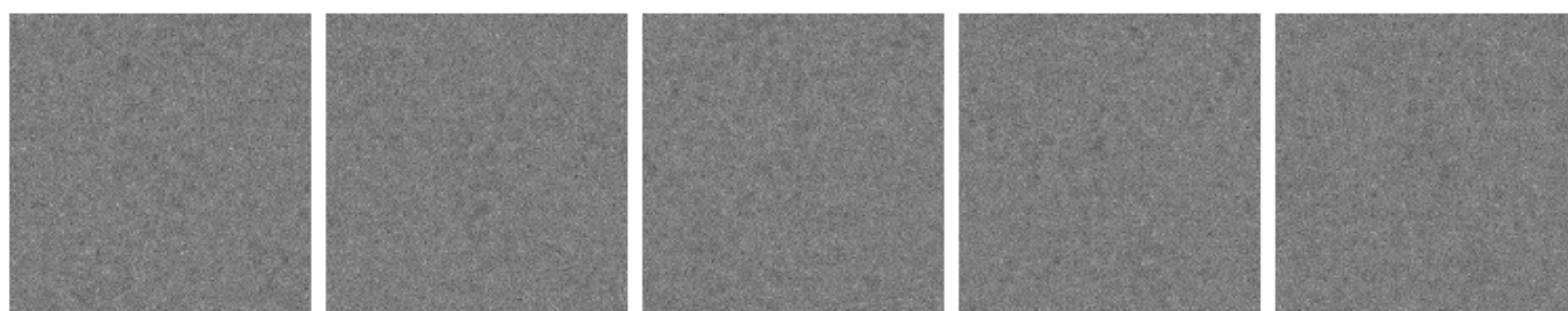


Research Question

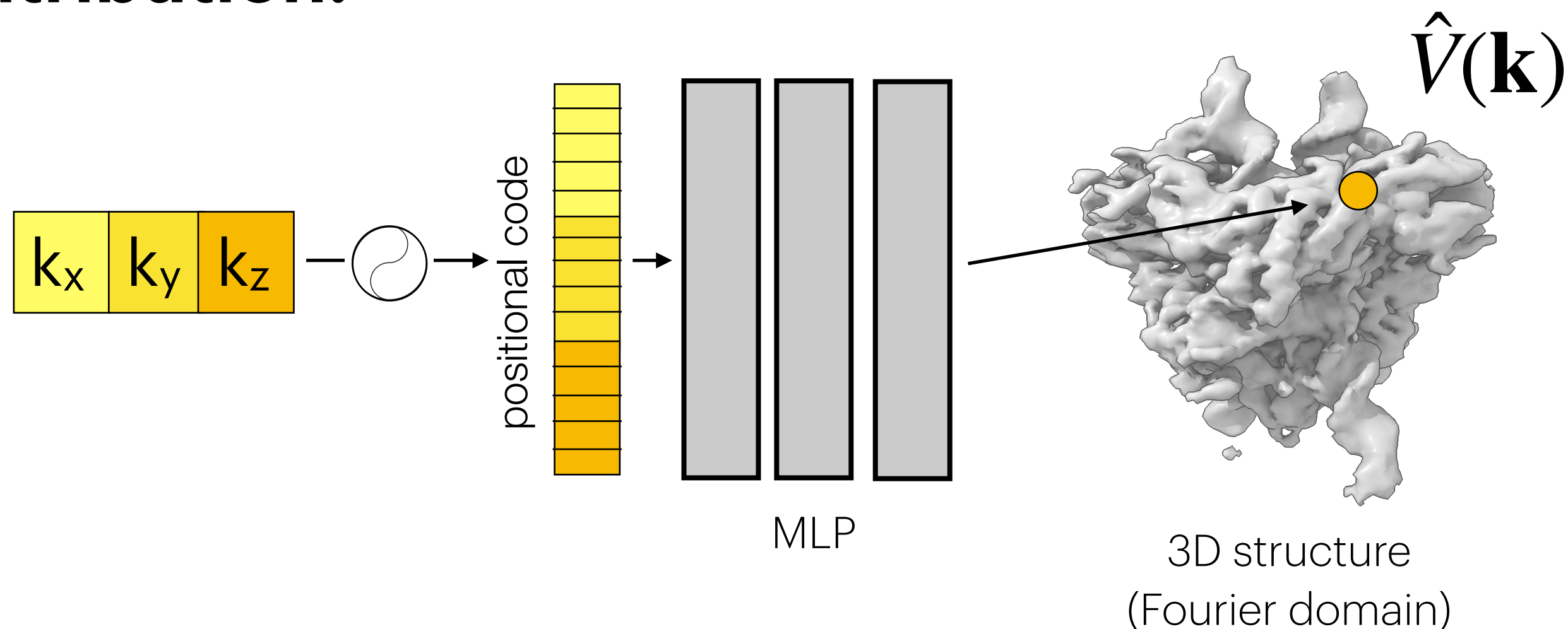
Can we design a modern machine learning method for heterogeneous cryo-EM reconstruction?

CryoDRGN : Deep Reconstructing Generative Networks

Task: 3D reconstruction from unlabeled 2D images



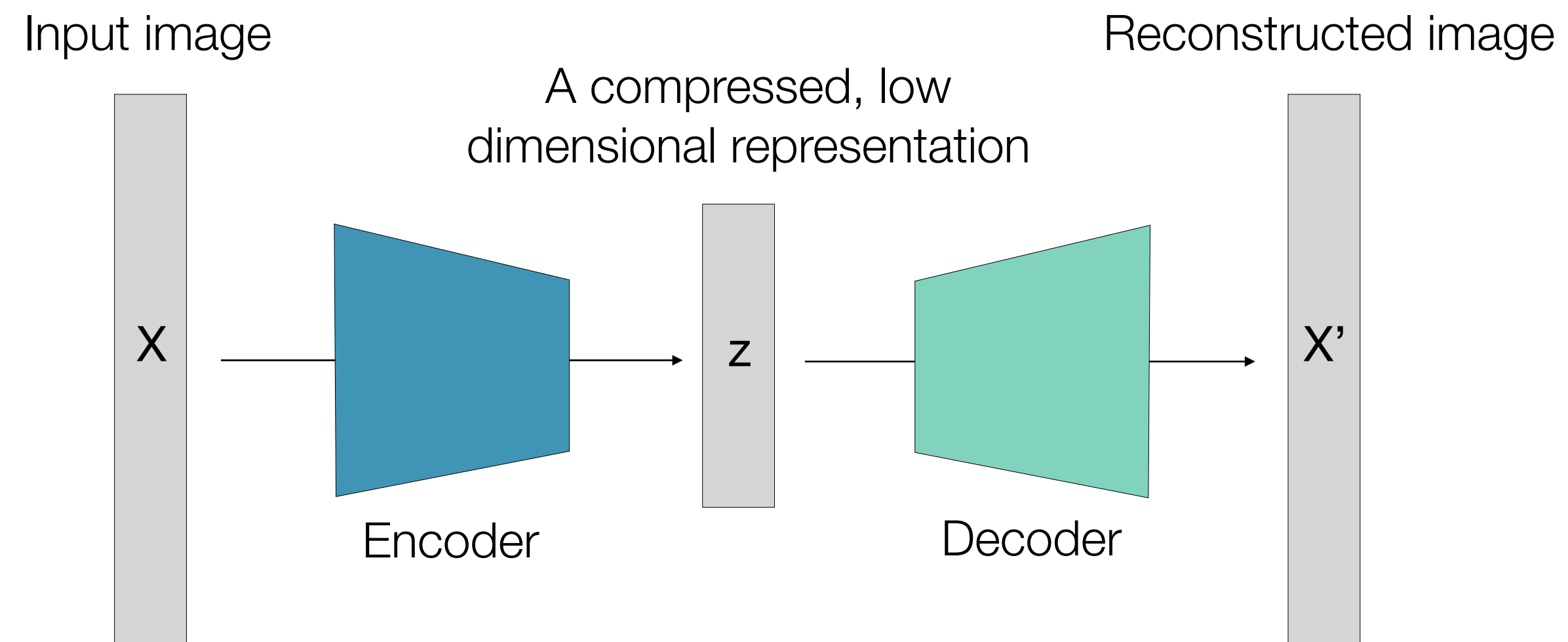
Contribution:



- A new paradigm for heterogeneous cryo-EM reconstruction based on *deep generative models*
- Addresses a major open problem in the field of reconstructing *continuous heterogeneity*
- Introduced a *neural field* representation of 3D structure that has shown broad applicability in computer vision (e.g. NeRF)

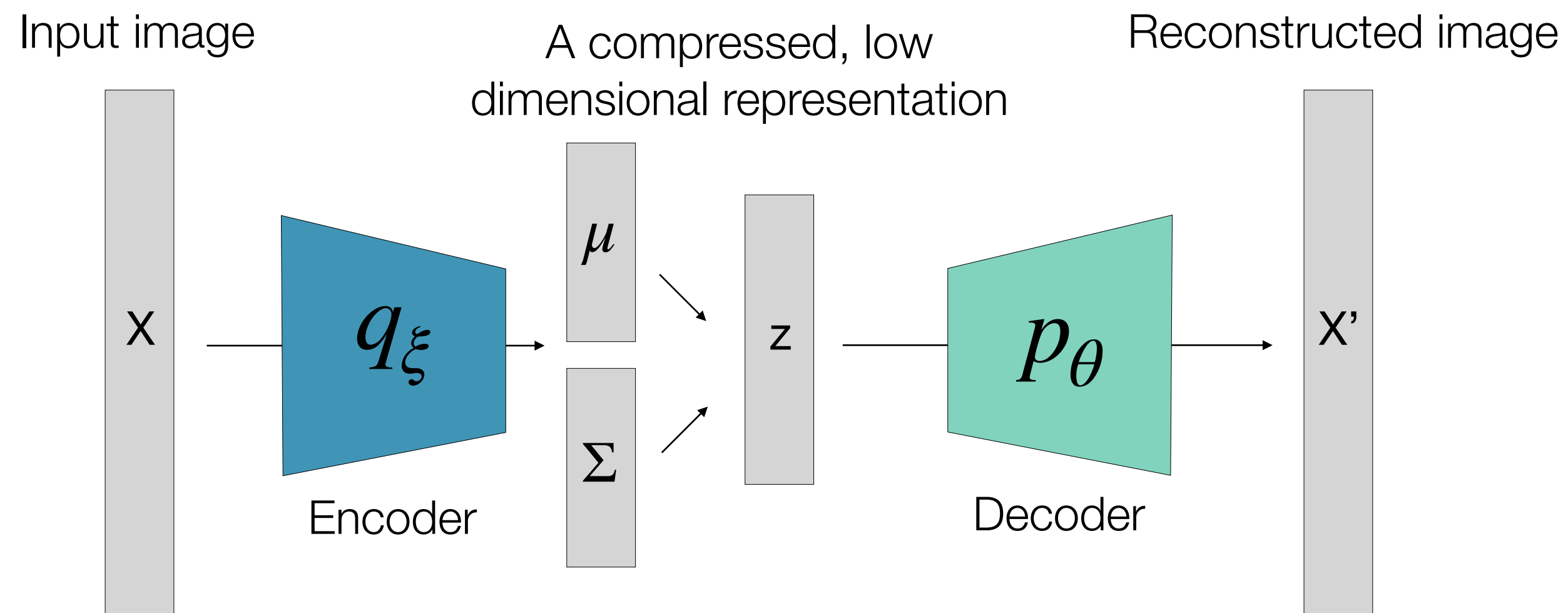
Autoencoders and Variational Autoencoders (VAEs)

- The autoencoder is a nonlinear, dimensionality reduction technique



Autoencoders and Variational Autoencoders (VAEs)

- The autoencoder is a nonlinear, dimensionality reduction technique



- The VAE extends the AE as inference of a probabilistic model – “a regularized autoencoder”

$$\mathcal{L}_{VAE}(X; \theta, \xi) = \underbrace{\mathbb{E}_{q_{\xi}(z|X)}[\log p_{\theta}(X|z)]}_{\text{Reconstruction error}} - \underbrace{KL(q_{\xi}(z|X) || p(z))}_{\text{Regularization}}$$

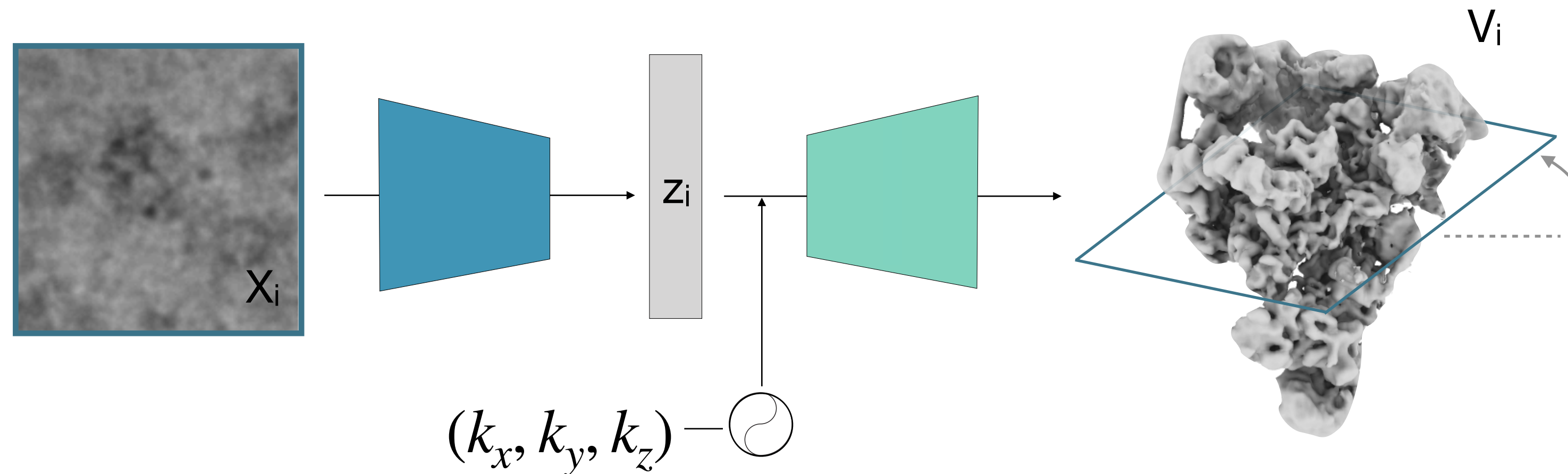
Reconstruction error

Regularization

CryoDRGN ❄️🐉: Deep reconstructing generative networks

Unsupervised learning of a **deep generative model** of 3D biomolecular structures from 2D cryo-EM images

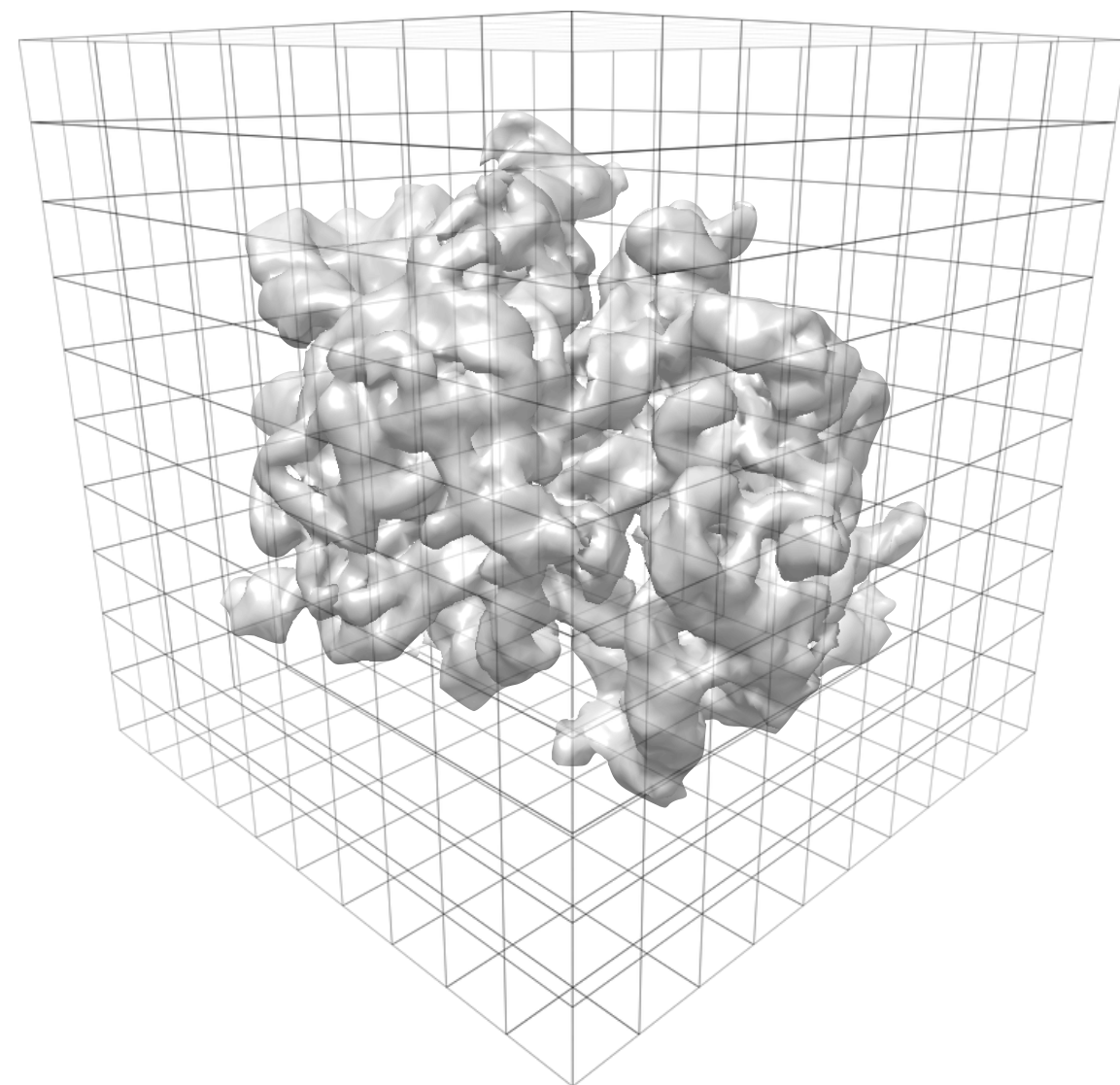
1. We develop **coordinate-based neural networks** to directly approximate the 3D structure
2. Fourier space **image encoder-volume decoder** architecture based on the variational autoencoder (VAE)
3. Exact inference for pose and variational inference for heterogeneity



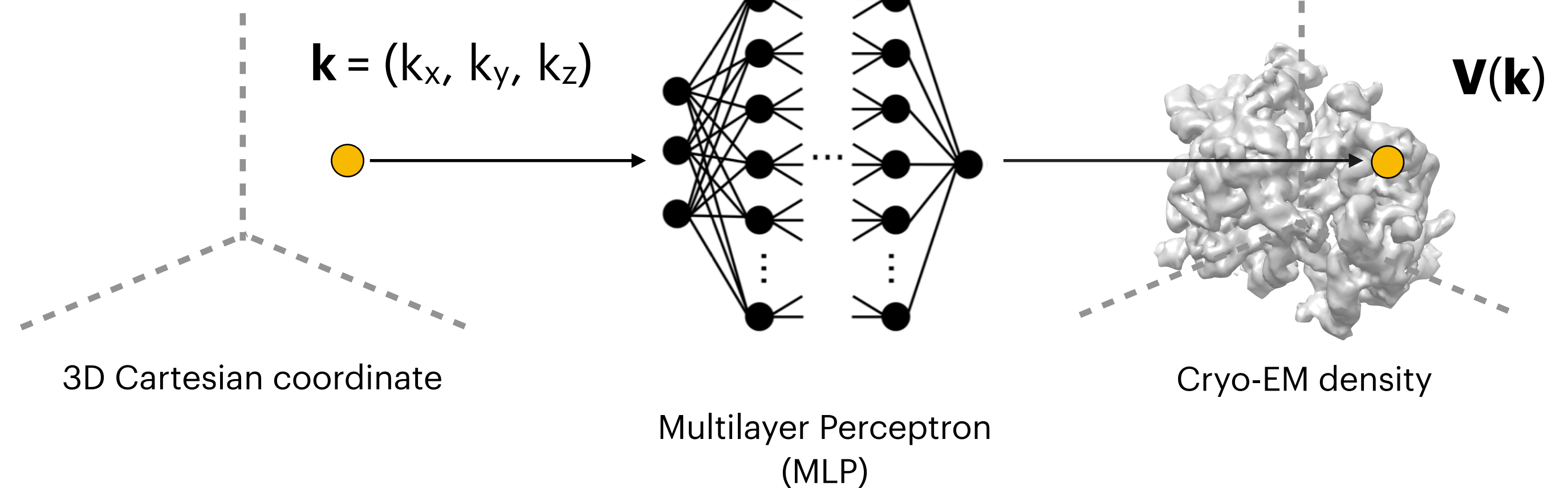
Coordinate-based neural networks for 3D volumes

- **Key idea:** Instead of representing the structure as discrete points on a 3D lattice, learn a *continuous* function, $V : \mathbb{R}^3 \rightarrow \mathbb{R}$

Traditional algorithms



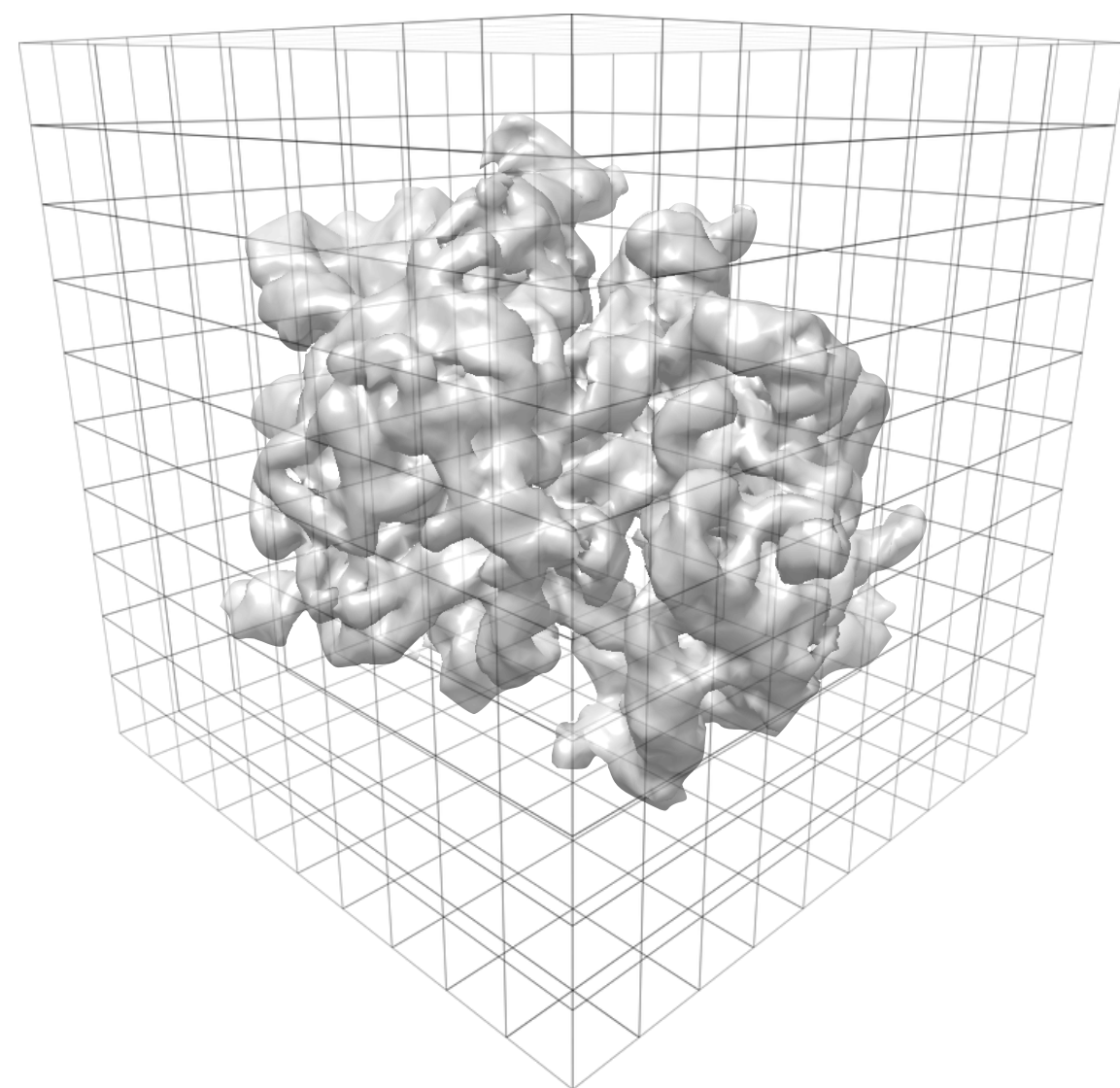
cryoDRGN



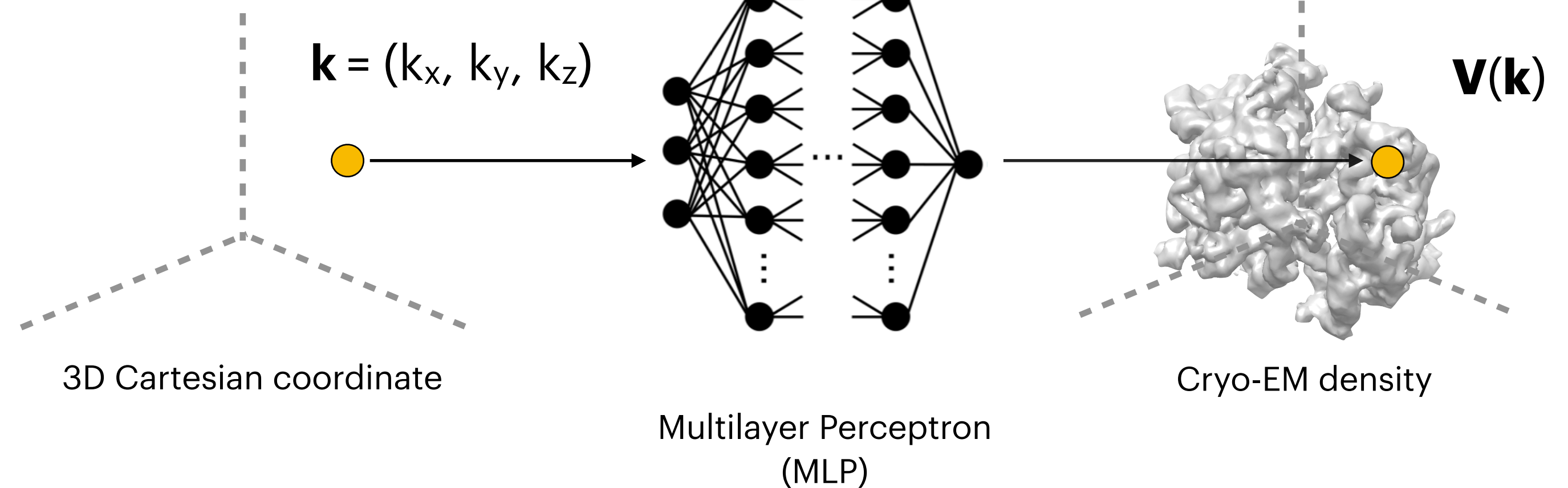
Coordinate-based neural networks for 3D volumes

- **Key idea:** Instead of representing the structure as discrete points on a 3D lattice, learn a *continuous* function, $V : \mathbb{R}^3 \rightarrow \mathbb{R}$

Traditional algorithms



cryoDRGN

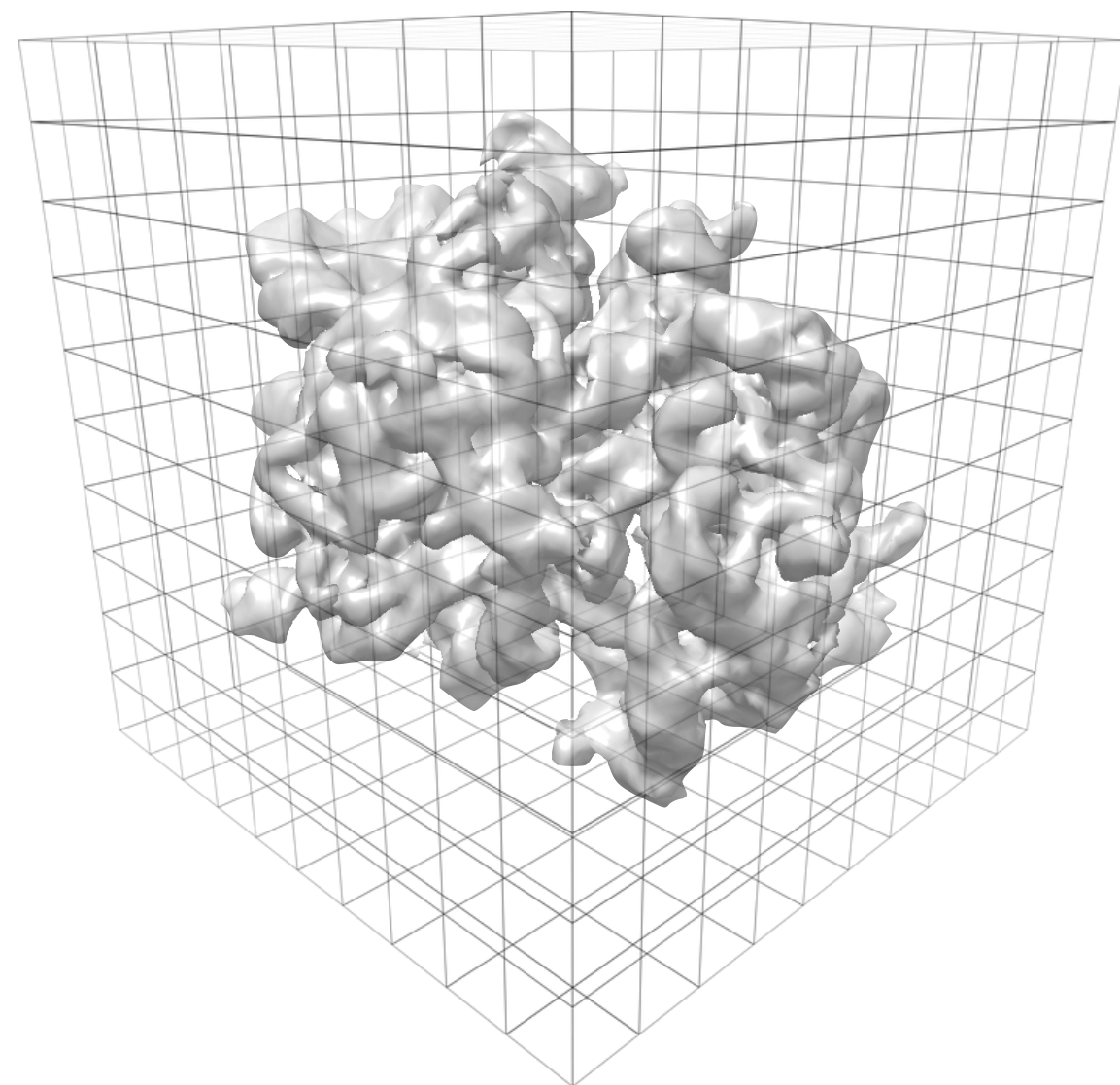


CryoDRGN structures are parameterized as a **neural network** instead of a voxel array

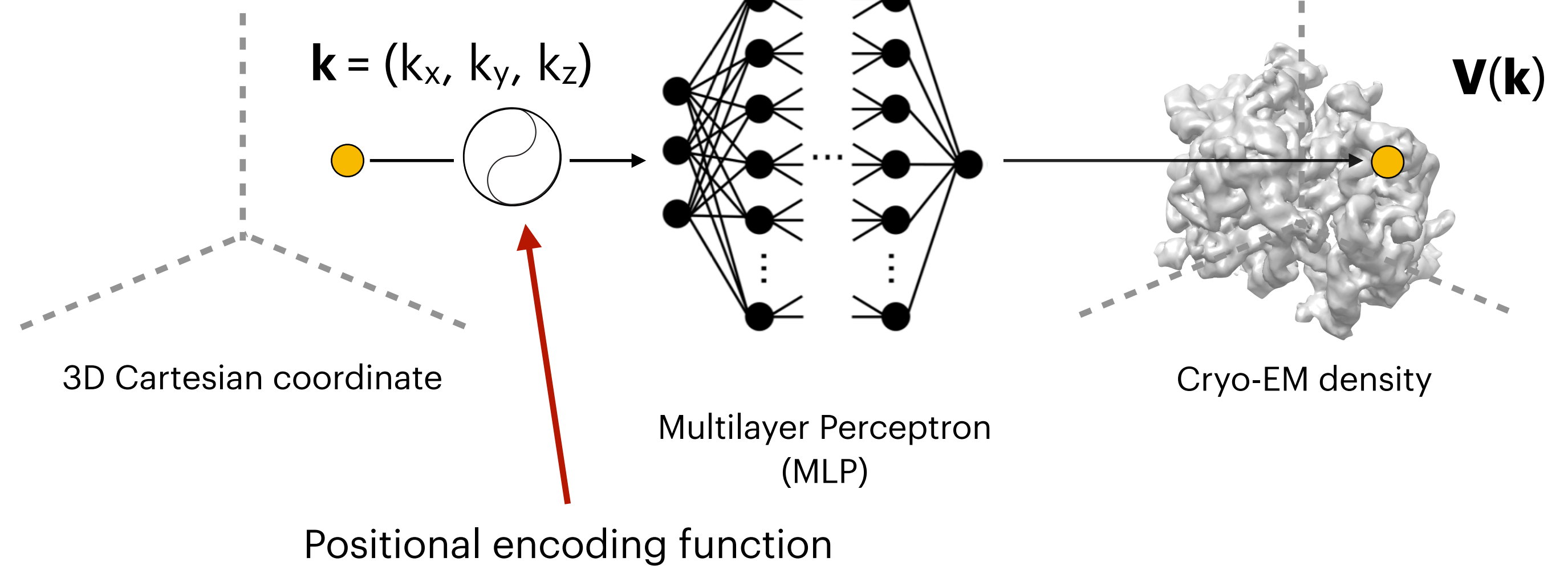
Coordinate-based neural networks for 3D volumes

- **Key idea:** Instead of representing the structure as discrete points on a 3D lattice, learn a *continuous* function, $V : \mathbb{R}^3 \rightarrow \mathbb{R}$

Traditional algorithms



cryoDRGN

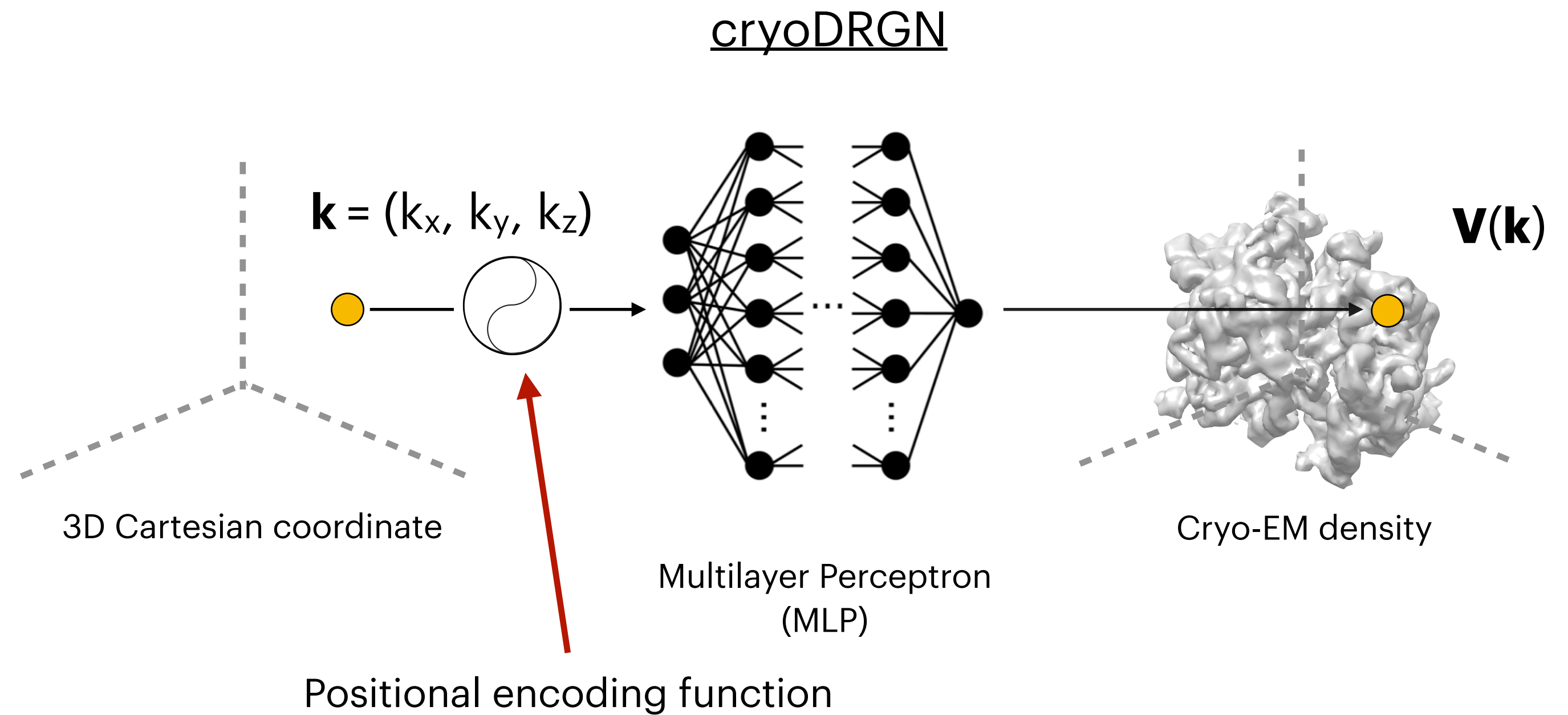
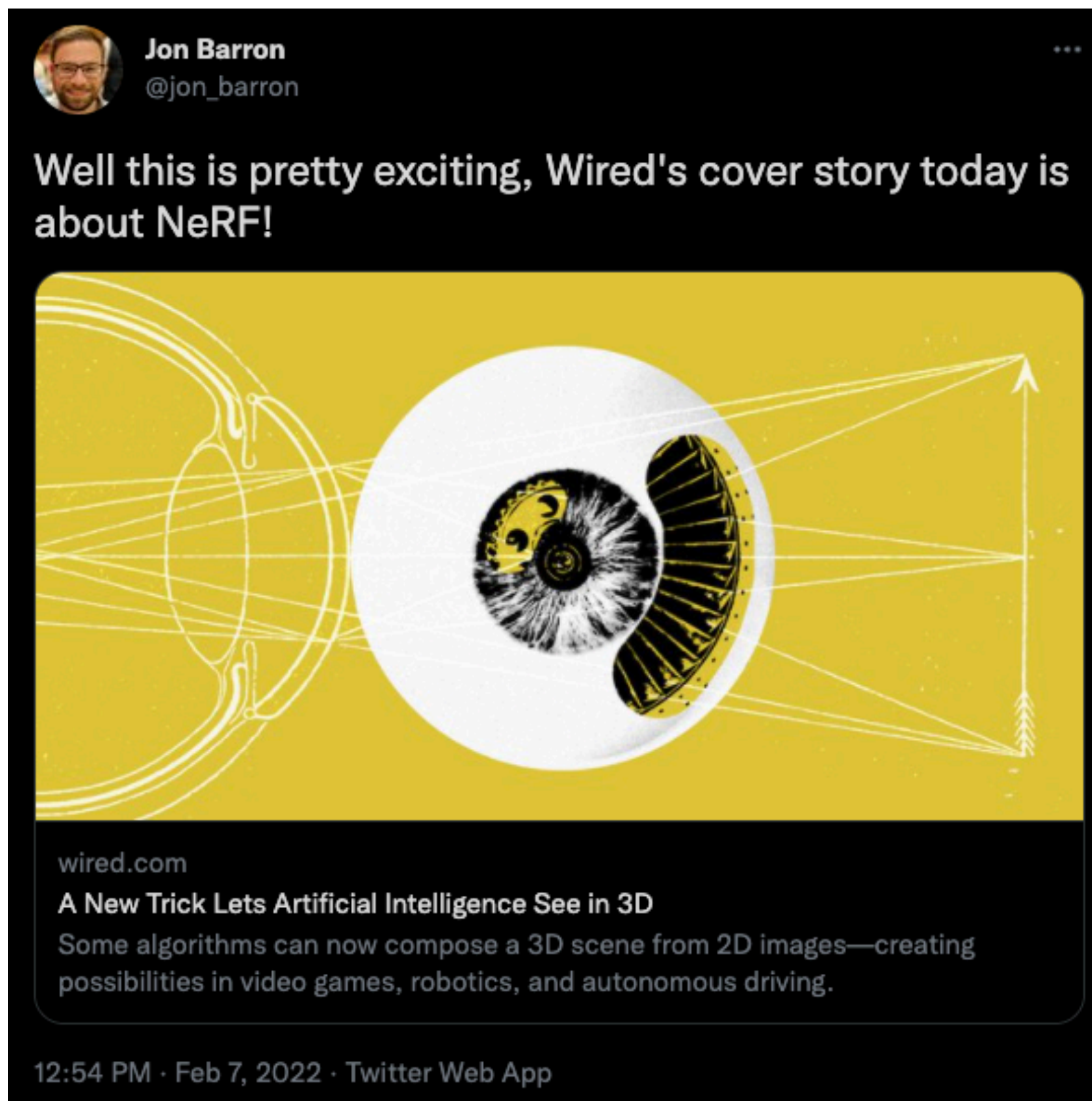


$$pe^{(2i)}(k_j) = \sin(k_j D \pi (2/D)^{2i/D}), \quad i = 1, \dots, D/2; k_j \in k$$

$$pe^{(2i+1)}(k_j) = \cos(k_j D \pi (2/D)^{2i/D}), \quad i = 1, \dots, D/2; k_j \in k$$

Coordinate-based neural networks for 3D volumes

- **Key idea:** Instead of representing the structure as discrete points on a 3D lattice, learn a *continuous* function, $V : \mathbb{R}^3 \rightarrow \mathbb{R}$



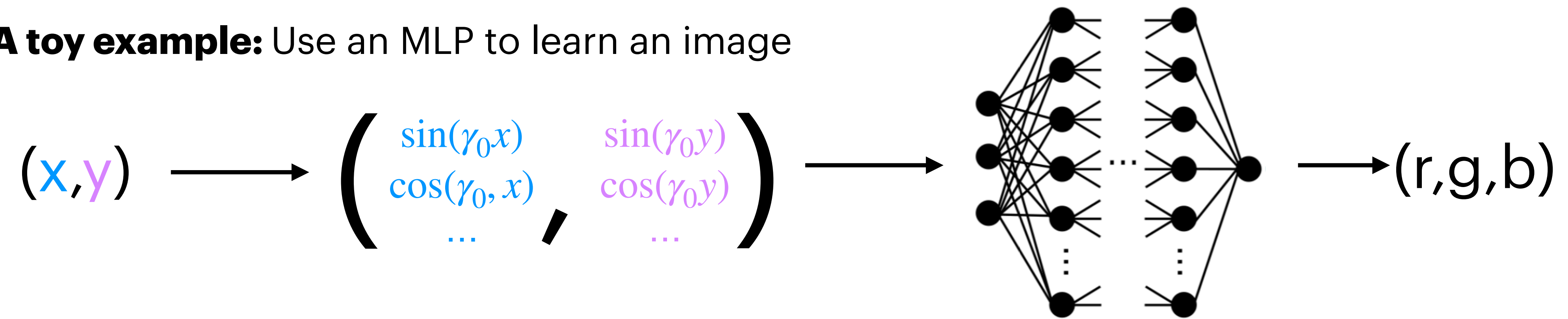
$$pe^{(2i)}(k_j) = \sin(k_j D \pi (2/D)^{2i/D}), \quad i = 1, \dots, D/2; k_j \in k$$

$$pe^{(2i+1)}(k_j) = \cos(k_j D \pi (2/D)^{2i/D}), \quad i = 1, \dots, D/2; k_j \in k$$

Also see:
NeRF, Mildenhall et al. ECCV 2020

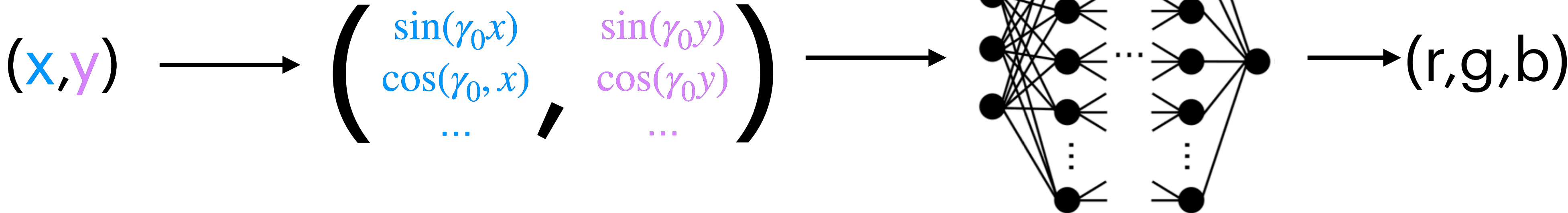
A sinusoidal encoding to featurize input coordinates

A toy example: Use an MLP to learn an image



A sinusoidal encoding to featurize input coordinates

A toy example: Use an MLP to learn an image

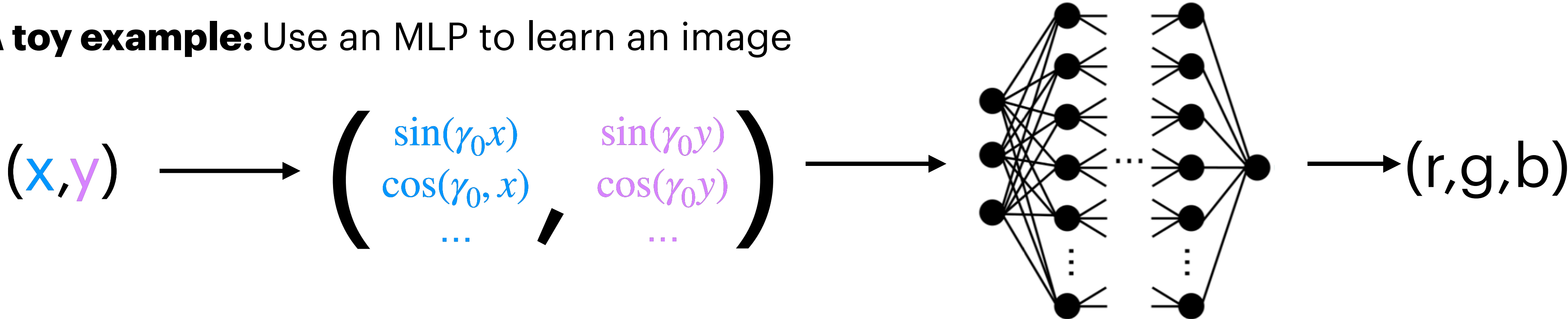


Ground truth



A sinusoidal encoding to featurize input coordinates

A toy example: Use an MLP to learn an image



Ground truth

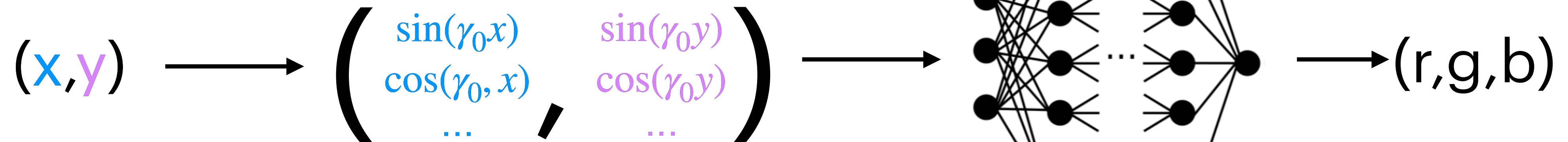


No sinusoidal encoding



A sinusoidal encoding to featurize input coordinates

A toy example: Use an MLP to learn an image



Ground truth



No sinusoidal encoding



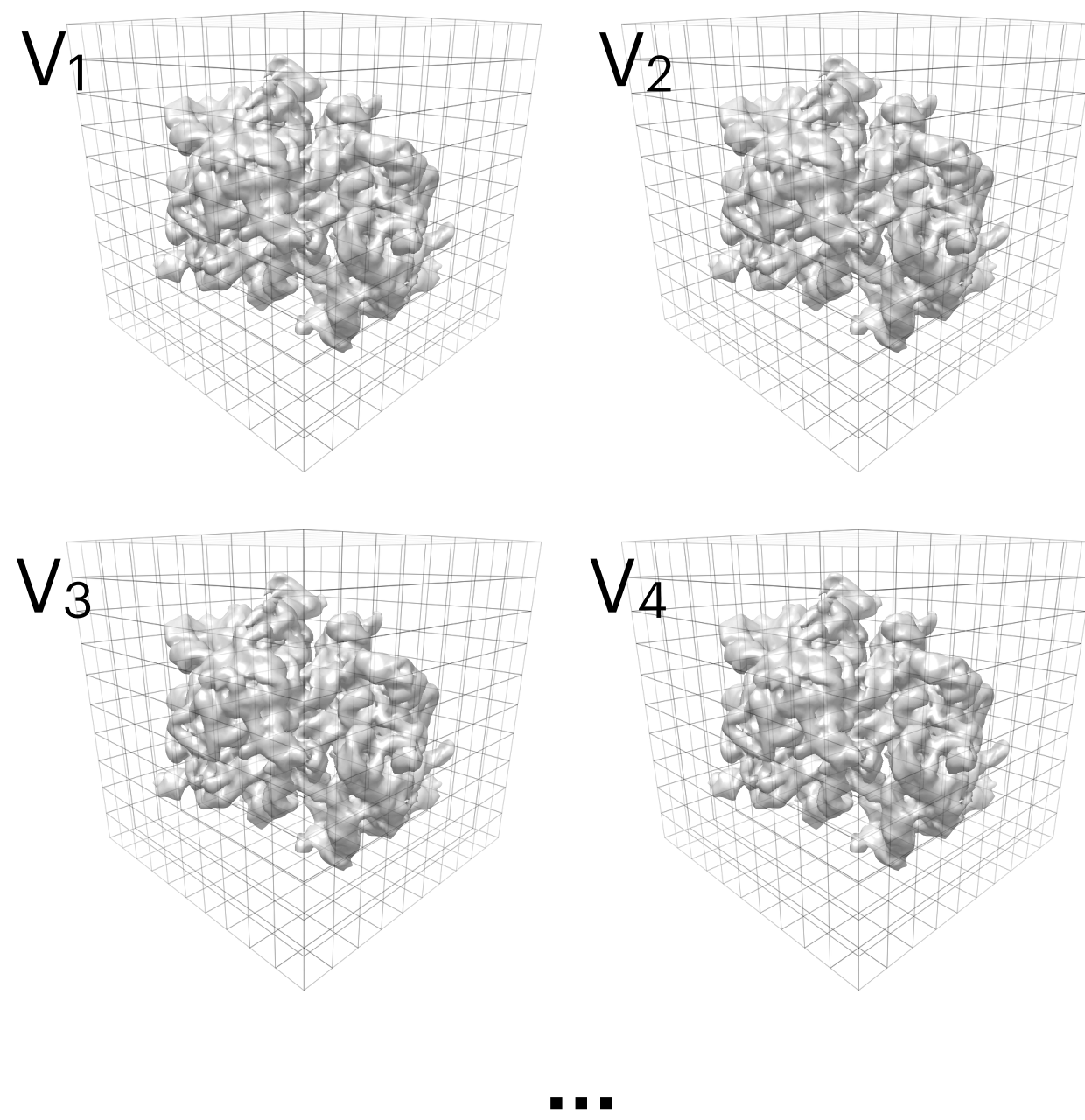
With sinusoidal encoding



Latent variable models for heterogeneous structures

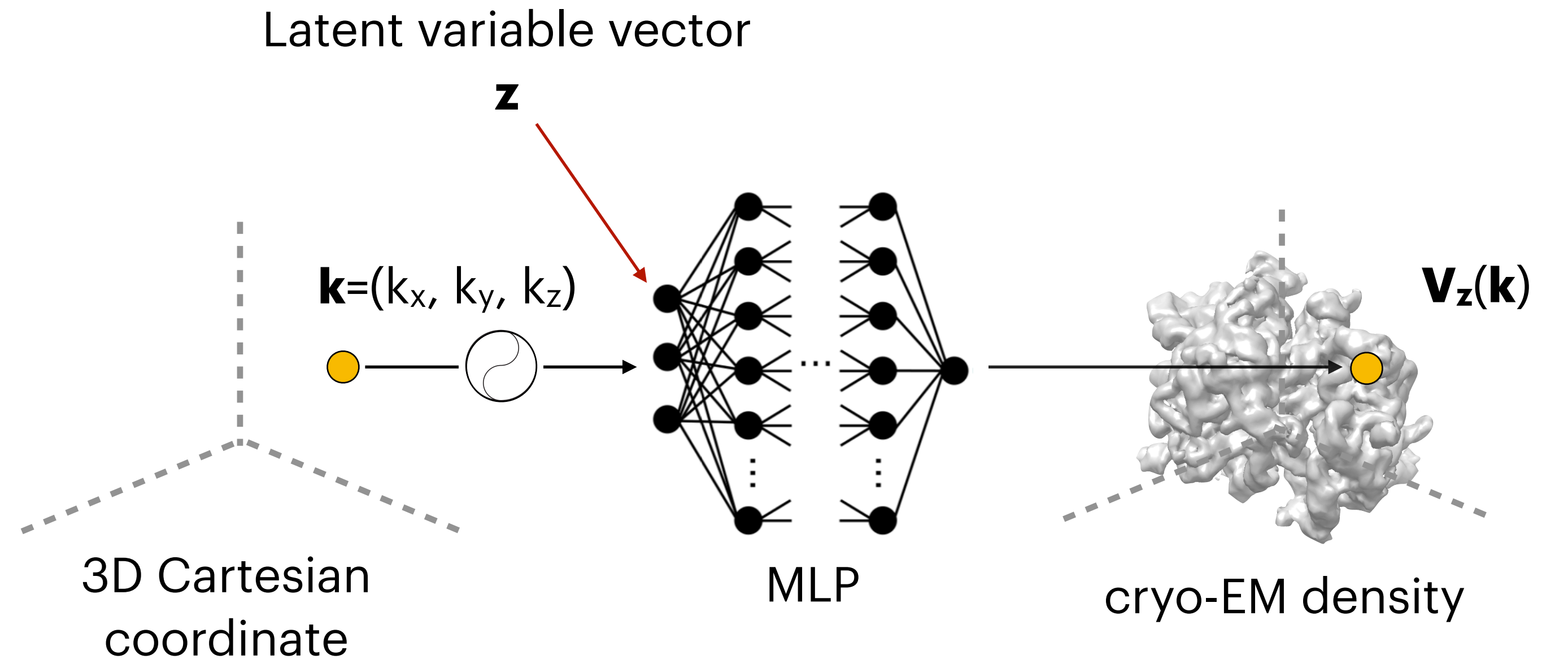
Multiclass refinement

V_z , where z in $\{1,2,3\dots,K\}$



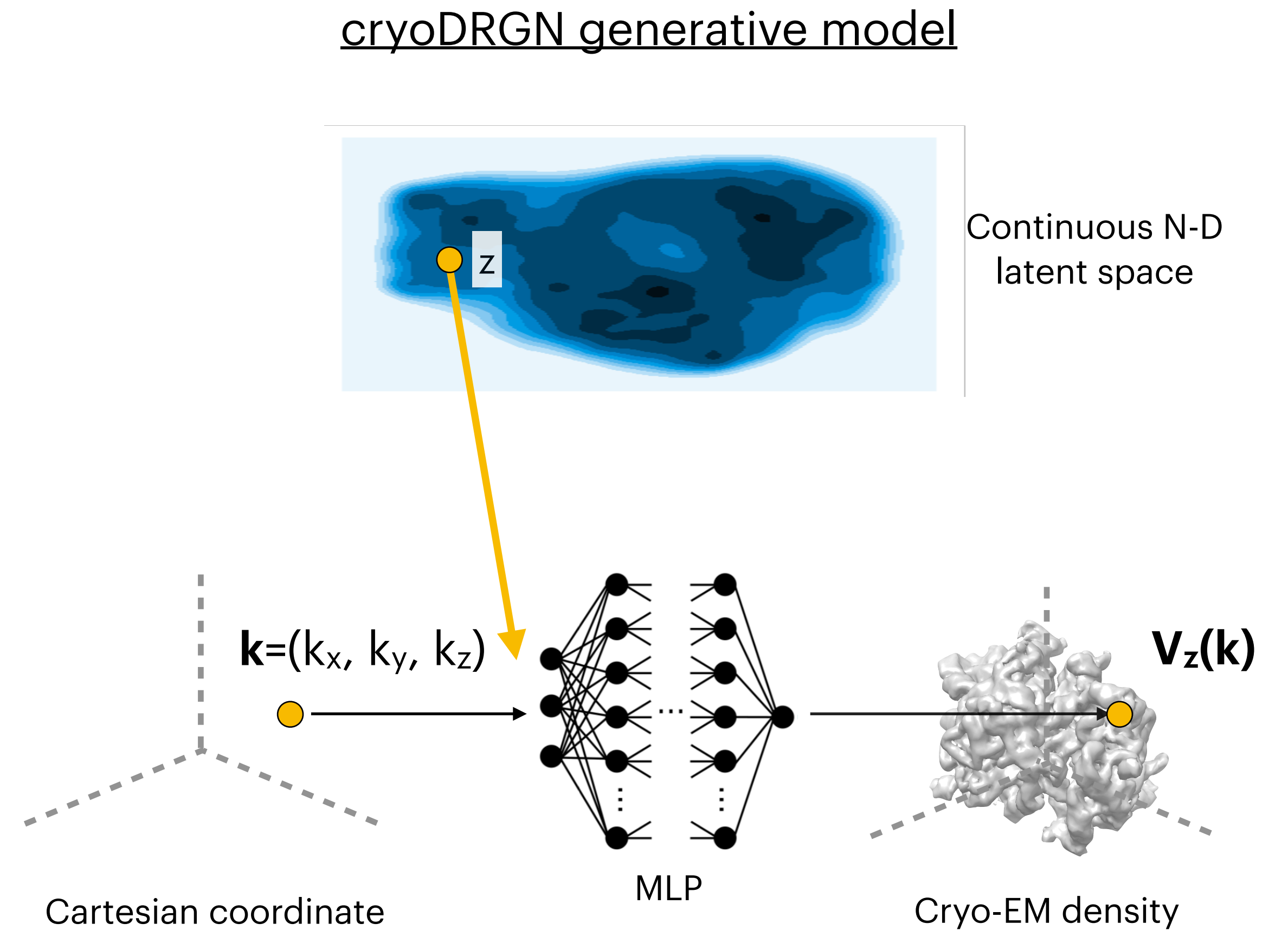
- Manual selection of K and initial volumes
- Typically, $K < 10$

cryoDRGN



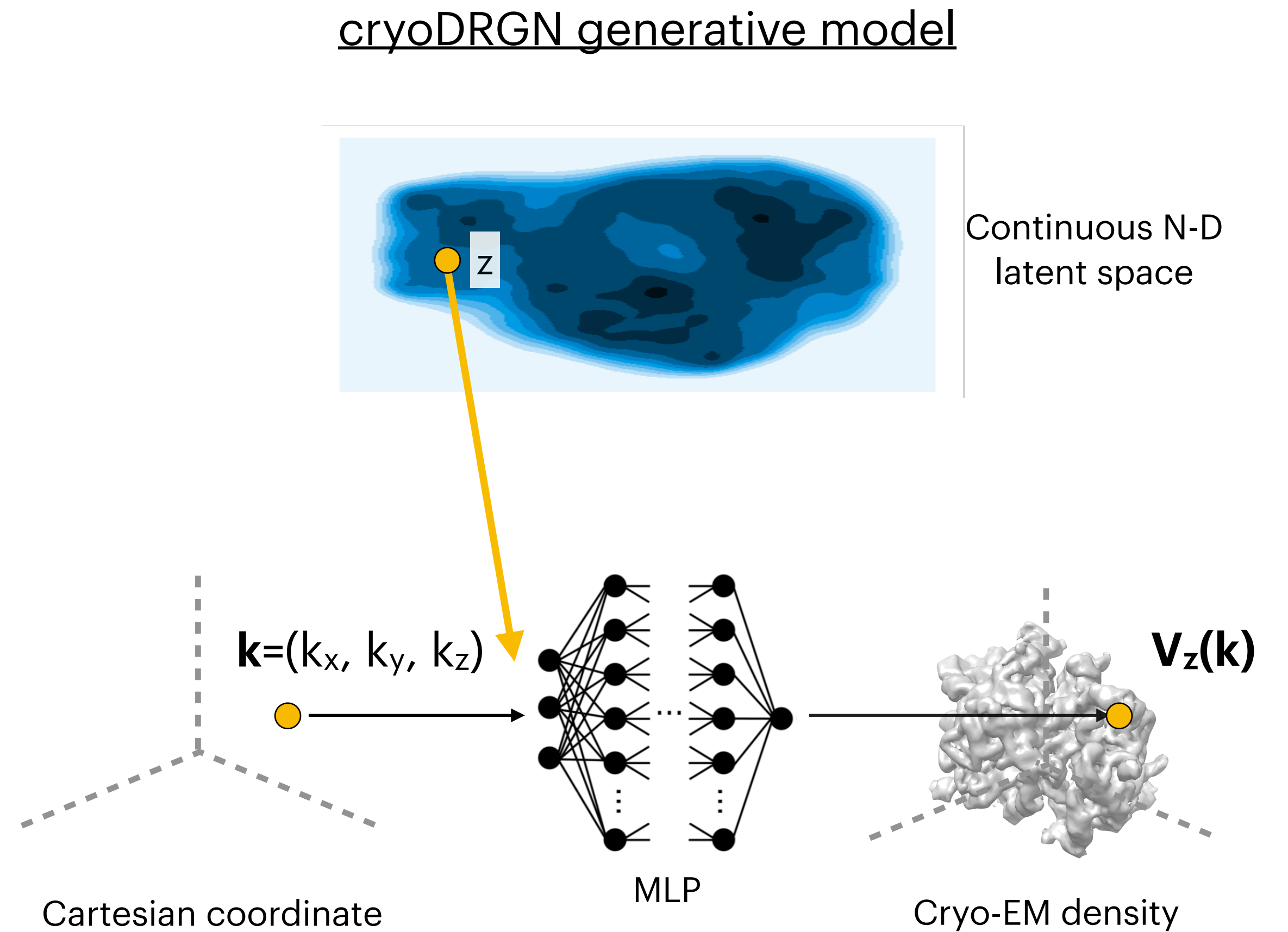
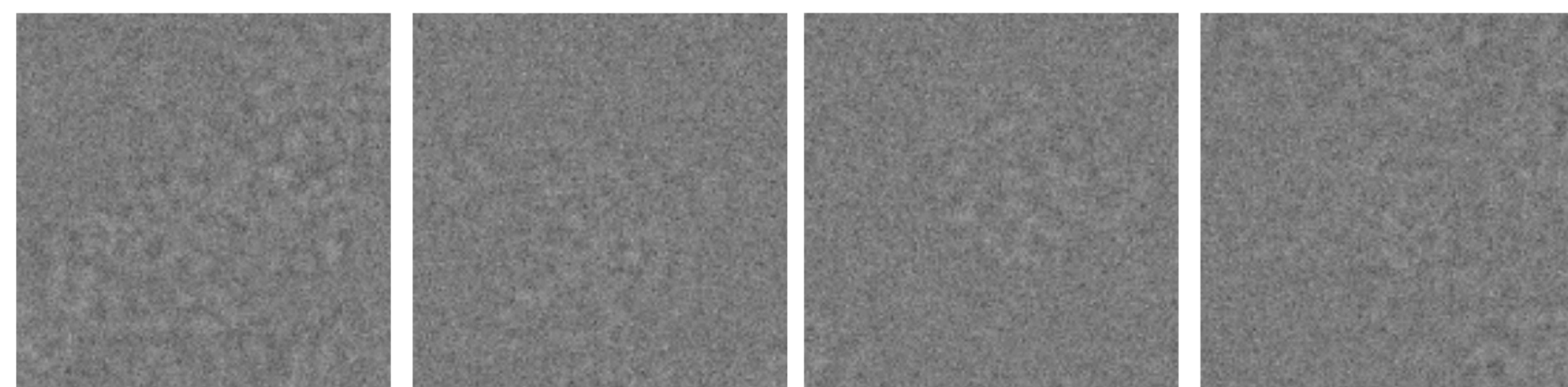
CryoDRGN's continuous latent variable model

- Extend the neural representation of volume with a conditional latent variable model
- Encodes a N-dimensional *continuous* distribution over structures



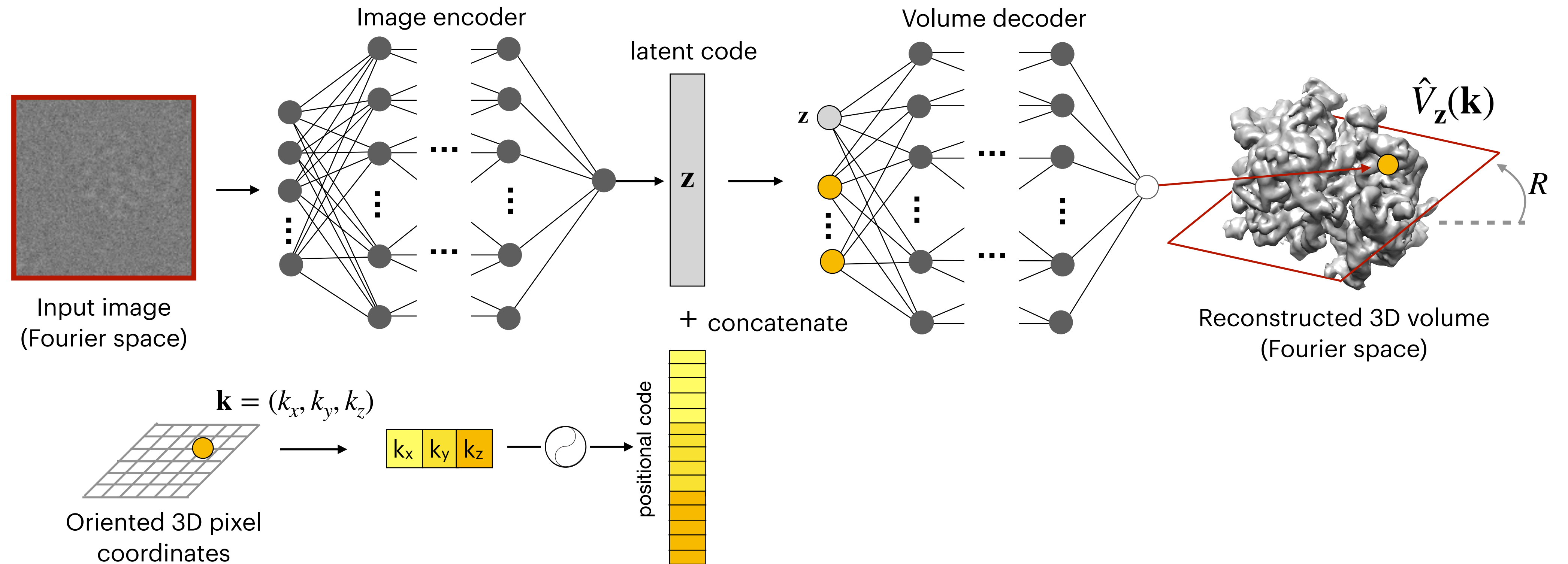
CryoDRGN's continuous latent variable model

- Extend the neural representation of volume with a conditional latent variable model
- Encodes a N-dimensional *continuous* distribution over structures
- How to learn such a model from data?



CryoDRGN's overall architecture

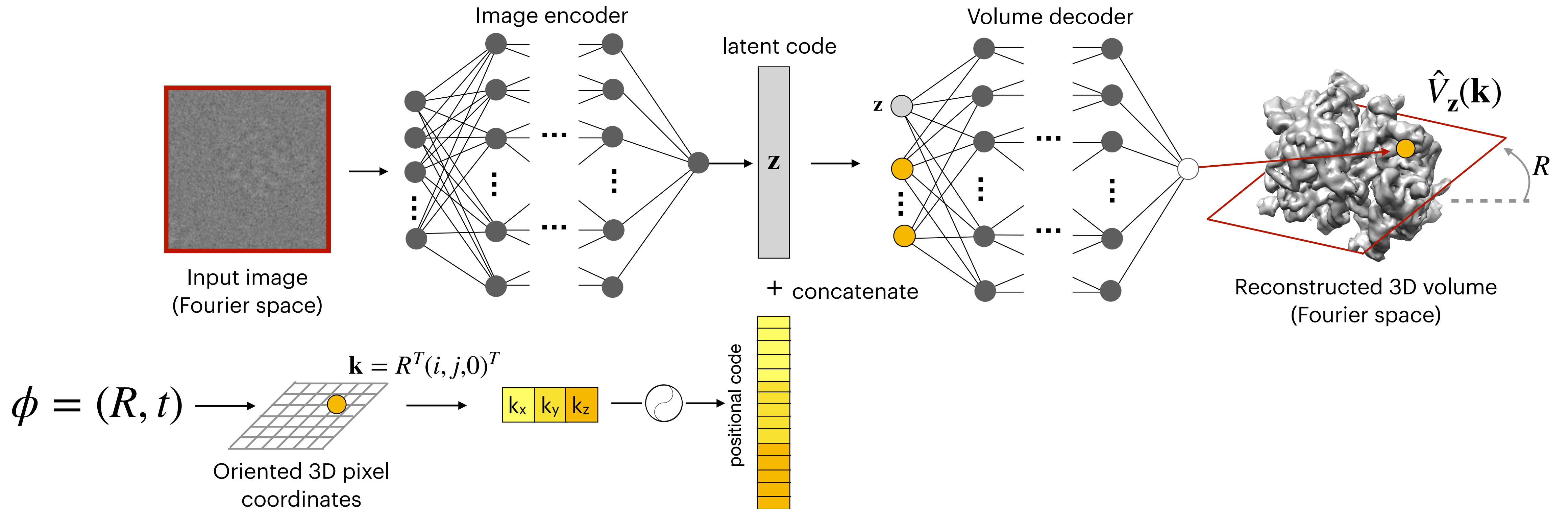
- We propose a Fourier domain **image encoder - volume decoder** architecture based on the VAE



- The decoder reconstructs an image pixel-by-pixel given \mathbf{z} and the 3D coordinates of the pixels
- Coordinate-based volume architecture enforces geometric consistency between 2D views (Fourier slice theorem)

CryoDRGN's overall architecture

- We propose a Fourier domain **image encoder - volume decoder** architecture based on the VAE



- To obtain oriented 3D pixel coordinates, a coordinate lattice on the x-y plane is rotated by R
- For each image, we need to approximate its pose $\phi = (R, t)$. How?

Possible paradigms for pose inference

- Amortized variational inference [1]

$$\phi_i \sim q_{\xi}(\phi | X)$$

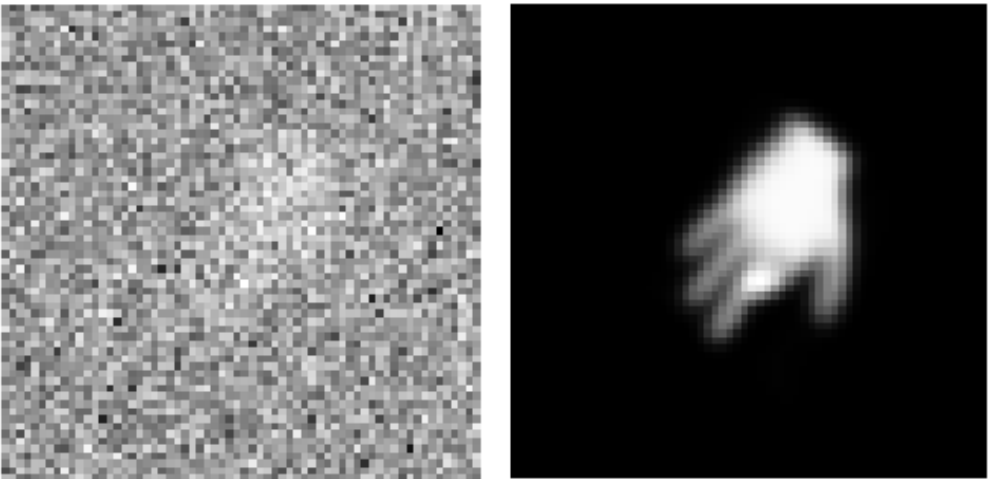
- Gradient descent [2]

$$\phi^{(n+1)} = \phi^{(n)} - \alpha \nabla_{\phi} \mathcal{L}(\phi)$$

- Distribution matching/GANs [3]

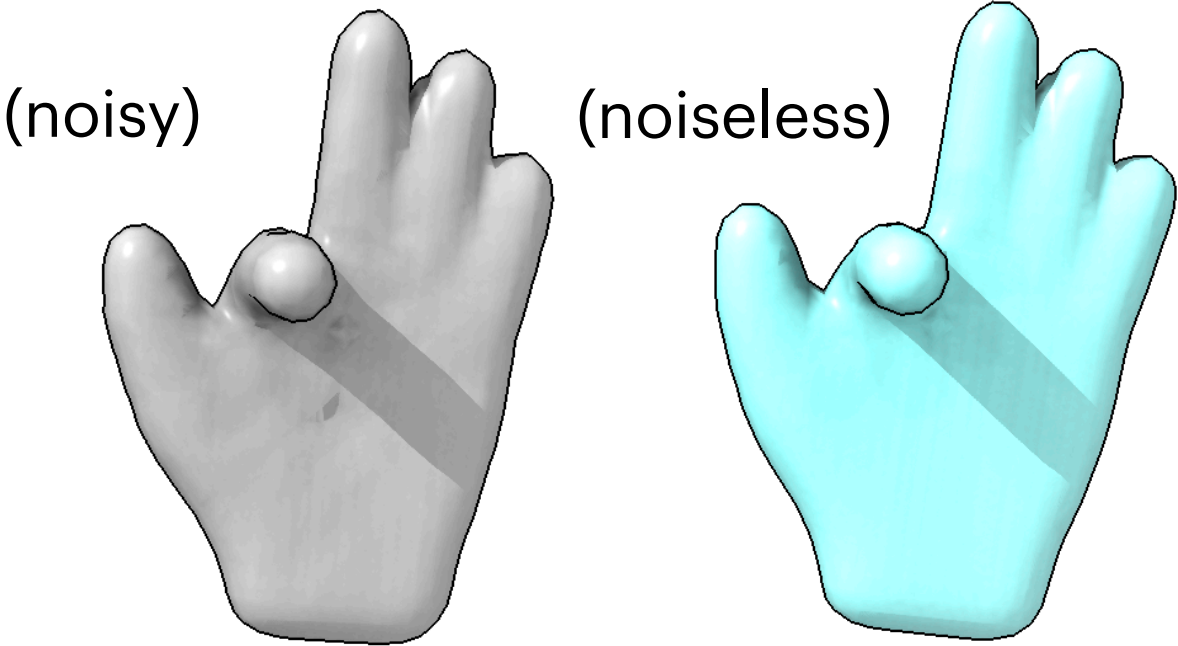
$$\underset{V}{\operatorname{argmin}} D(p_{sim}(X | V), p_{data}(X))$$

Spurious local minima in the training objective

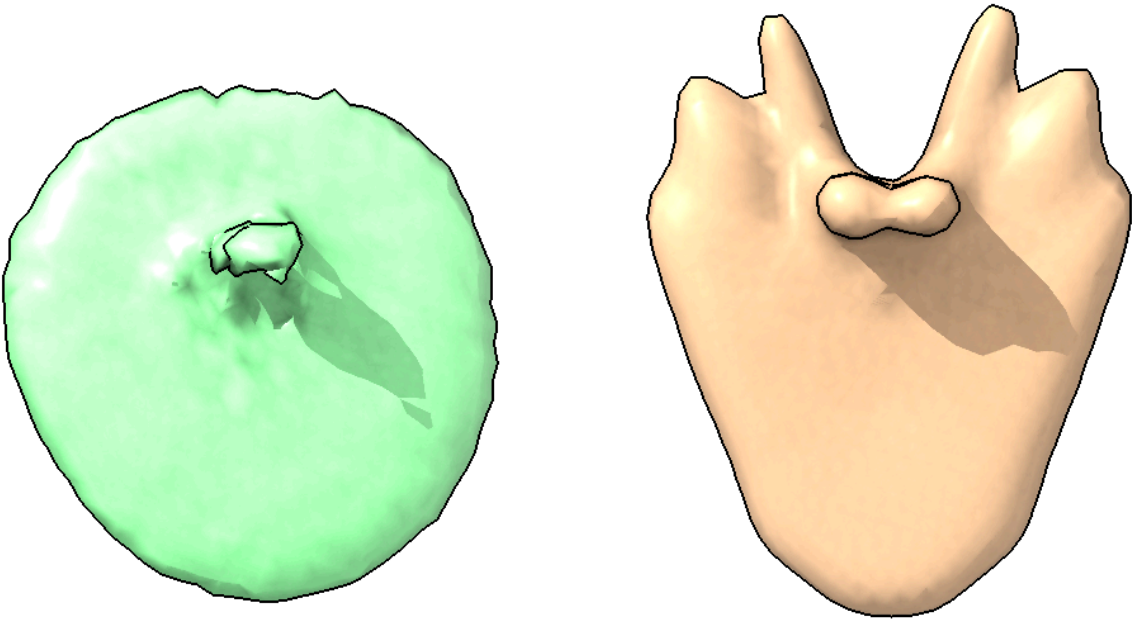


Example image

Ground truth poses

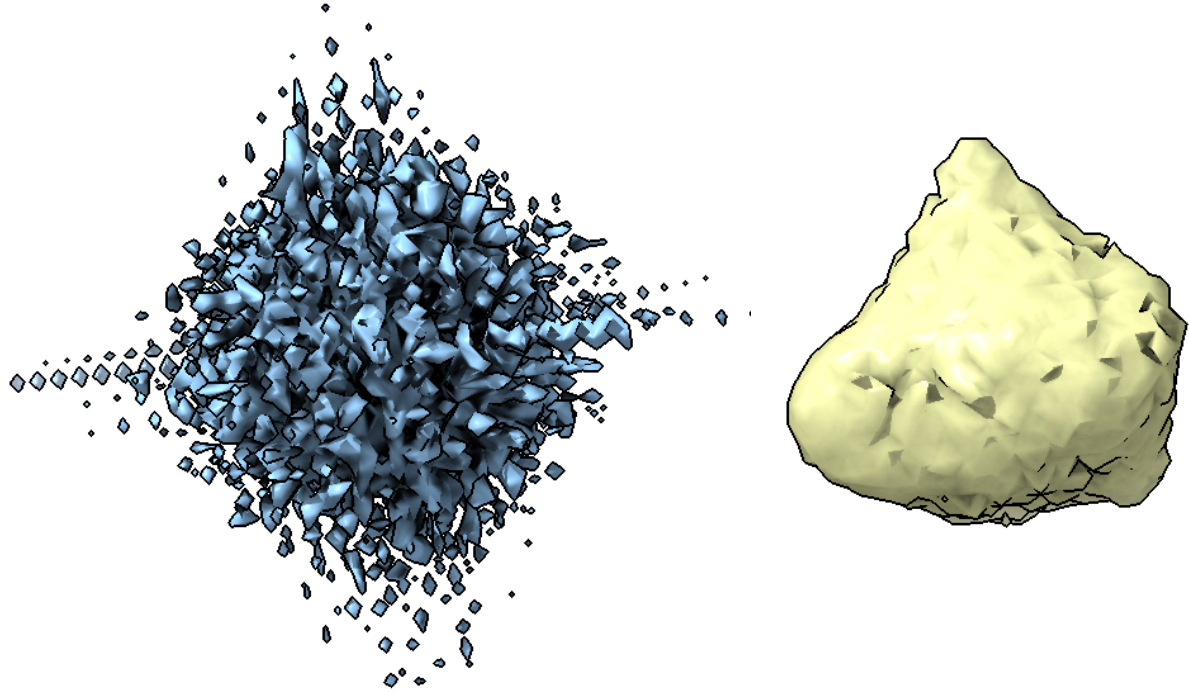


Amortized Variational Inference



$$\phi_i \sim q_{\xi}(\phi | X)$$

Pose SGD

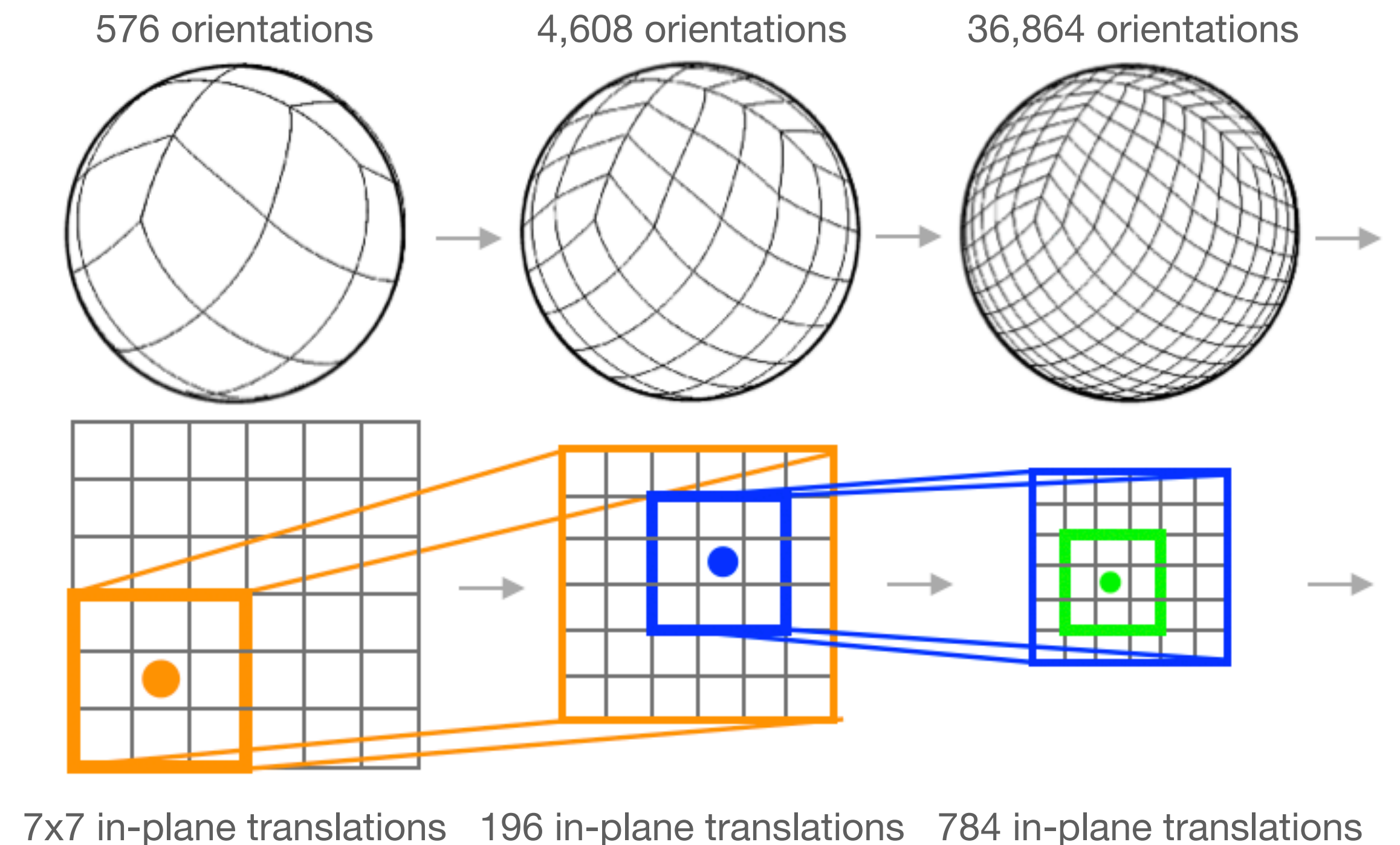


$$\phi^{(n+1)} = \phi^{(n)} - \alpha \nabla_{\phi} \mathcal{L}(\phi)$$

Search algorithms for inference of image pose

- Instead we perform a global search over a discretization of $SO(3) \times \mathbb{R}^2$ for the MLE pose for each image X_i given the current decoder V_θ
- A hierarchical search procedure:
 - Start with an exhaustive search over a discretization of the 5D space of poses
 - A uniform discretization of $SO(3)$ with the Hopf fibration, regular 2D grid for in-plane translations
 - Iteratively refine the poses by keeping the top K poses that minimize the reconstruction loss
 - Choose K via a branch-and-bound procedure¹

$$\operatorname{argmax}_{\phi_i} p(X_i | V_\theta)$$

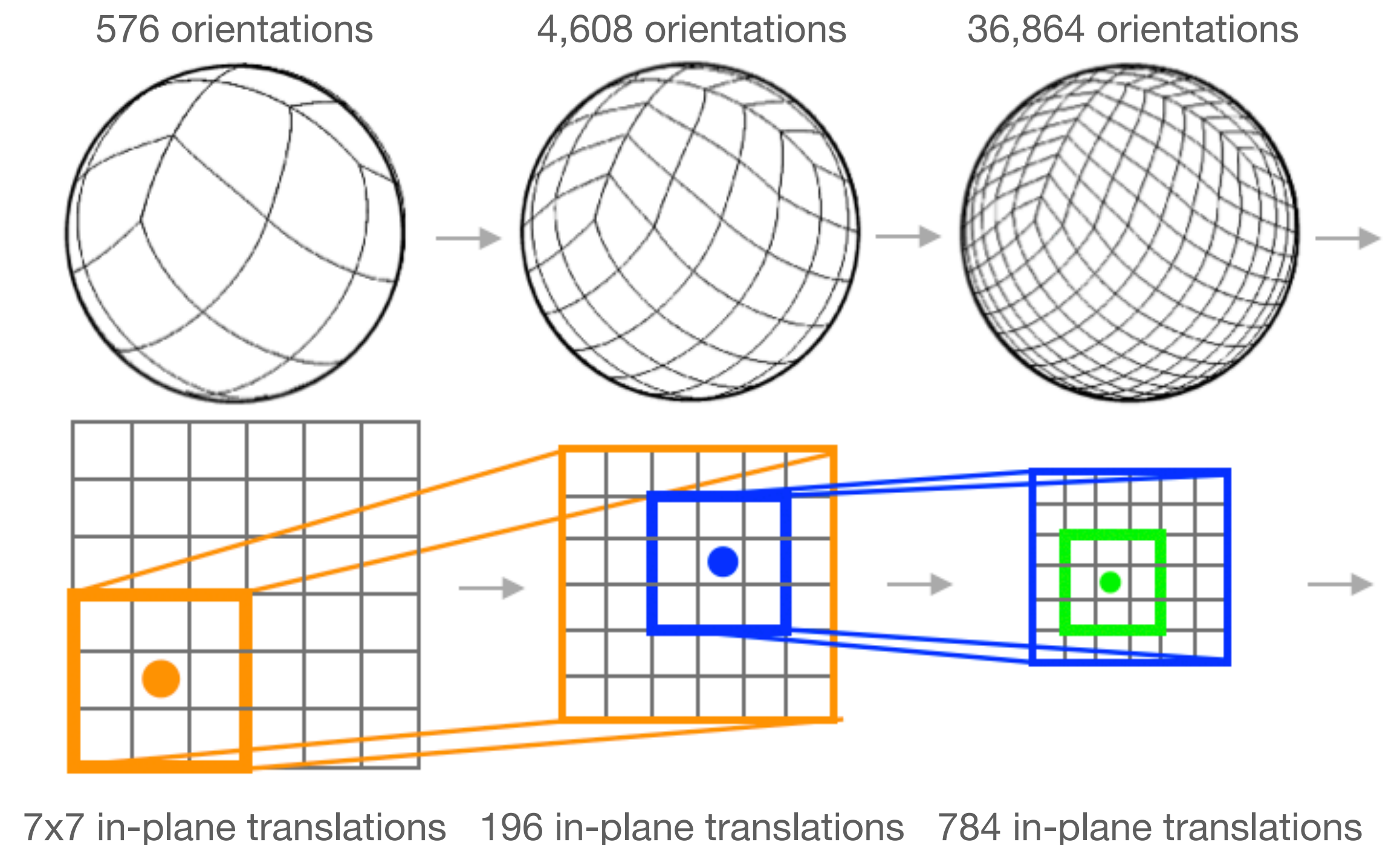


Search algorithms for inference of image pose

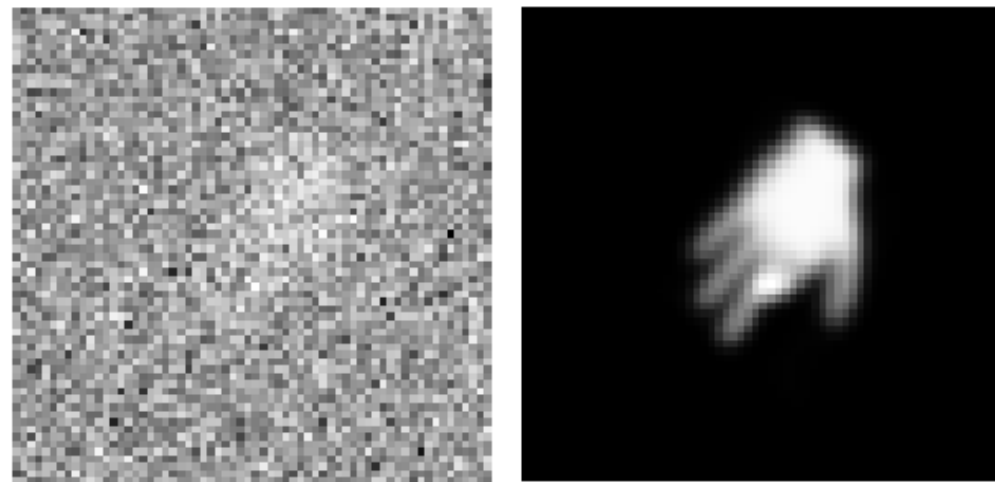
- Instead we perform a global search over a discretization of $SO(3) \times \mathbb{R}^2$ for the MLE pose for each image X_i given the current decoder V_θ

$$\operatorname{argmax}_{\phi_i} p(X_i | V_\theta)$$

- Frequency marching¹:
 - Band limit the loss function to low frequency components
 - Benefit 1: Computational efficiency
 - Benefit 2: Prevent overfitting

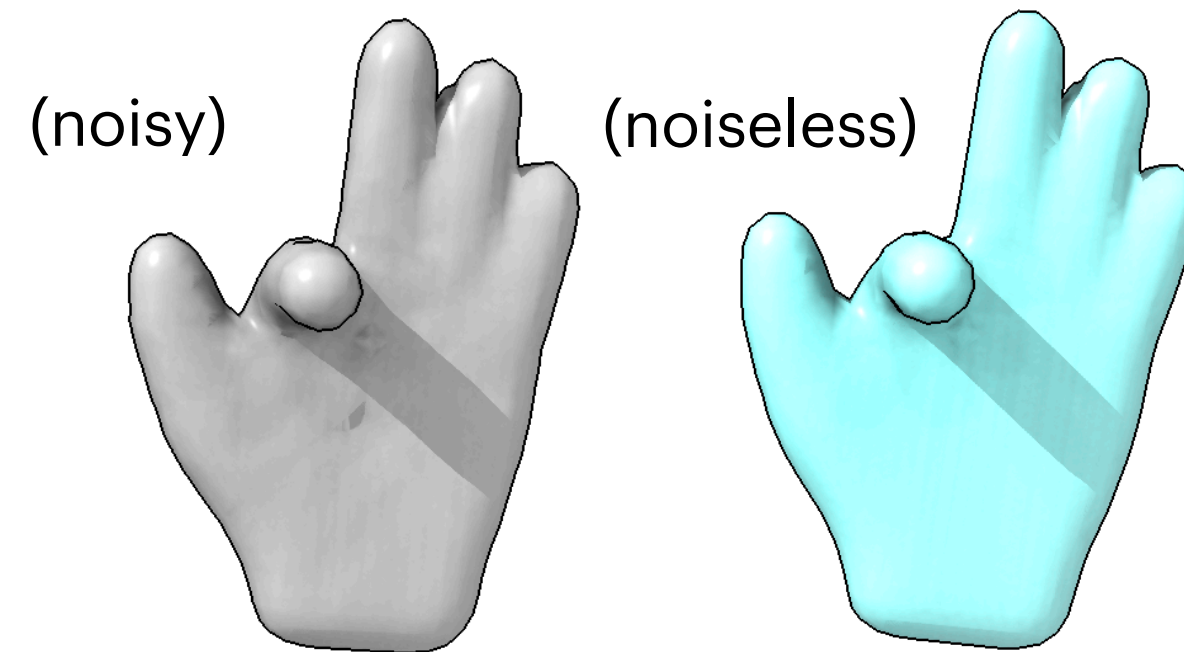


Spurious local minima in the training objective

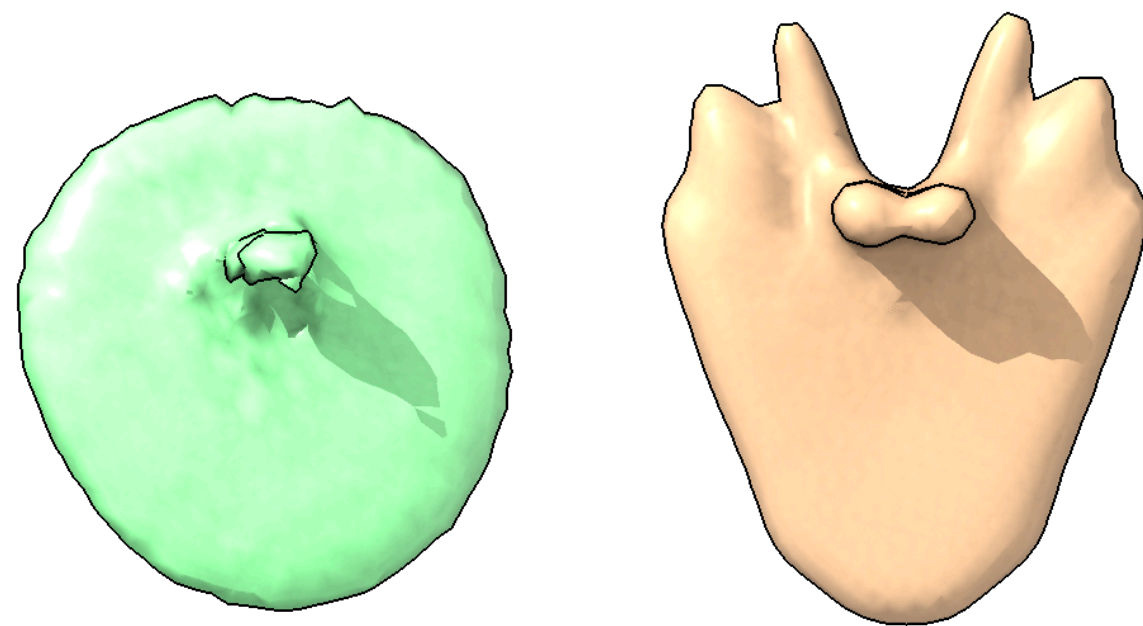


Example image

Ground truth poses

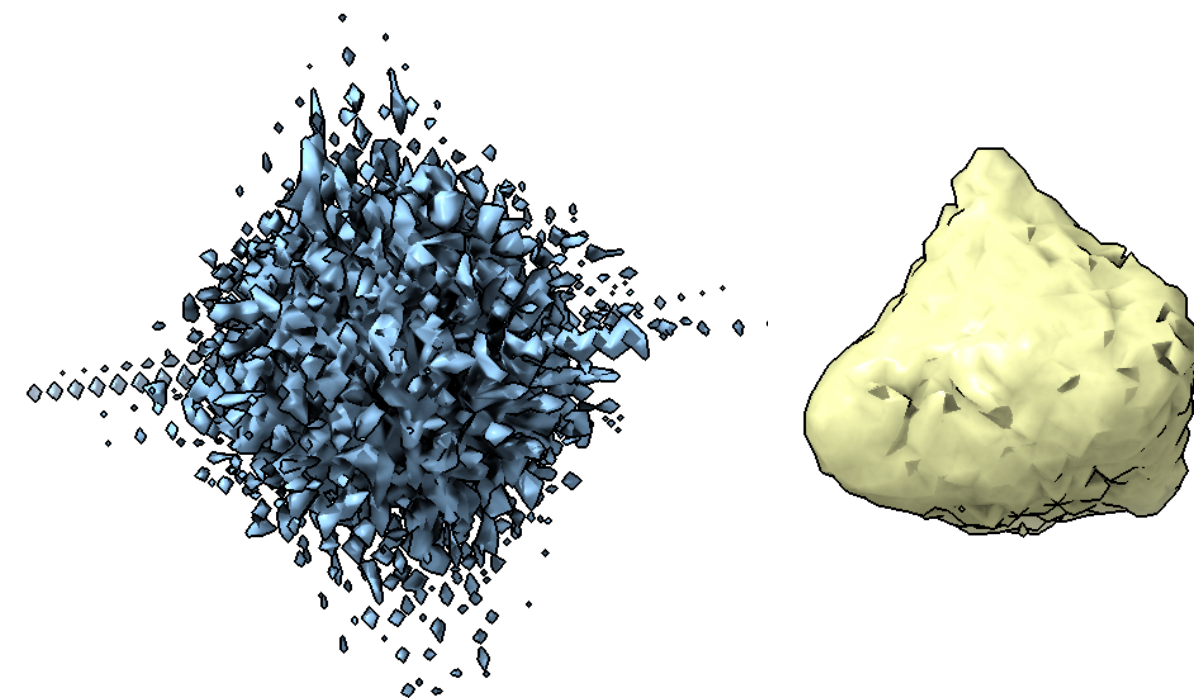


Amortized variational inference



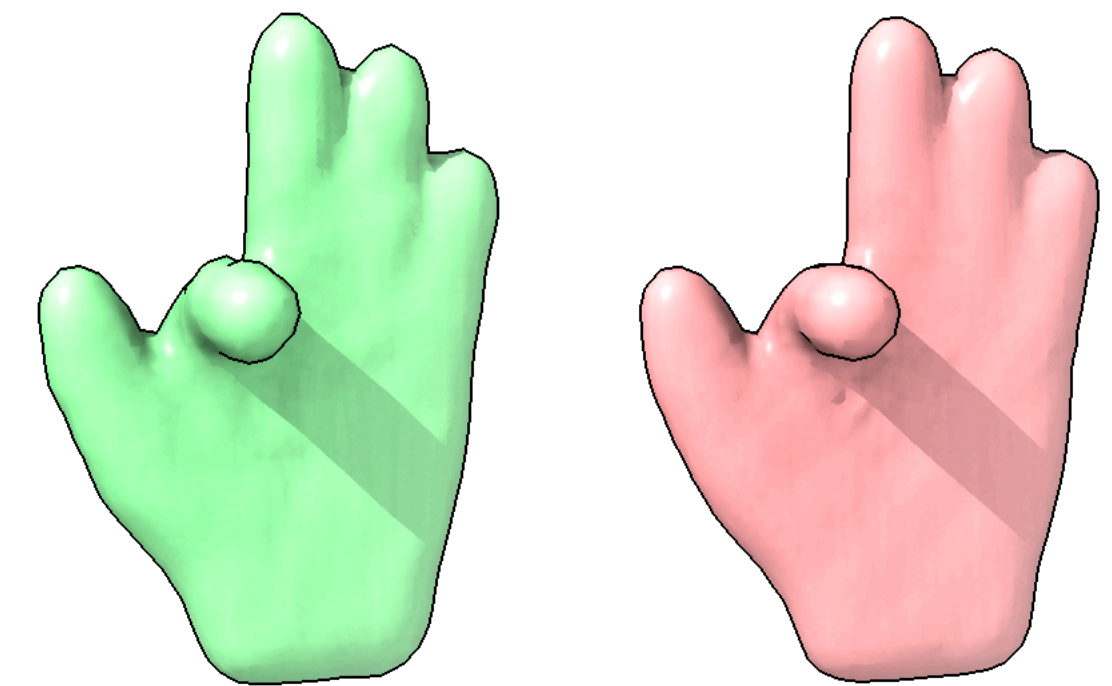
$$\phi_i \sim q_\xi(\phi | X)$$

Pose SGD



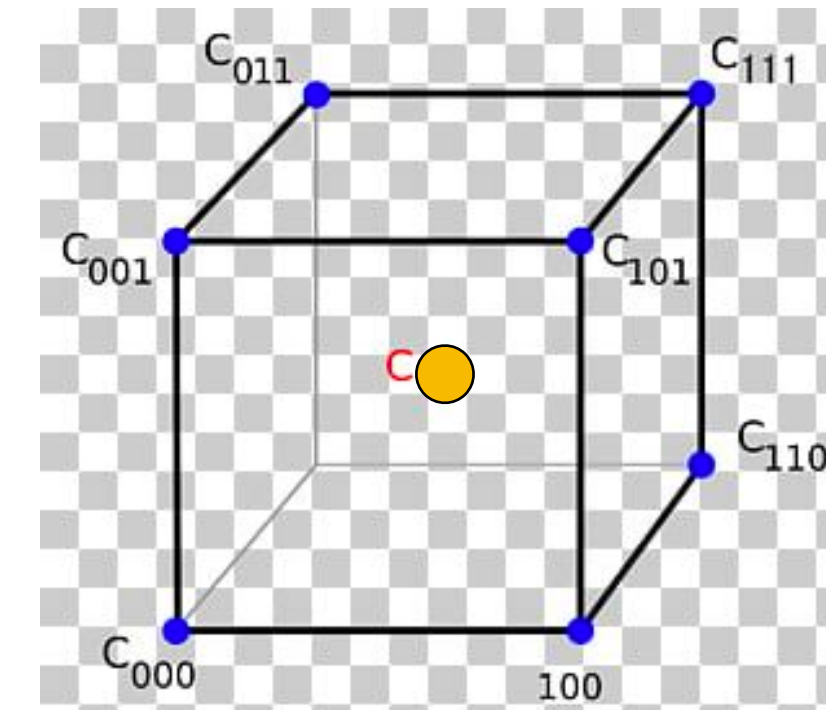
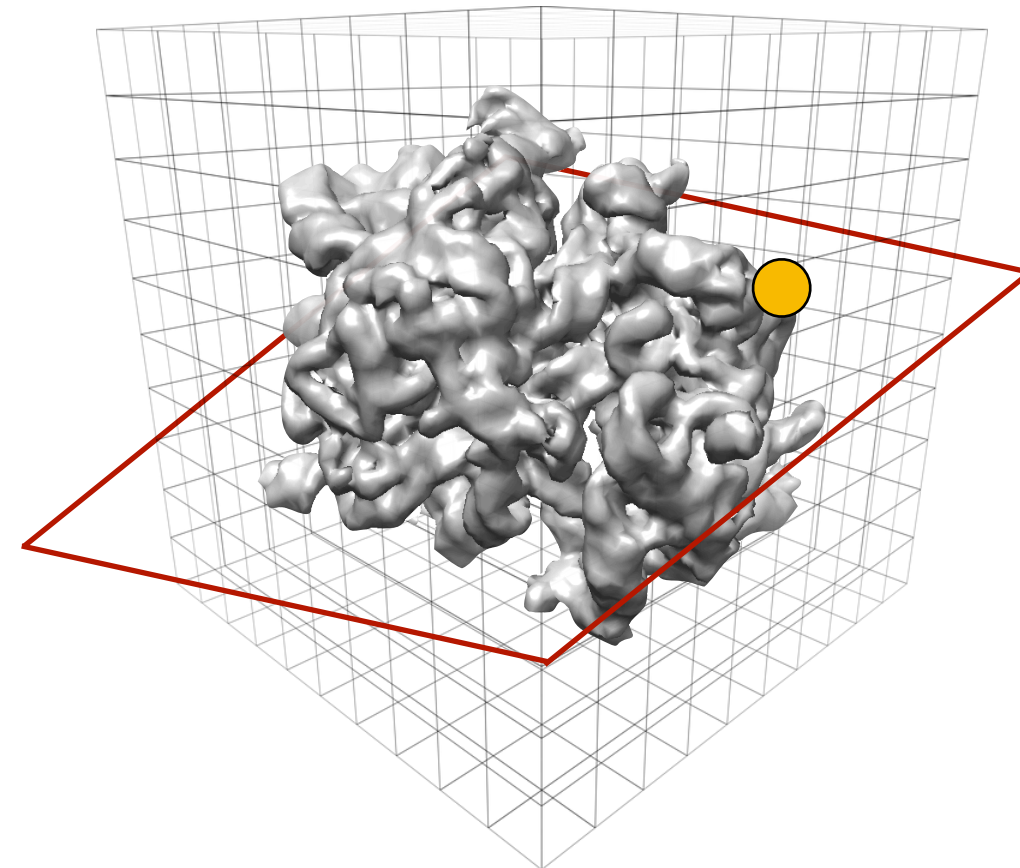
$$\phi^{(n+1)} = \phi^{(n)} - \alpha \nabla_\phi \mathcal{L}(\phi)$$

Hierarchical search



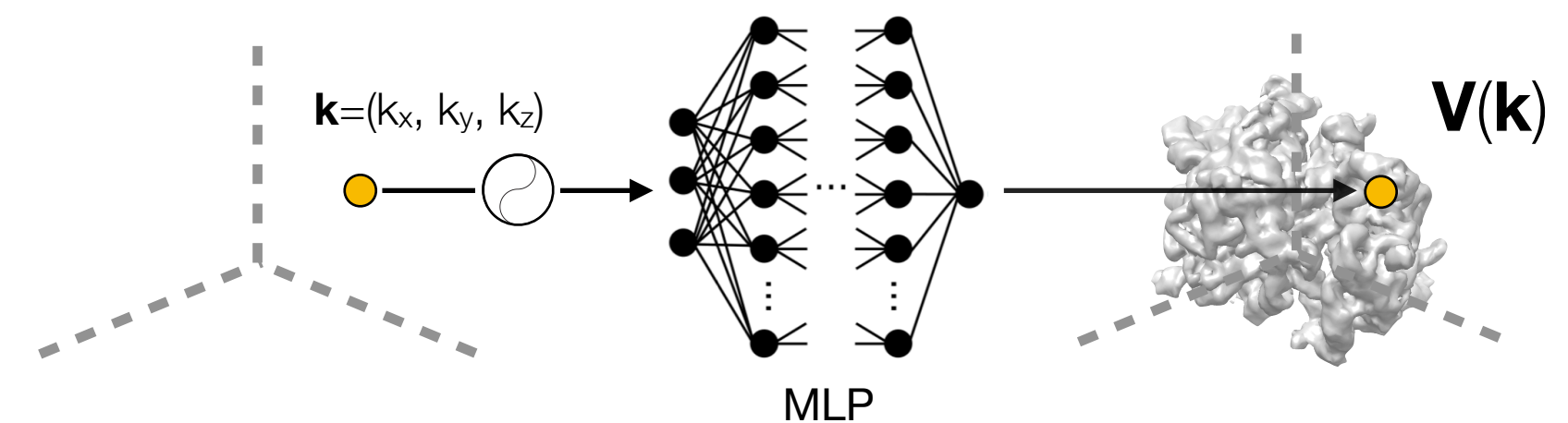
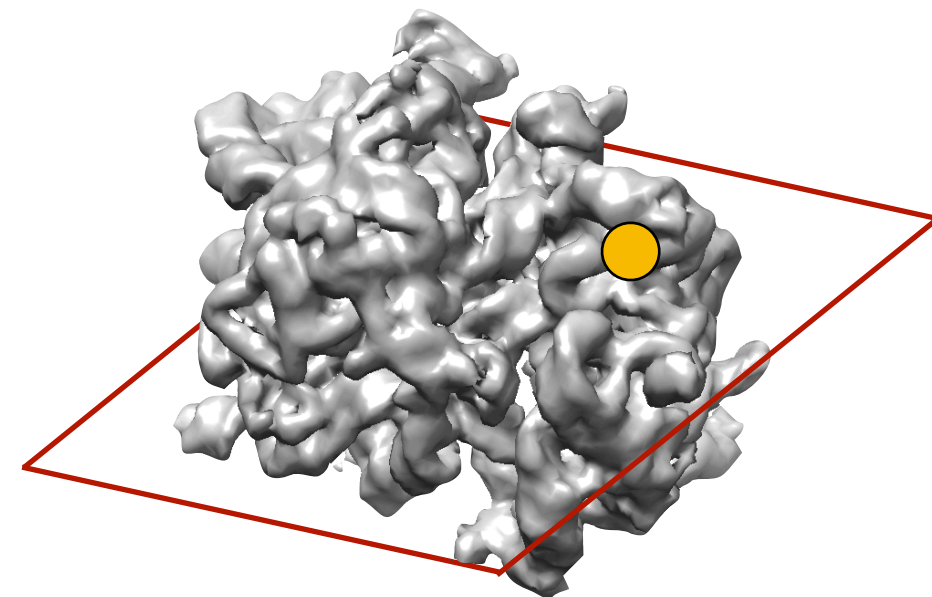
Pose search: Traditional vs. neural

Traditional



Each off-voxel point is computed as the weighted average of its 8 spatially closest neighbors

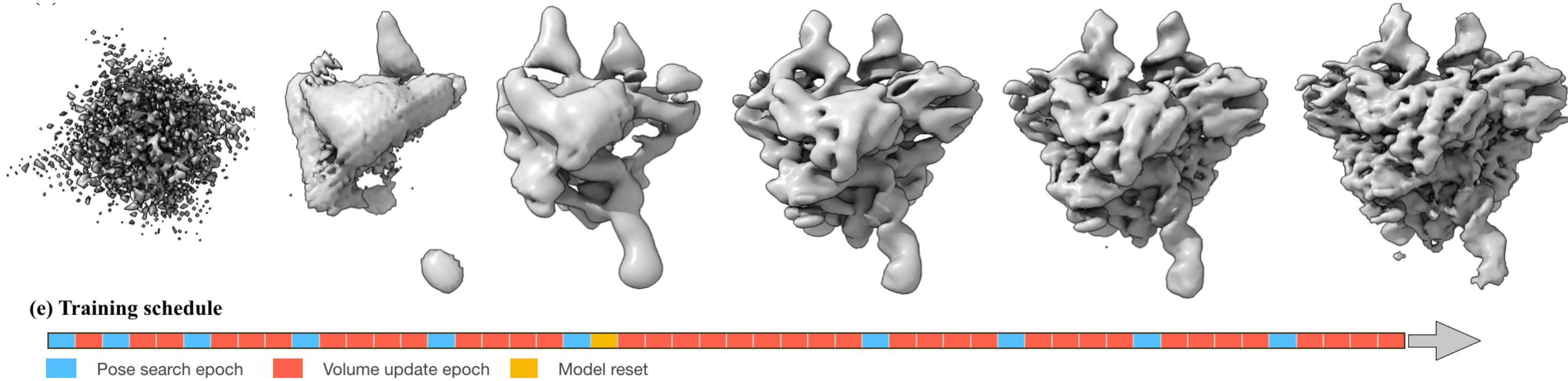
cryoDRGN



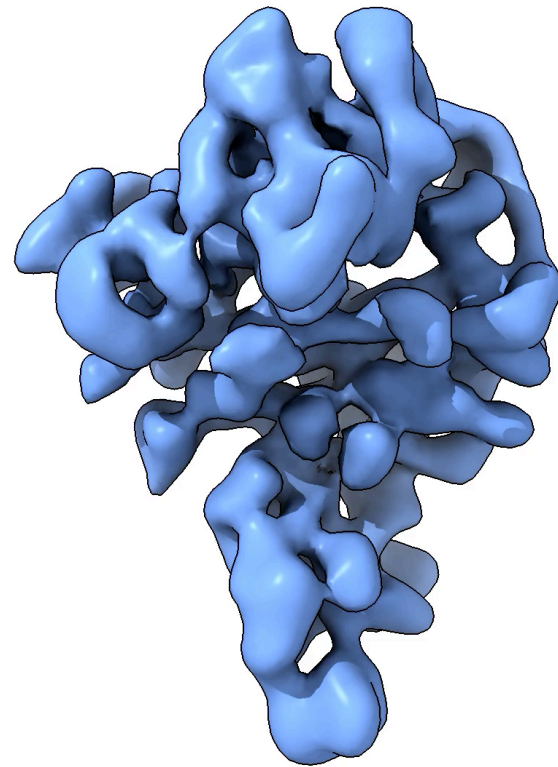
Each off-voxel point is computed by evaluating the MLP

cryoDRGN models are much more expensive to evaluate

CryoDRGN2: *Ab initio* heterogeneous reconstruction of real data

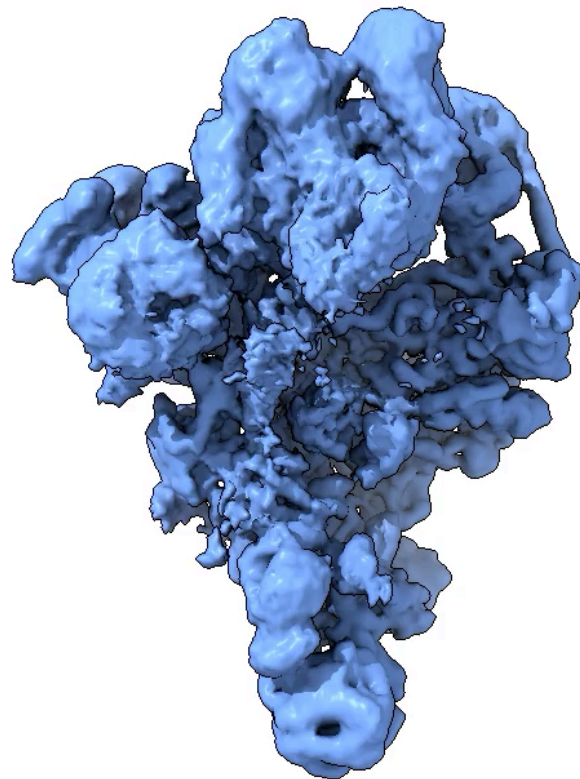


cryoDRGN1



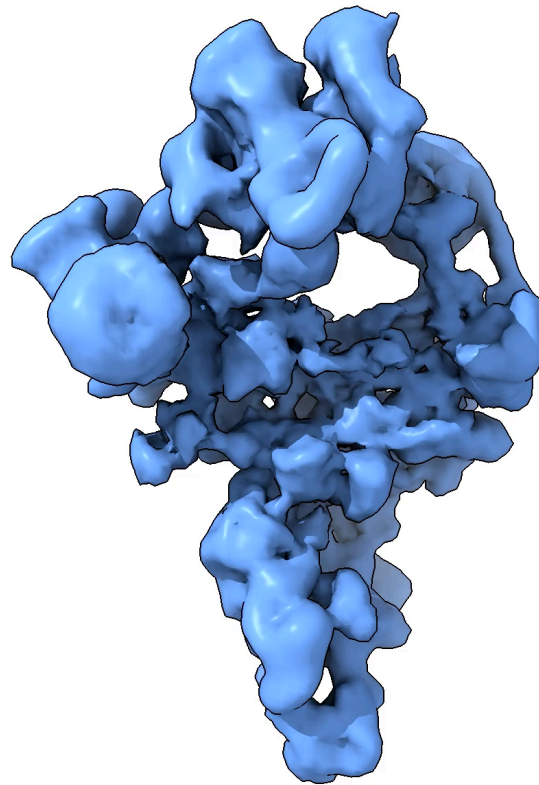
Training time: 38.3 hours (1 GPU)

cryoDRGN, pose supervision



3.8 h

cryoDRGN2

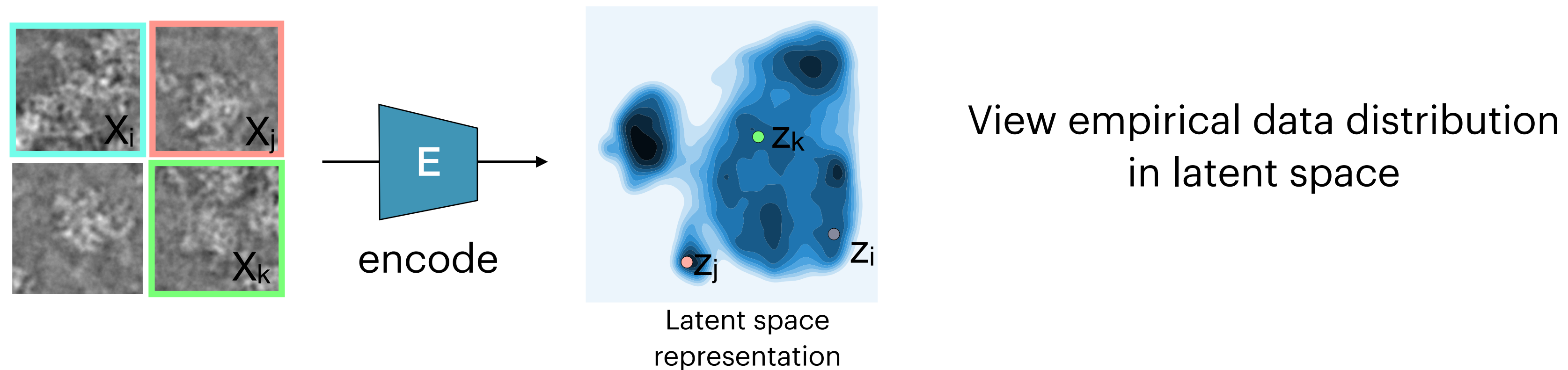


11.8 h

Analyzing the generative model

CryoDRGN at test time:

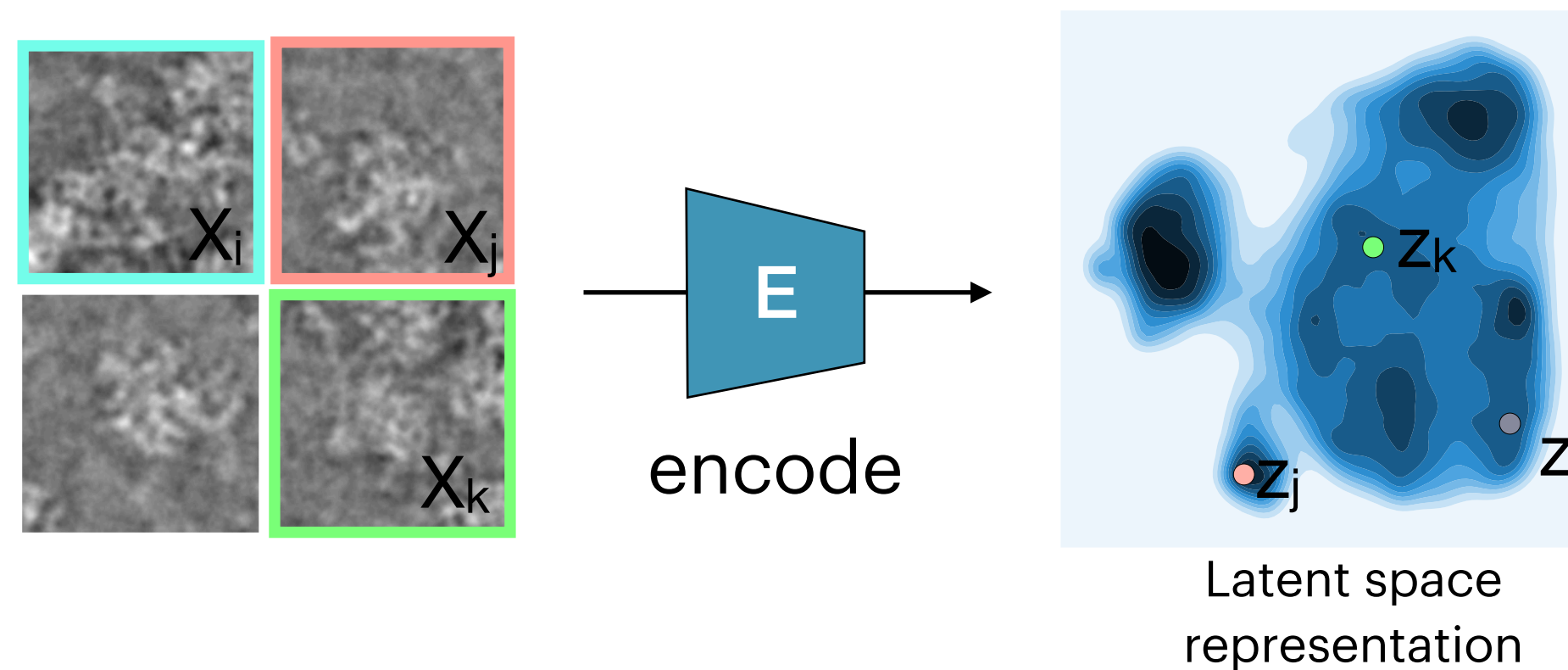
- Use the encoder network to evaluate the latent embedding \mathbf{z} for each image



Analyzing the generative model

CryoDRGN at test time:

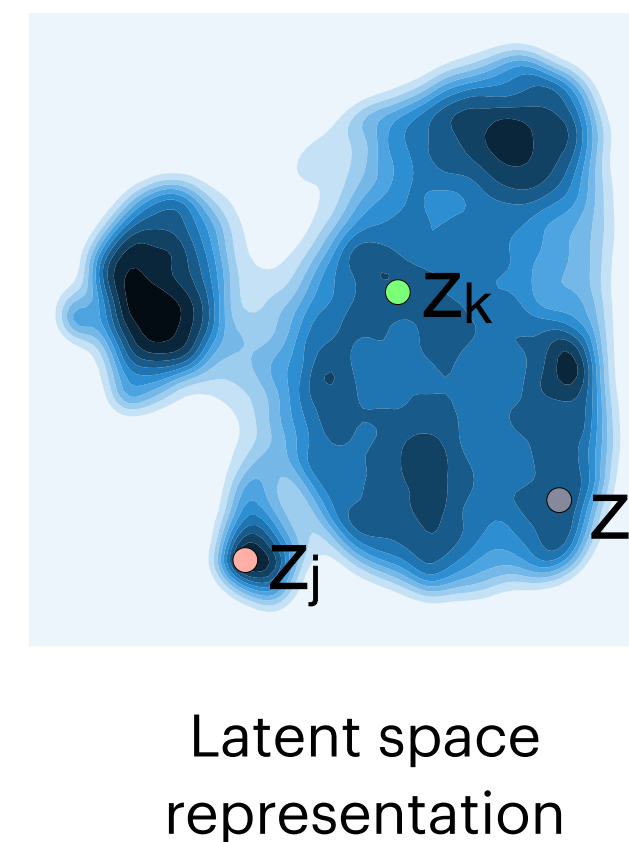
- Use the encoder network to evaluate the latent embedding \mathbf{z} for each image



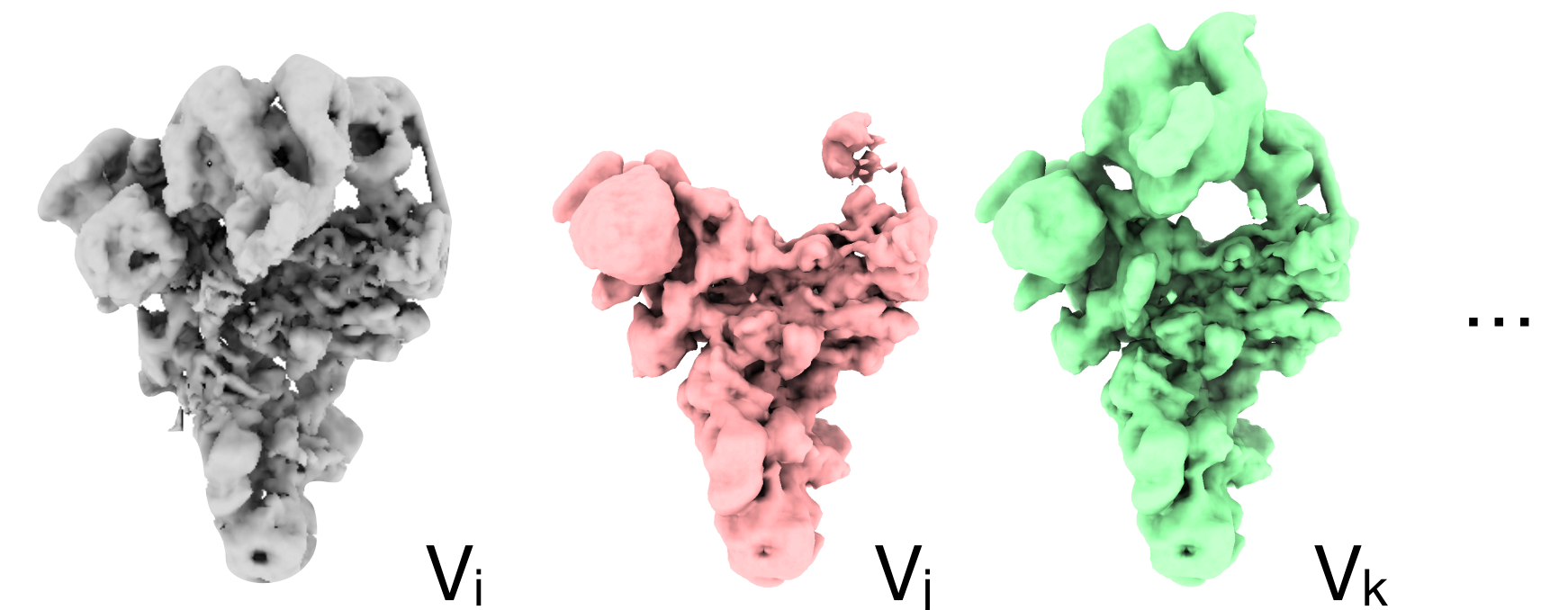
View empirical data distribution in latent space

- Use the decoder network to generate \mathbf{V} at different values of \mathbf{z}

View the structural ensemble and generate movies from trajectories in latent space



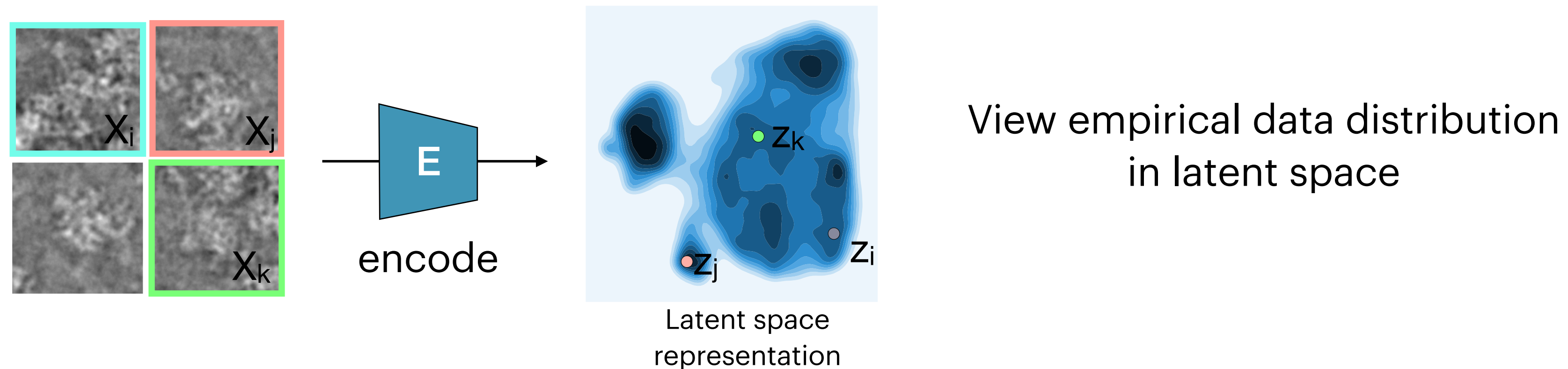
Representative samples



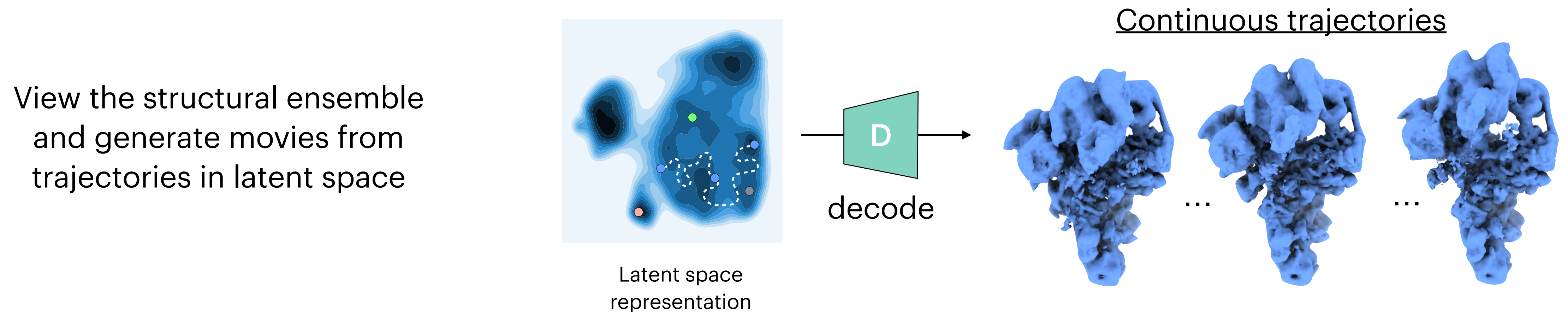
Analyzing the generative model

CryoDRGN at test time:

- Use the encoder network to evaluate the latent embedding \mathbf{z} for each image



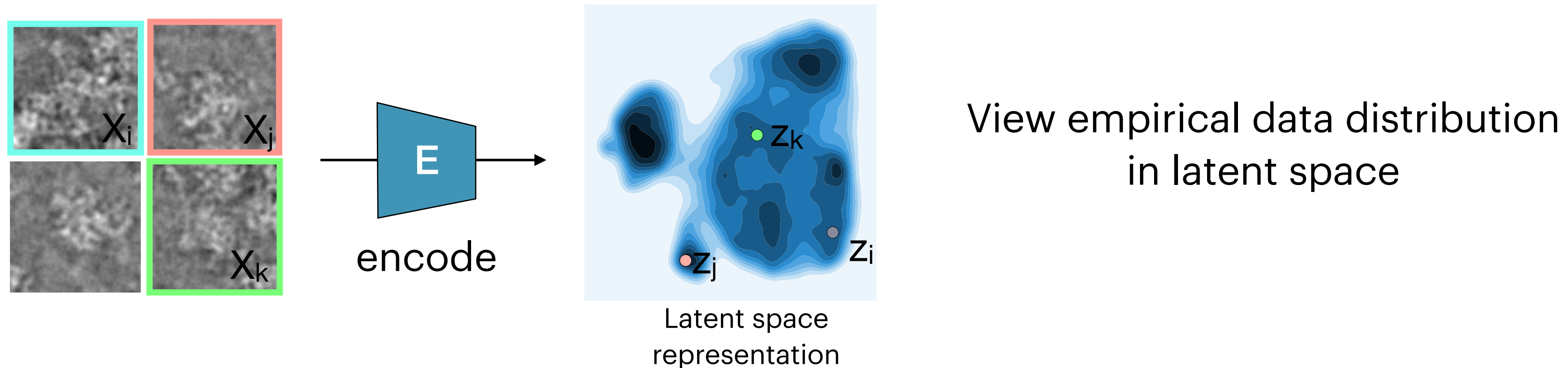
- Use the decoder network to generate \mathbf{V} at different values of \mathbf{z}



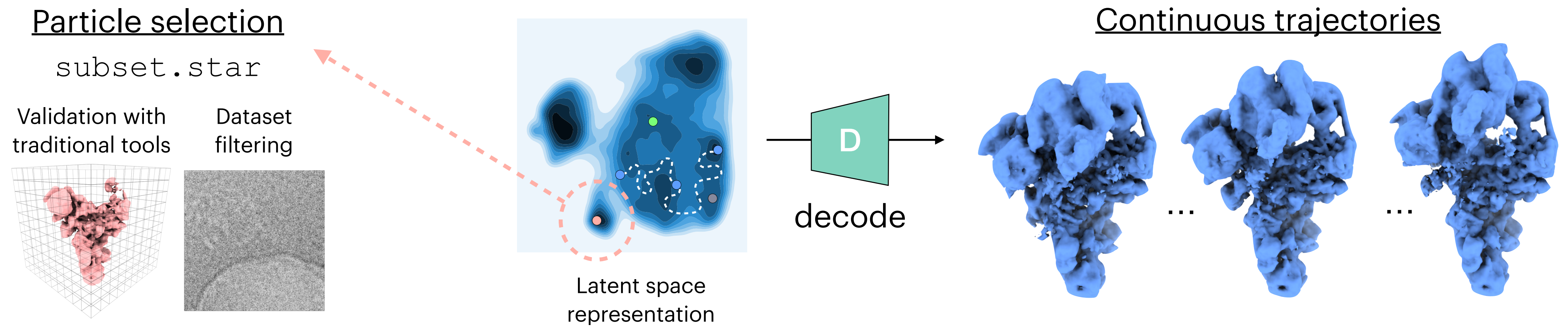
Analyzing the generative model

CryoDRGN at test time:

- Use the encoder network to evaluate the latent embedding \mathbf{z} for each image



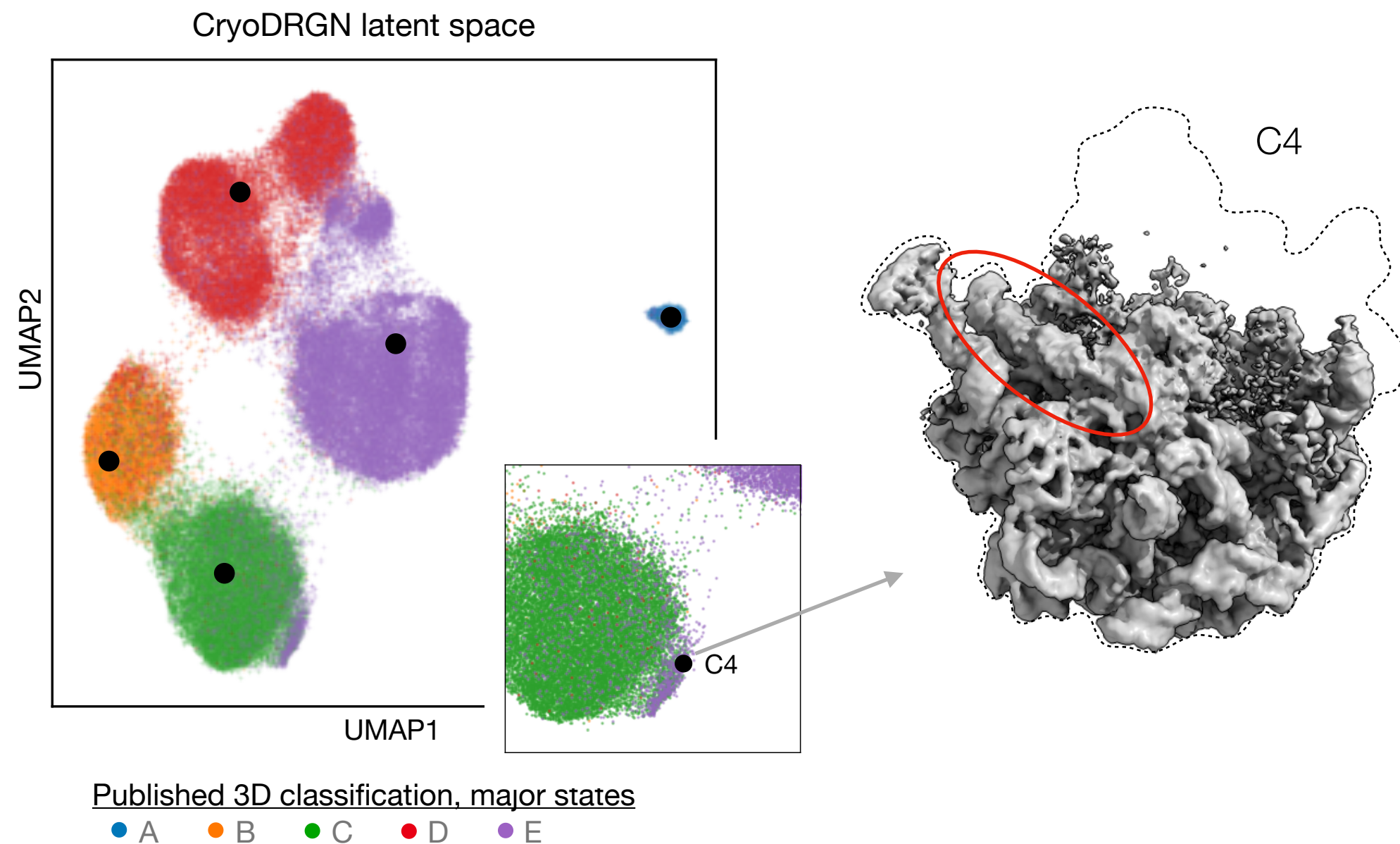
- Use the decoder network to generate \mathbf{V} at different values of \mathbf{z}



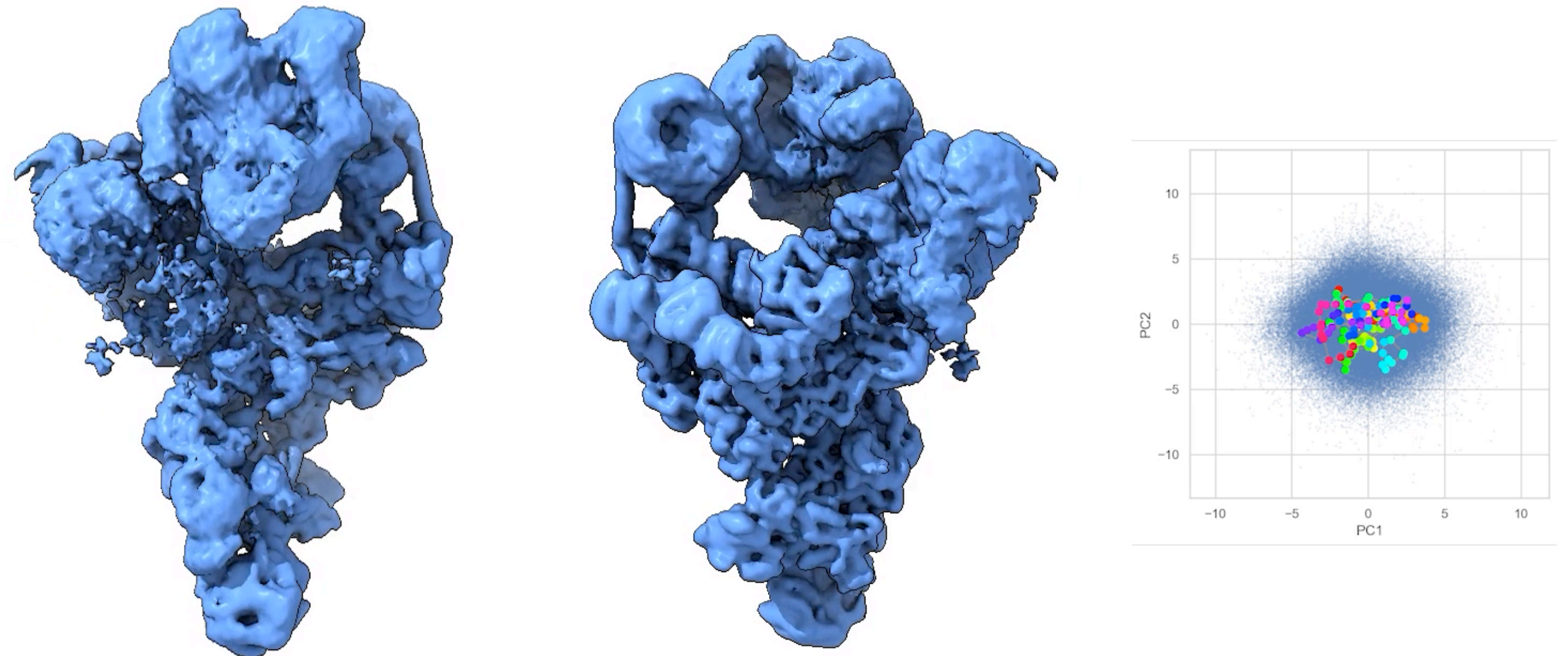
CryoDRGN: Applications and software

- Released as an open source software tool for training cryoDRGN models and interpreting results

Discovery of new structures



Visualization of continuous dynamics



CryoDRGN software



Software Pipeline

1. Preprocess inputs

```
$ cryodrgn downsample -h  
$ cryodrgn parse_ctf_star -h  
$ cryodrgn parse_pose_star -h
```

2. Training

```
$ cryodrgn train_vae -h
```

3. Analysis

```
# Analysis pipeline  
$ cryodrgn analyze -h
```

```
# Making movies  
$ cryodrgn pc_traversal -h  
$ cryodrgn graph_traversal -h
```

Tutorial

cryoDRGN EMPIAR-10076 tutorial

Preparing cryoDRGN inputs

- [Step 1\) Obtain the dataset](#)
- [Step 2\) Consensus reconstruction \(optional\)](#)
- [Step 3\) Preprocess inputs](#)
 - [Step 3.1\) Convert poses to cryoDRGN format](#)
 - [Step 3.2\) Convert CTF parameters to cryoDRGN format](#)
 - [Step 3.3\) Downsample images](#)

CryoDRGN training

- [General recommended workflow](#)
- [Step 4\) CryoDRGN initial training](#)
- [Extending or restarting from a checkpoint](#)

Overview of cryoDRGN analysis

- [Step 5\) cryodrgn analyze](#)
- [What's in the analysis directory?](#)
- [Visualization of the latent space](#)
- [Sampled density maps](#)
- [PC trajectories](#)

Step 6) Particle filtering with the cryoDRGN Jupyter notebook

- [Step 6.1\) Accessing the jupyter notebook](#)
- [Step 6.2\) Run the jupyter-notebook for particle filtering](#)
 - [Baseline: Published filtering results](#)
- [Step 6.3\) Filtering by GMM cluster label](#)
 - [Alternative method: Filtering by z-score](#)
 - [Alternative method: Filtering with an interactive lasso tool](#)
 - [View the raw particles](#)
 - [Saving the selection](#)
 - [\(Additional Functionality\) Writing a new .star file](#)
 - [\(Additional Functionality\) Extracting a new particle stack](#)

Step 7) CryoDRGN high resolution training

Enhancements and integrations

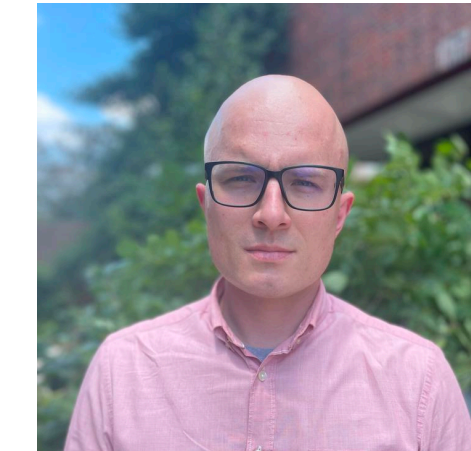
Suggestions wanted for cryoDRGN map and plot visualization in ChimeraX
#134

[Open](#) tomgoddard opened this issue 24 days ago · 6 comments

tomgoddard commented 24 days ago

I'm making a cryoDRGN visualization tool in ChimeraX and am interested in any suggestions users have about what it should do. So far it shows the umap plot and the maps computed by "cryodrgn analyze" on as points that plot and you can click on the points to see the map in the 3D view. You can cycle through the precomputed maps with a slider. I think it will be nice to allow computing new maps by clicking a point on the plot and morph between pairs of precomputed maps, and maybe make movies along paths drawn on the plot. Please add comments if you have other suggestions. Thanks!

The screenshot shows the ChimeraX interface. On the left, a 3D map of a protein complex is displayed in a grey, semi-transparent style. On the right, a UMAP plot is shown with points colored in a gradient from blue to red. The plot has several points labeled with numbers (1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20). The ChimeraX menu bar is visible at the top, and the cryoDRGN Viewer window title is 'cryoDRGN results directory ics/data/cryodrgn/analyze.41'.



Vineet Bansal
Michal Grzadkowski
Princeton Research
Computing

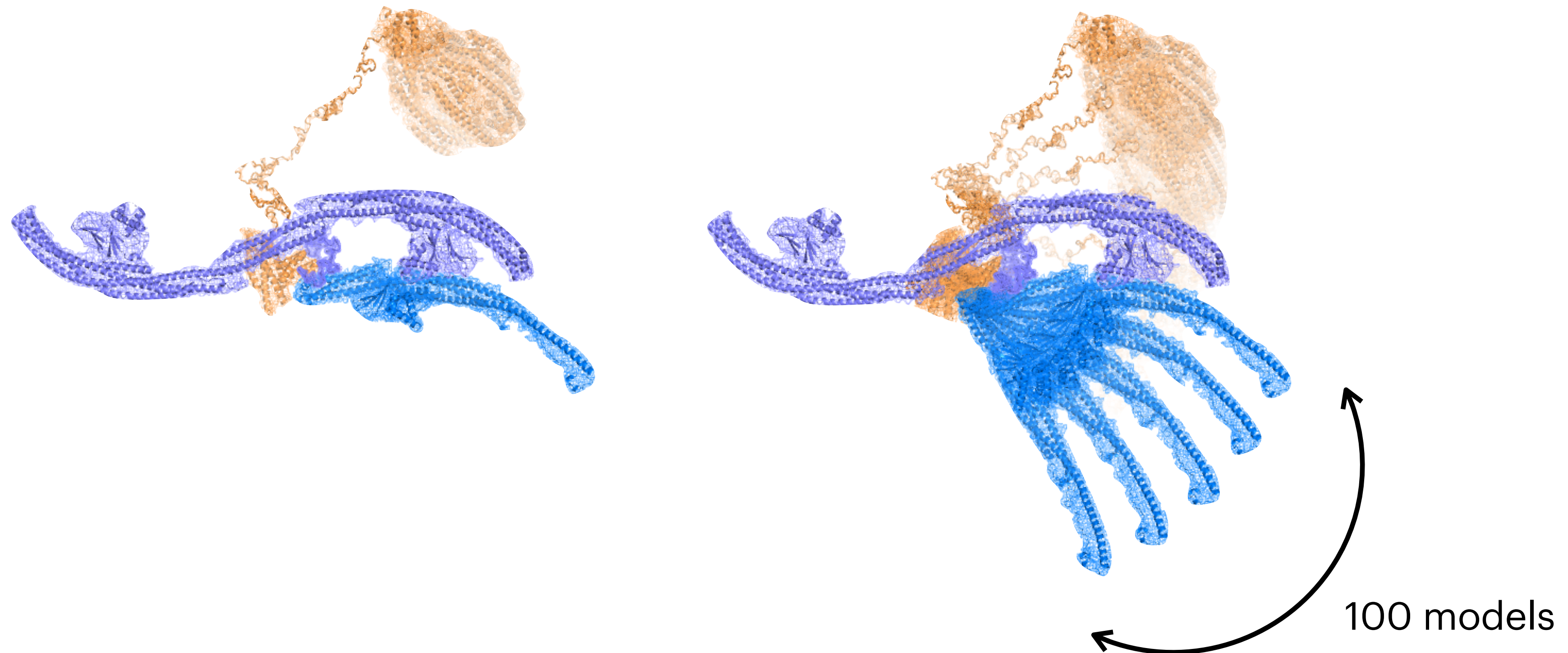
Now described in: Kinman, Powell, Zhong, Berger, Davis. Nature Protocols 2022.

Roadmap

- Motivation and background
- CryoDRGN: Deep Reconstructing Generative Networks
- **Validation on synthetic benchmarks**
- CryoDRGN reconstructions of real data
- Future vision

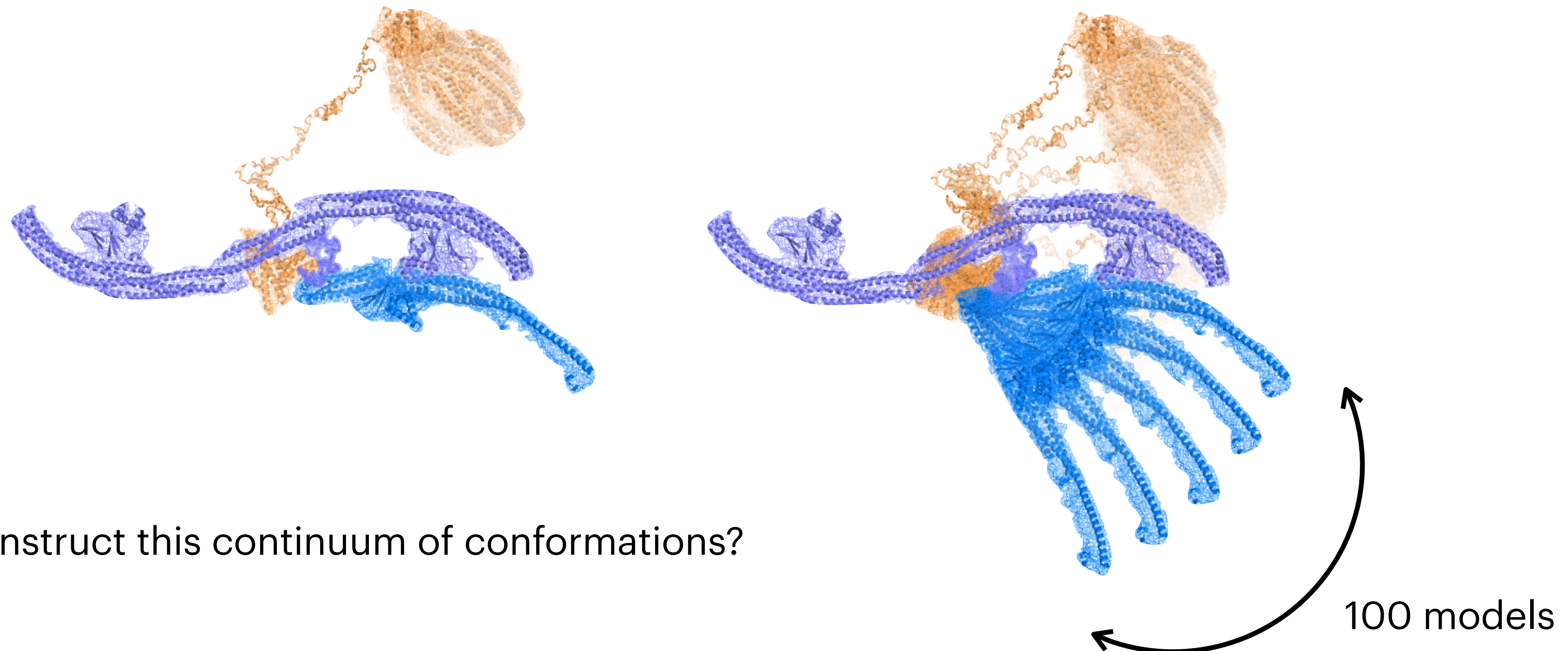
Heterogeneous reconstruction of a model protein complex containing 1 degree of freedom

- We generate a model protein complex containing one continuous degree of freedom
 - 100 atomic models varying one dihedral angle
 - 500 randomly oriented projections of each model, yielding a total of 50k projections

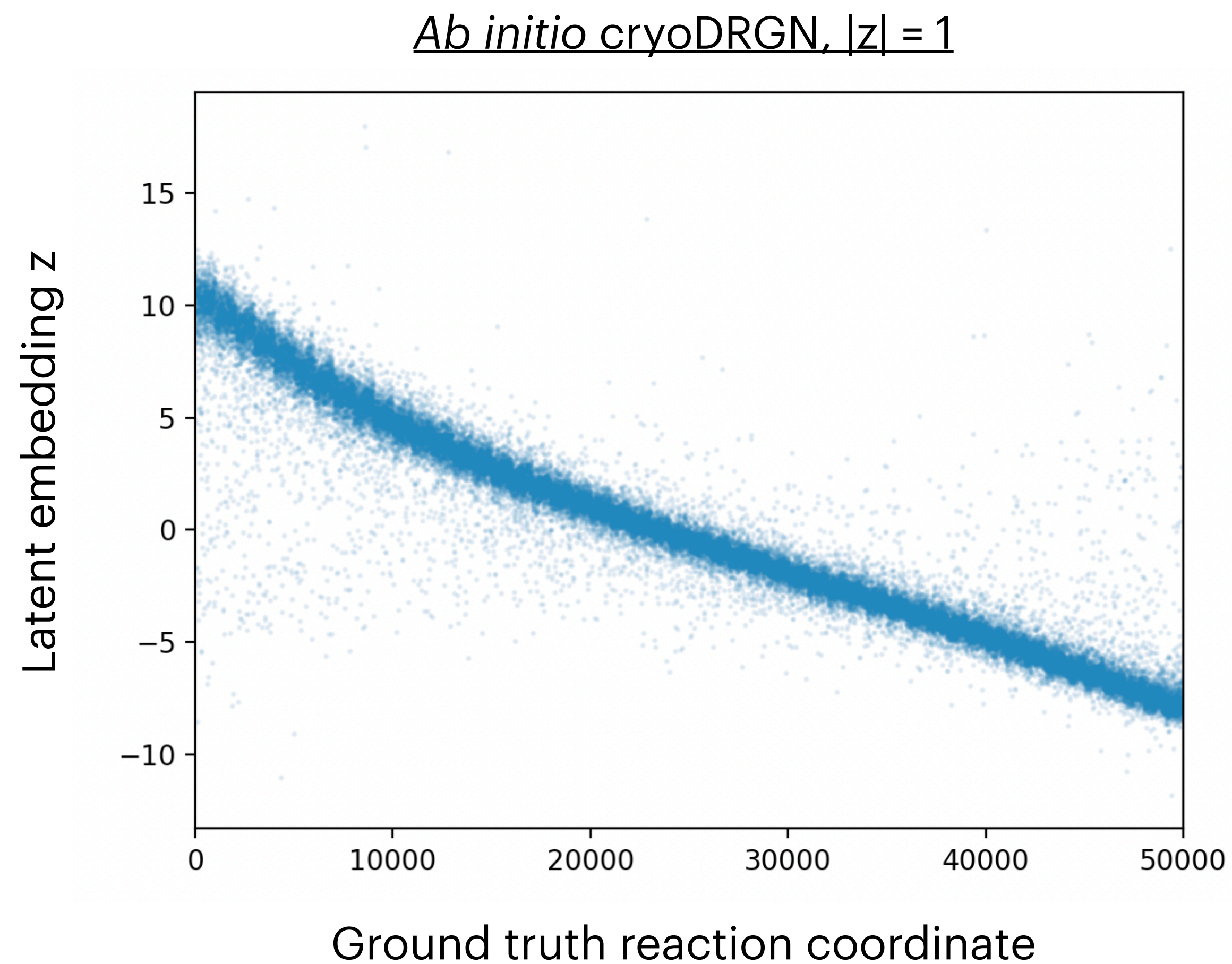


Heterogeneous reconstruction of a model protein complex containing 1 degree of freedom

- We generate a model protein complex containing one continuous degree of freedom
 - 100 atomic models varying one dihedral angle
 - 500 randomly oriented projections of each model, yielding a total of 50k projections



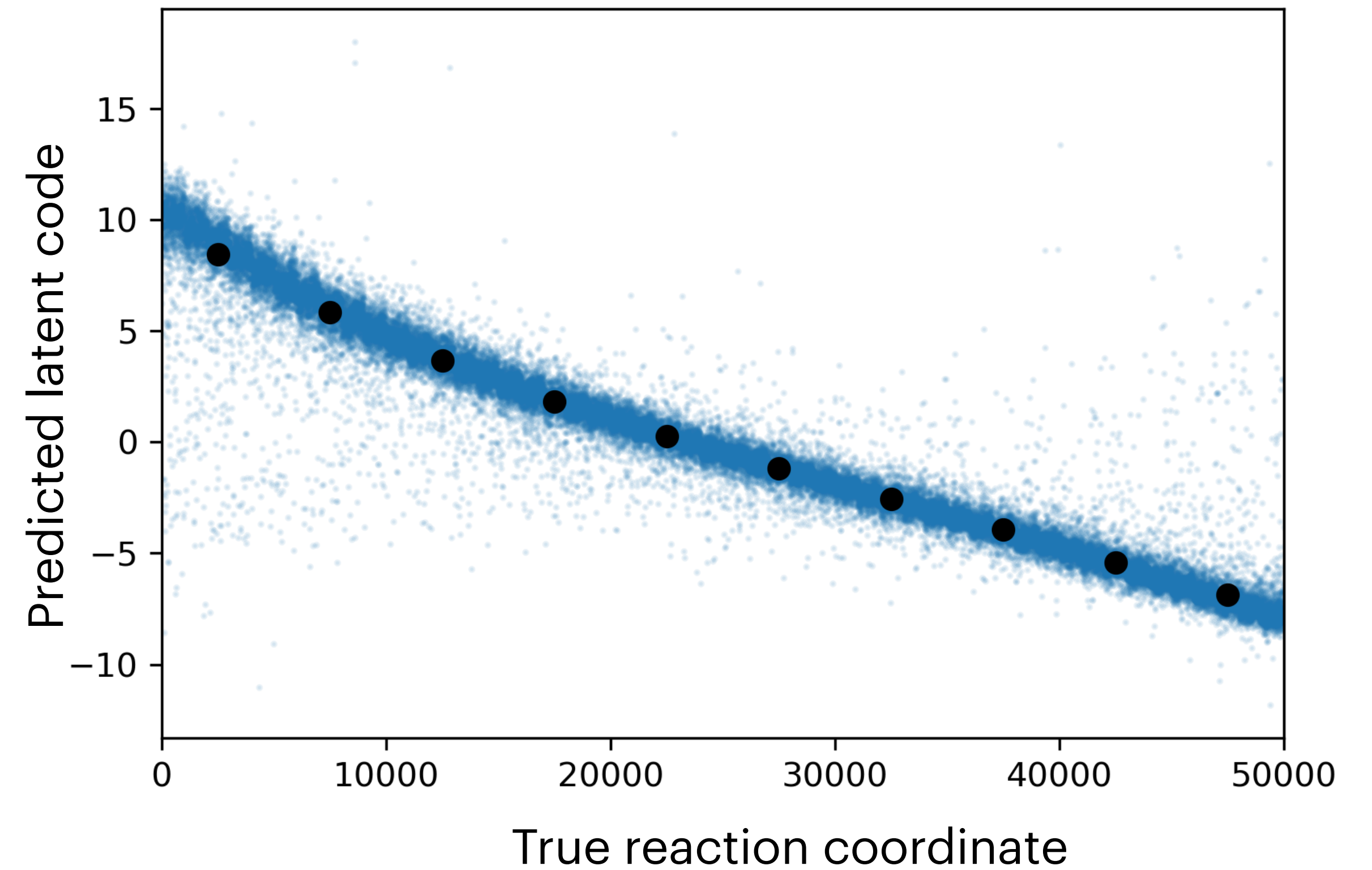
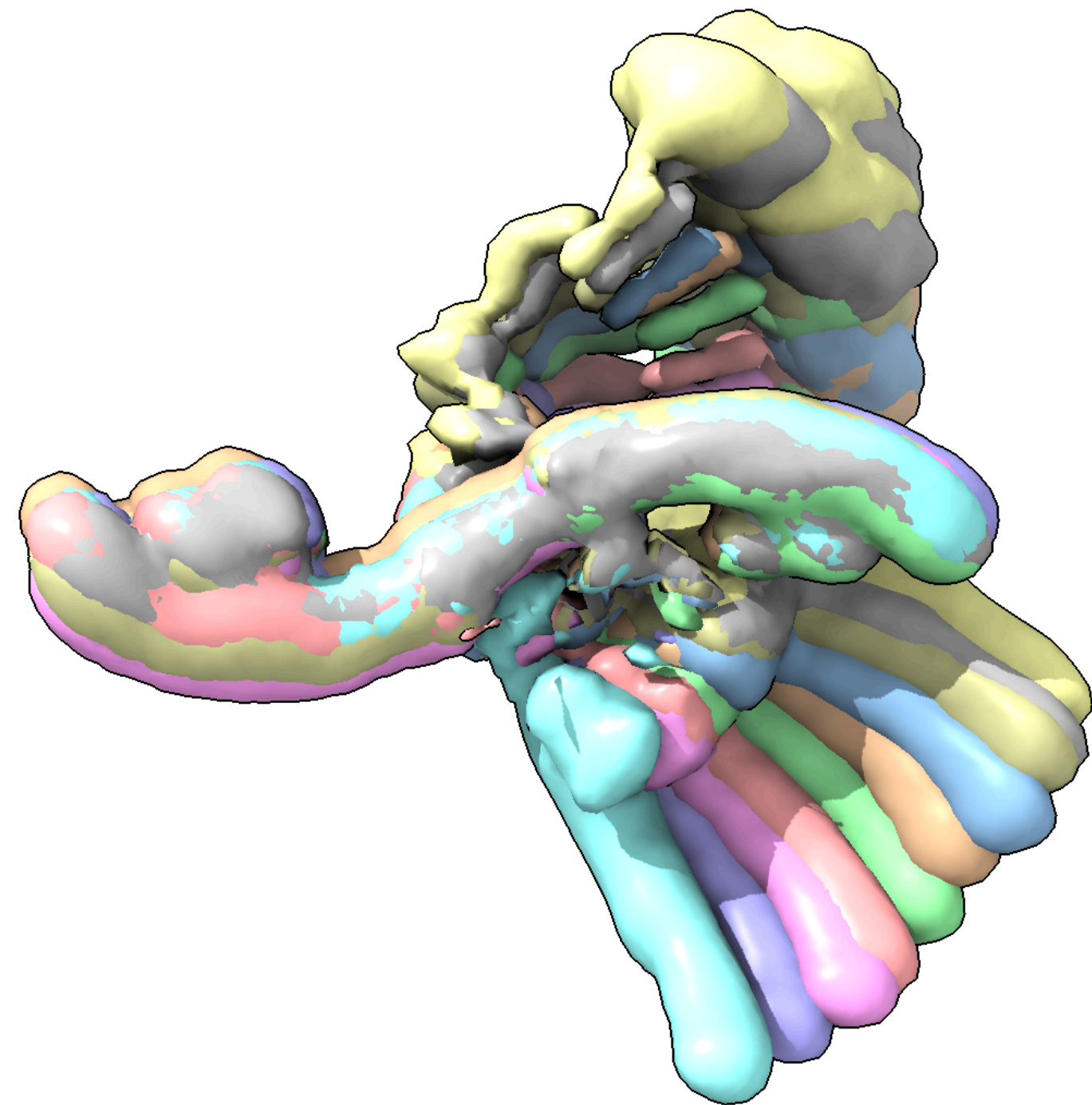
The predicted latent code correlates with the true reaction coordinate



CryoDRGN can reconstruct a continuum of structures along the true reaction coordinate

Ab initio cryoDRGN

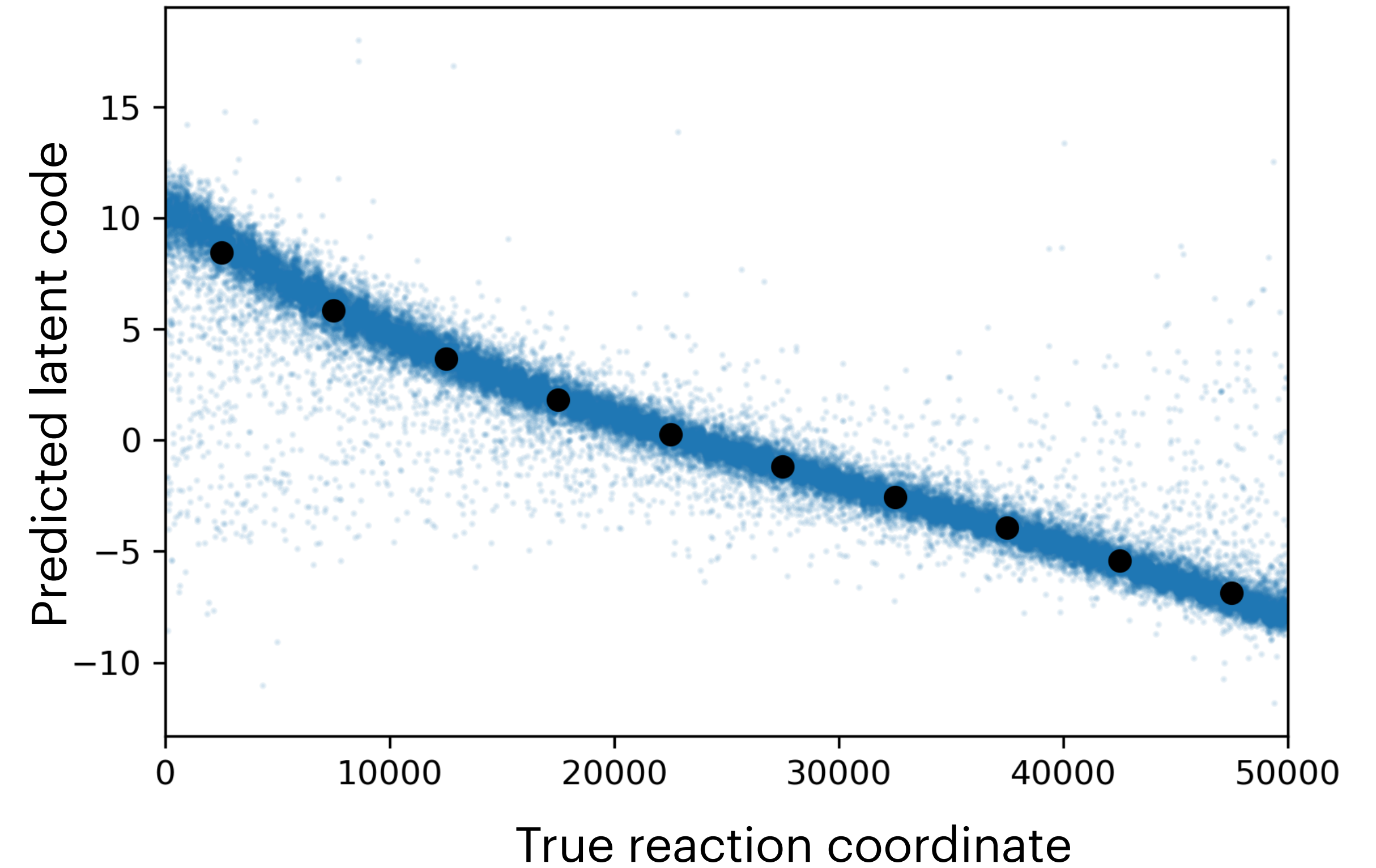
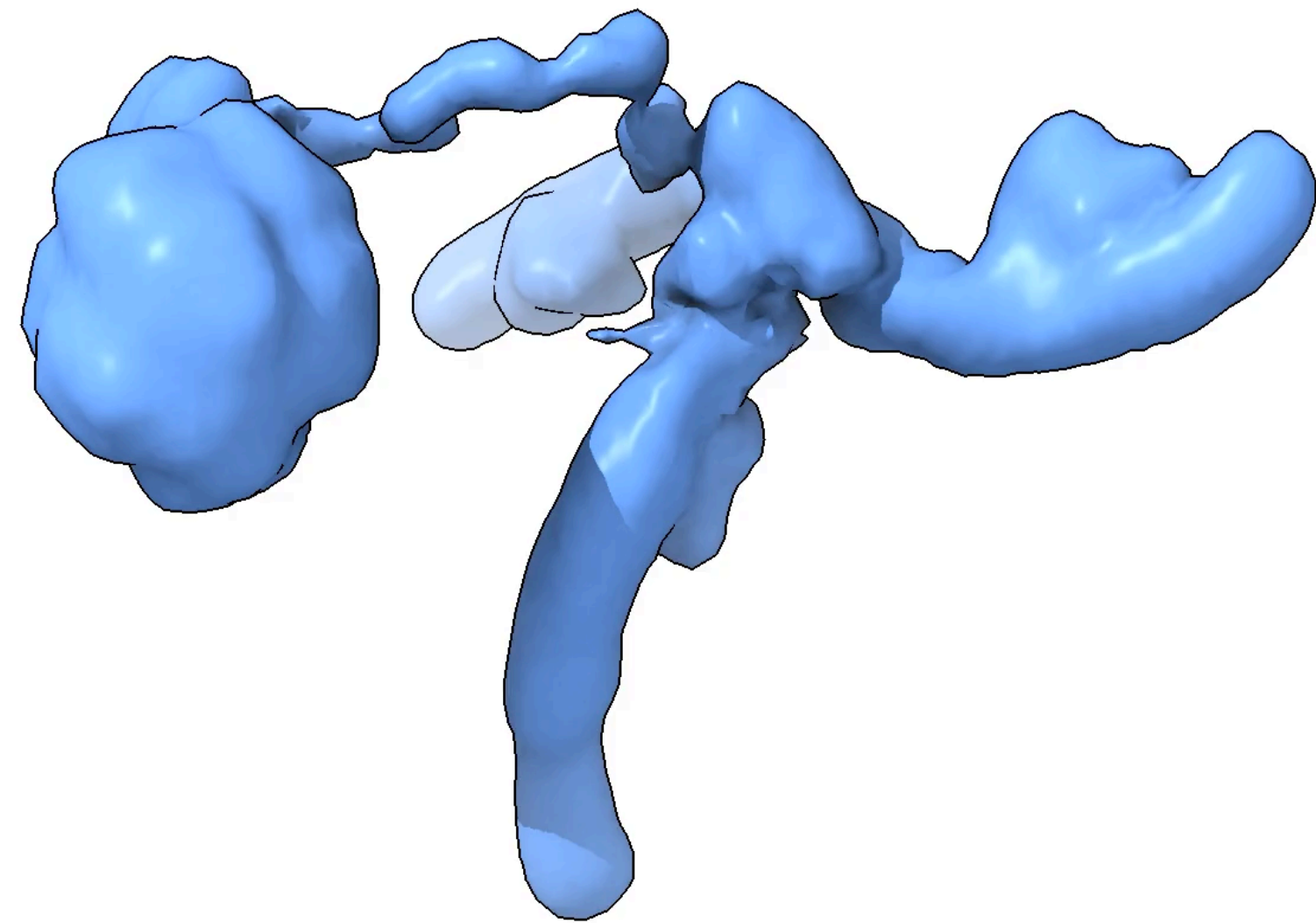
10 structures generated from latent representation



CryoDRGN can reconstruct a continuum of structures along the true reaction coordinate

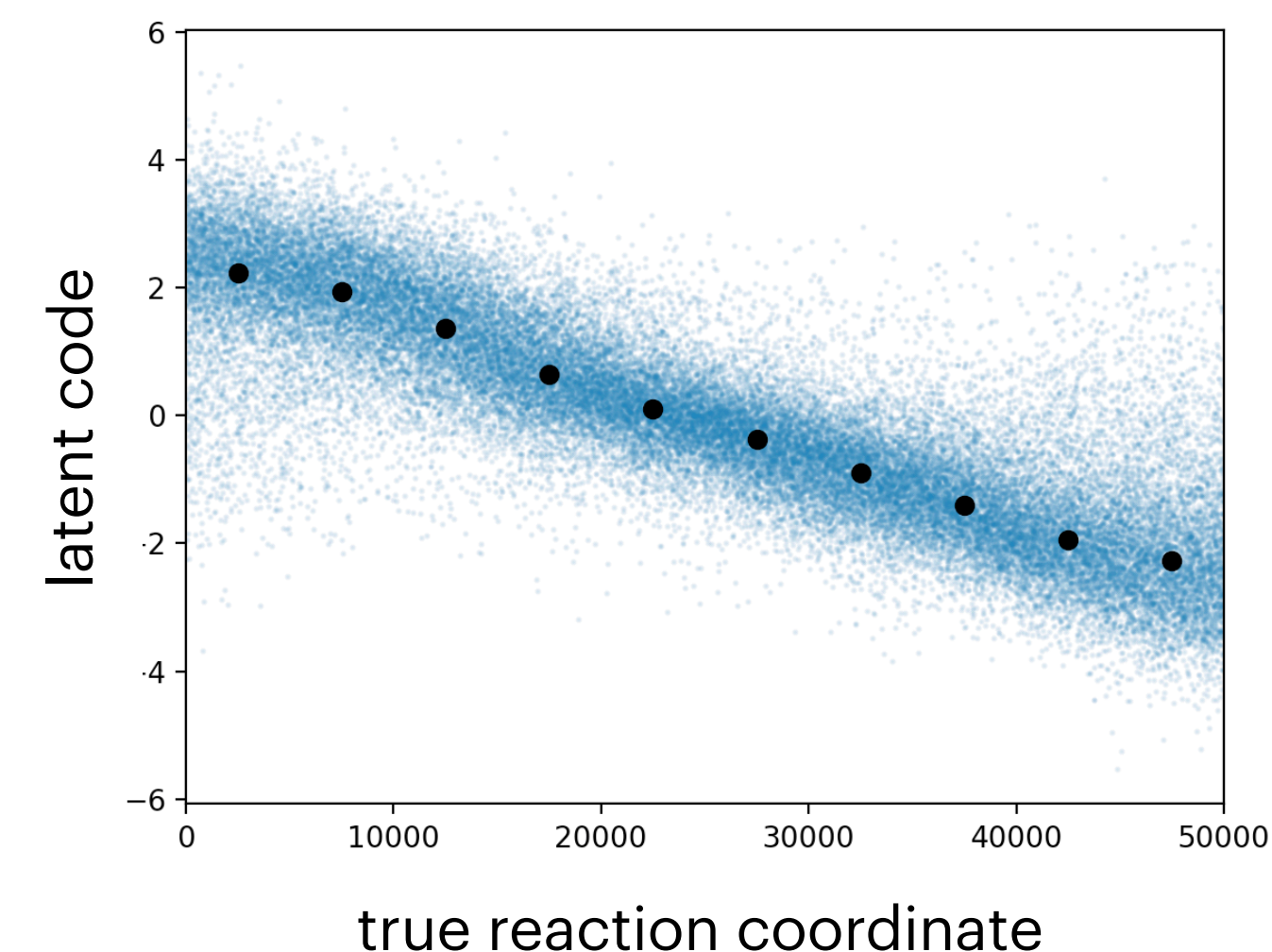
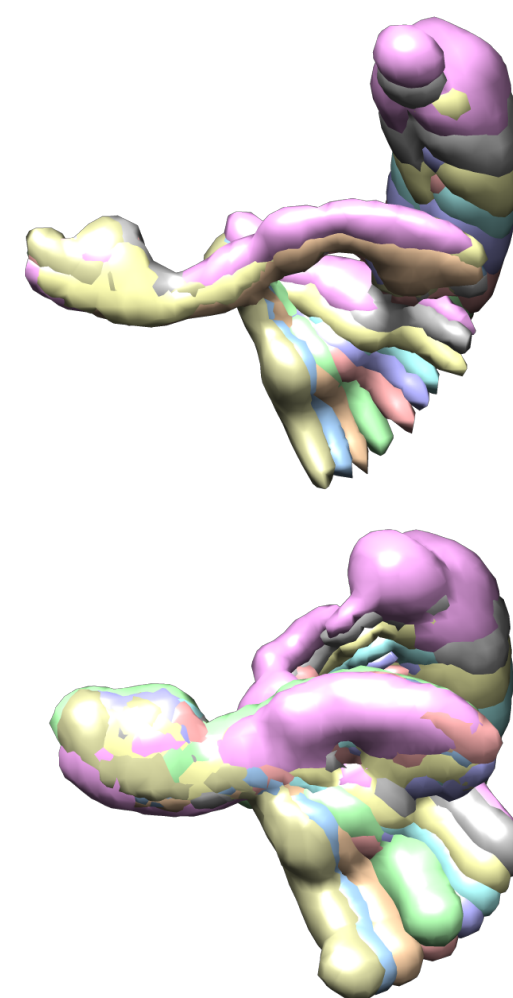
Ab initio cryoDRGN

10 structures generated from latent representation

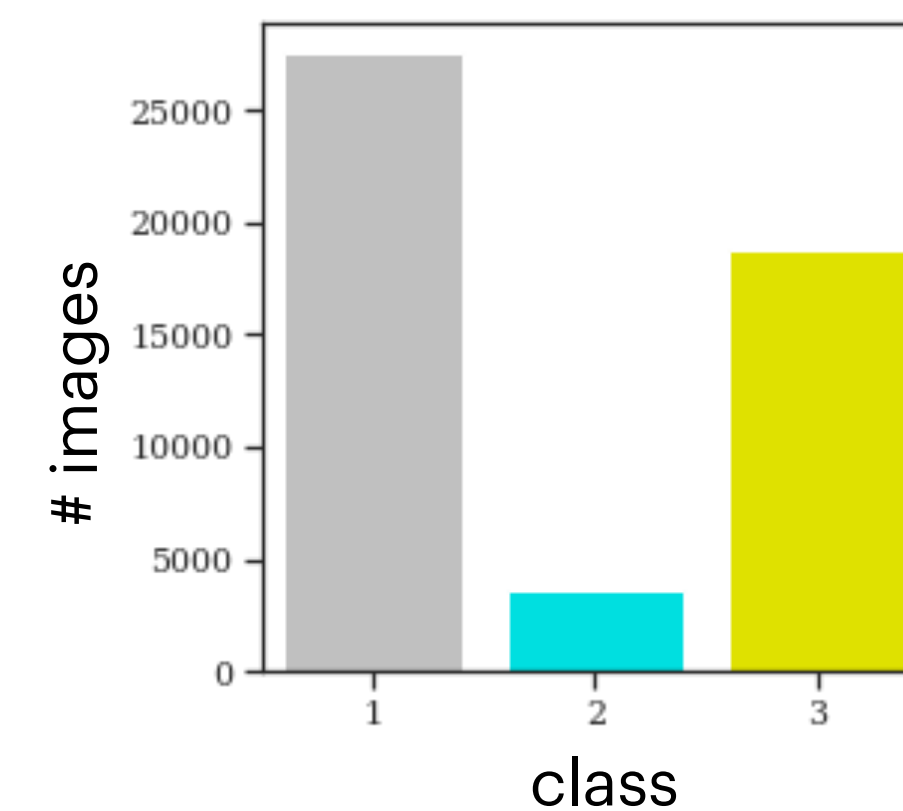
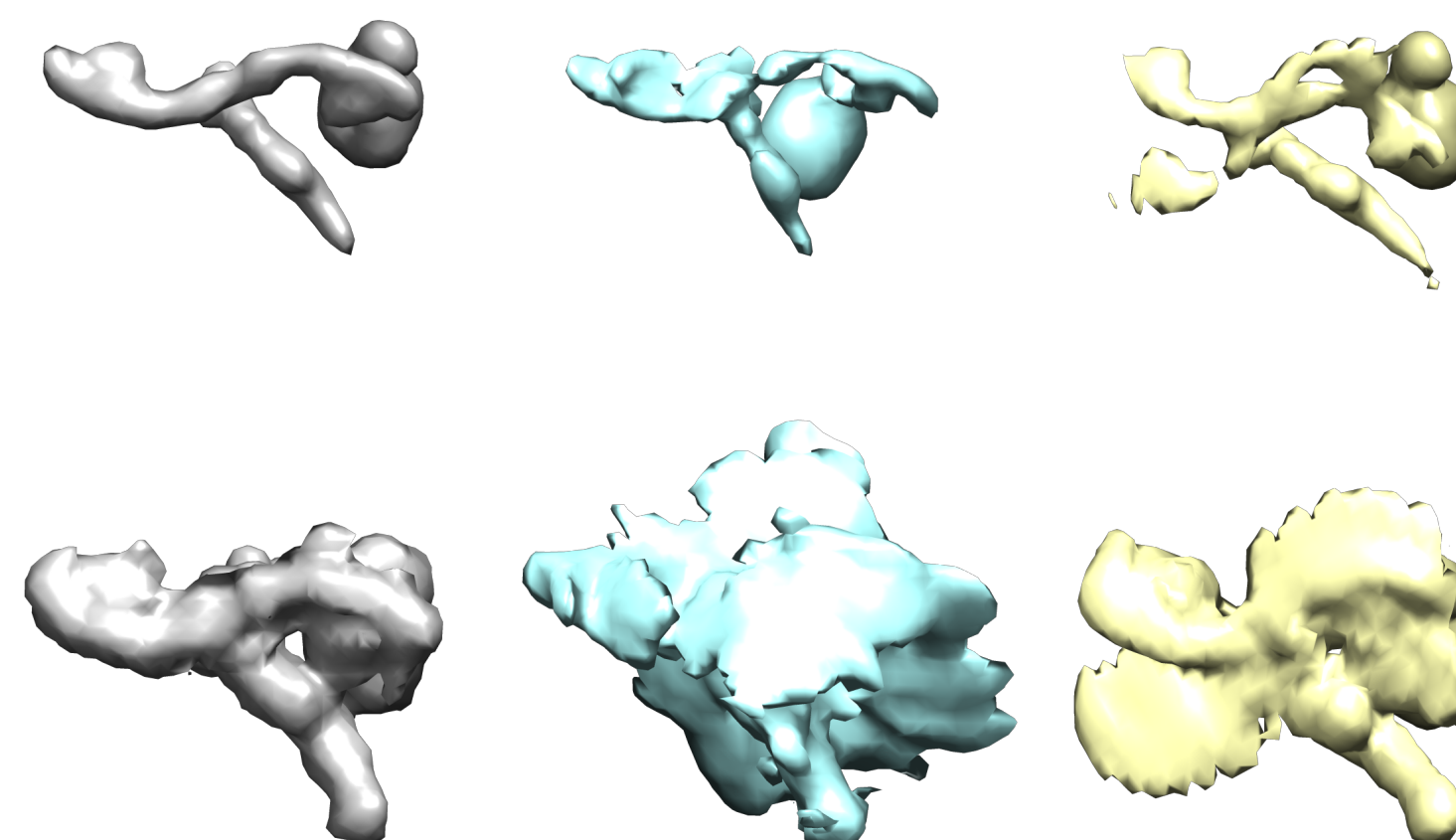


CryoDRGN can reconstruct a continuum of structures along the true reaction coordinate

Ab initio cryoDRGN
10 structures generated from latent representation



cryoSPARC discrete
multiclass reconstruction
K=3 classes

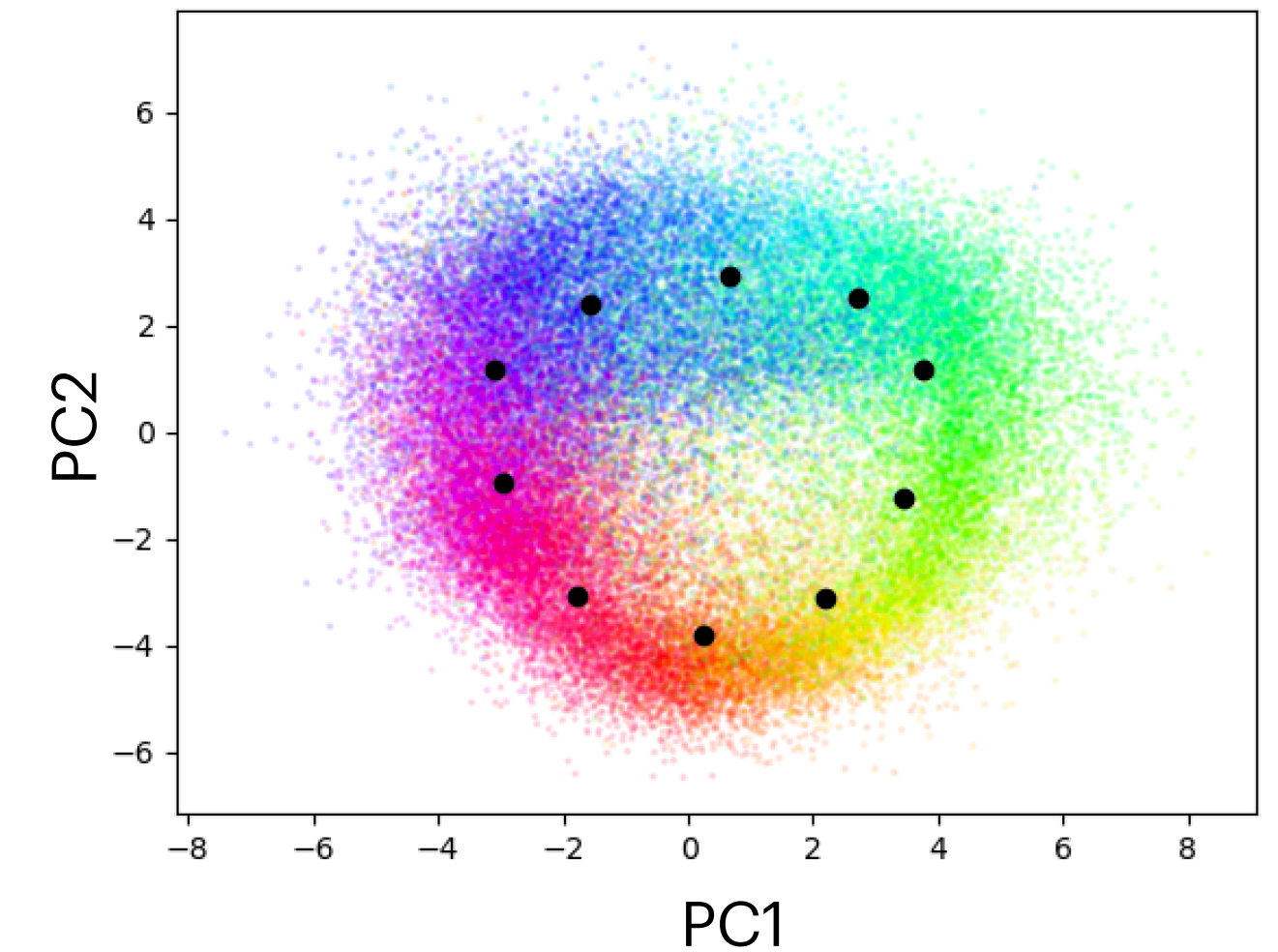
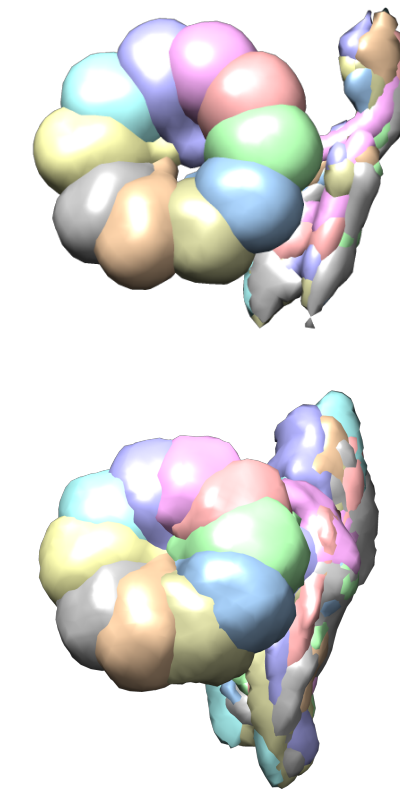


Heterogeneous reconstruction of a model protein complex with continuous motions

Ground truth

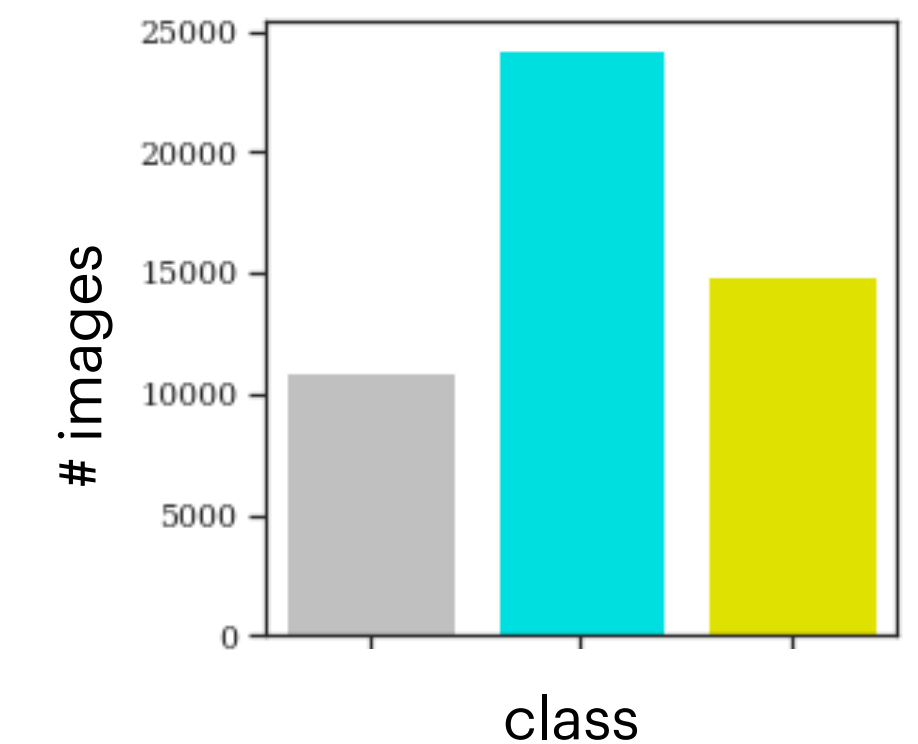
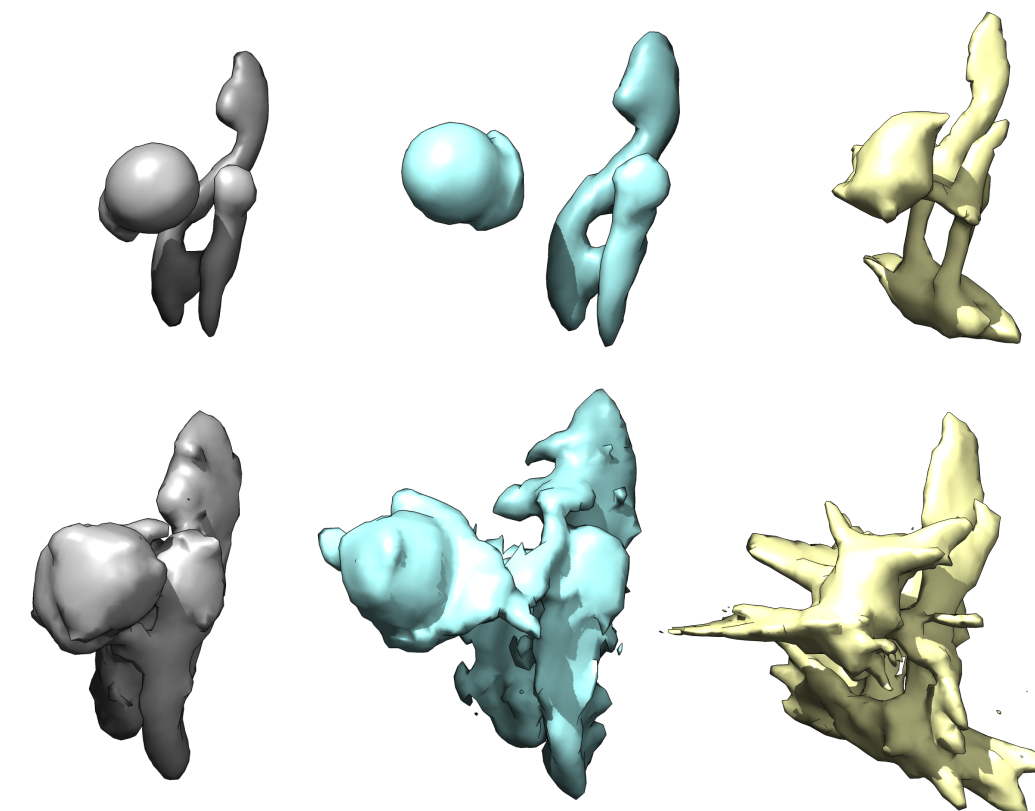


Ab initio cryoDRGN, $|z| = 10$
10 structures sampled
along latent space

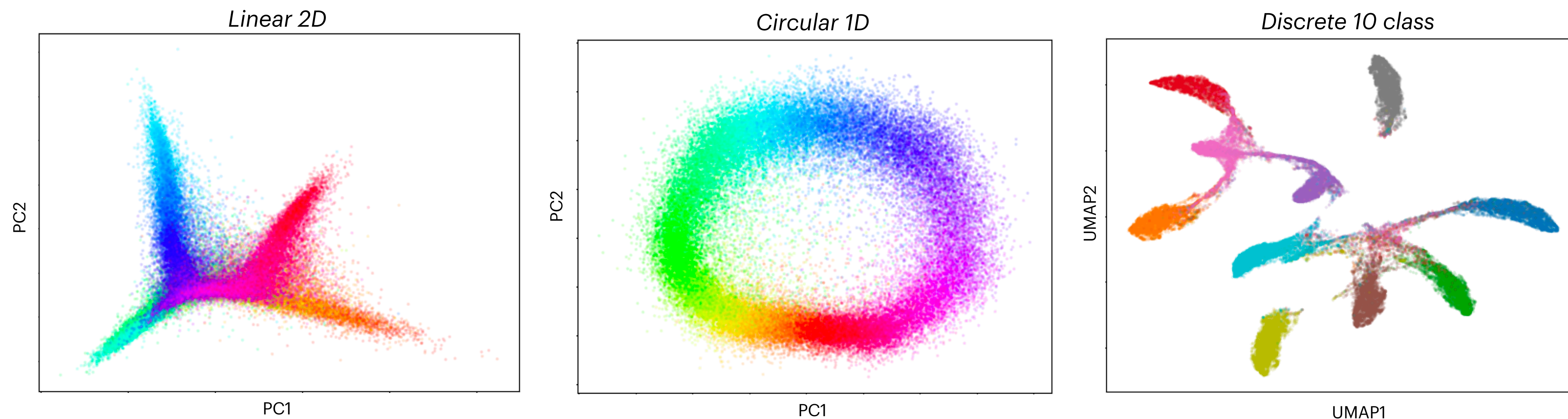


- 100 atomic models varying one dihedral angle
- 500 randomly oriented projections of each model
- 50k projection images

cryoSPARC discrete
multiclass reconstruction
K=3 classes



Additional datasets with more complex latent structure



Reconstruction accuracy quantified by an FSC=0.5 resolution metric between predicted and ground truth volume
(Lower is better; best possible is 2 pixels)

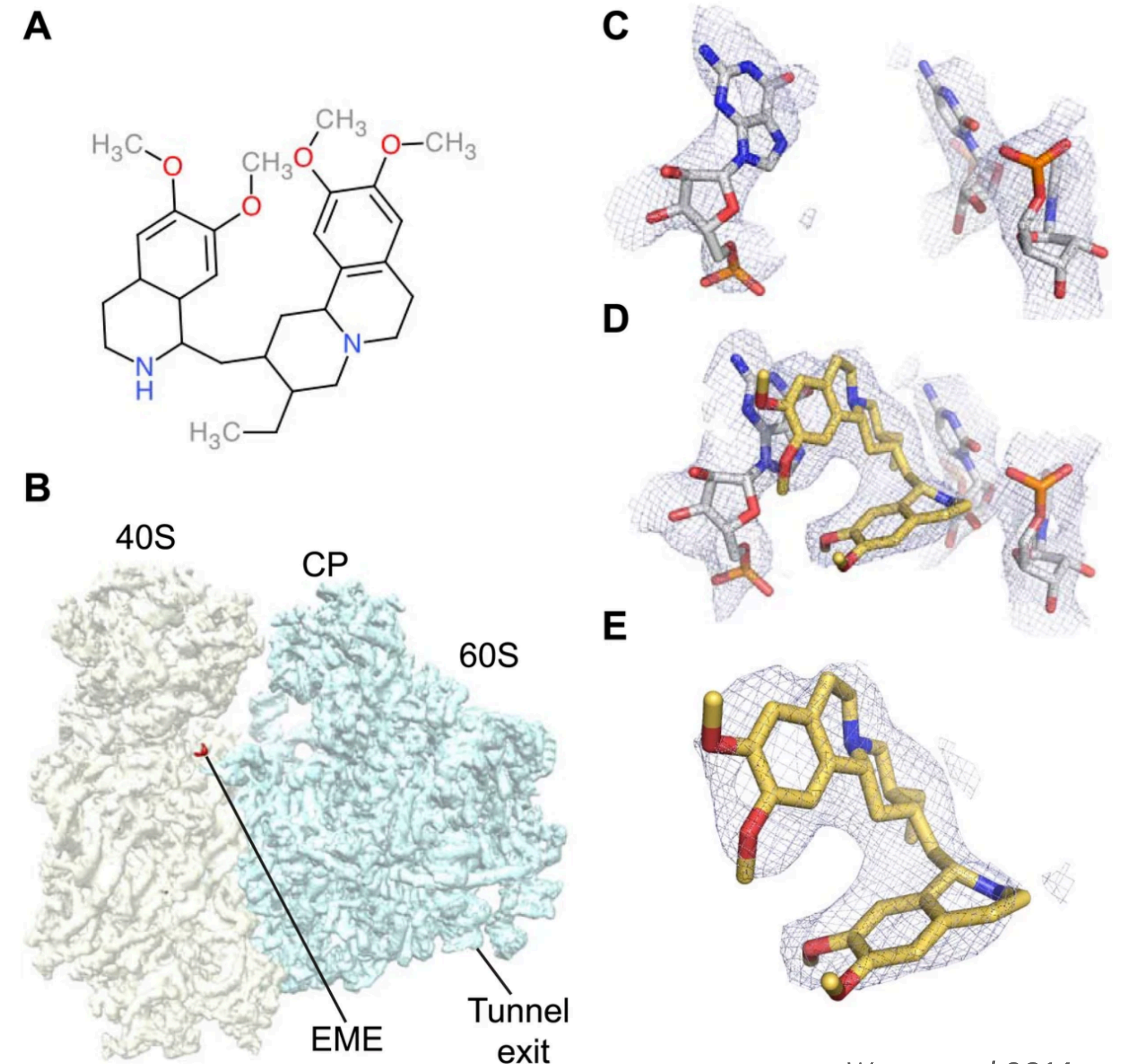
Dataset	cryoDRGN	cryoDRGN + tilt pairs	cryoSPARC
Linear 1D motion	2.50	2.35	3.60
Linear 2D motion	4.44	2.93	6.90
Circular 1D motion	3.86	2.63	4.87
Discrete 10 class	4.95	2.58	5.69

Roadmap

- Motivation and background
- CryoDRGN: Deep Reconstructing Generative Networks
- Validation on synthetic benchmarks
- **CryoDRGN reconstructions of real data**
 - **Uncovering residual heterogeneity in high resolution “homogeneous” datasets**
 - Discovering new states of the assembling ribosome
 - Reconstructing continuous motions of the pre-catalytic spliceosome
- Future vision

Homogeneous reconstruction of the *Pf*80S ribosome bound to the anti-protozoan drug emetine

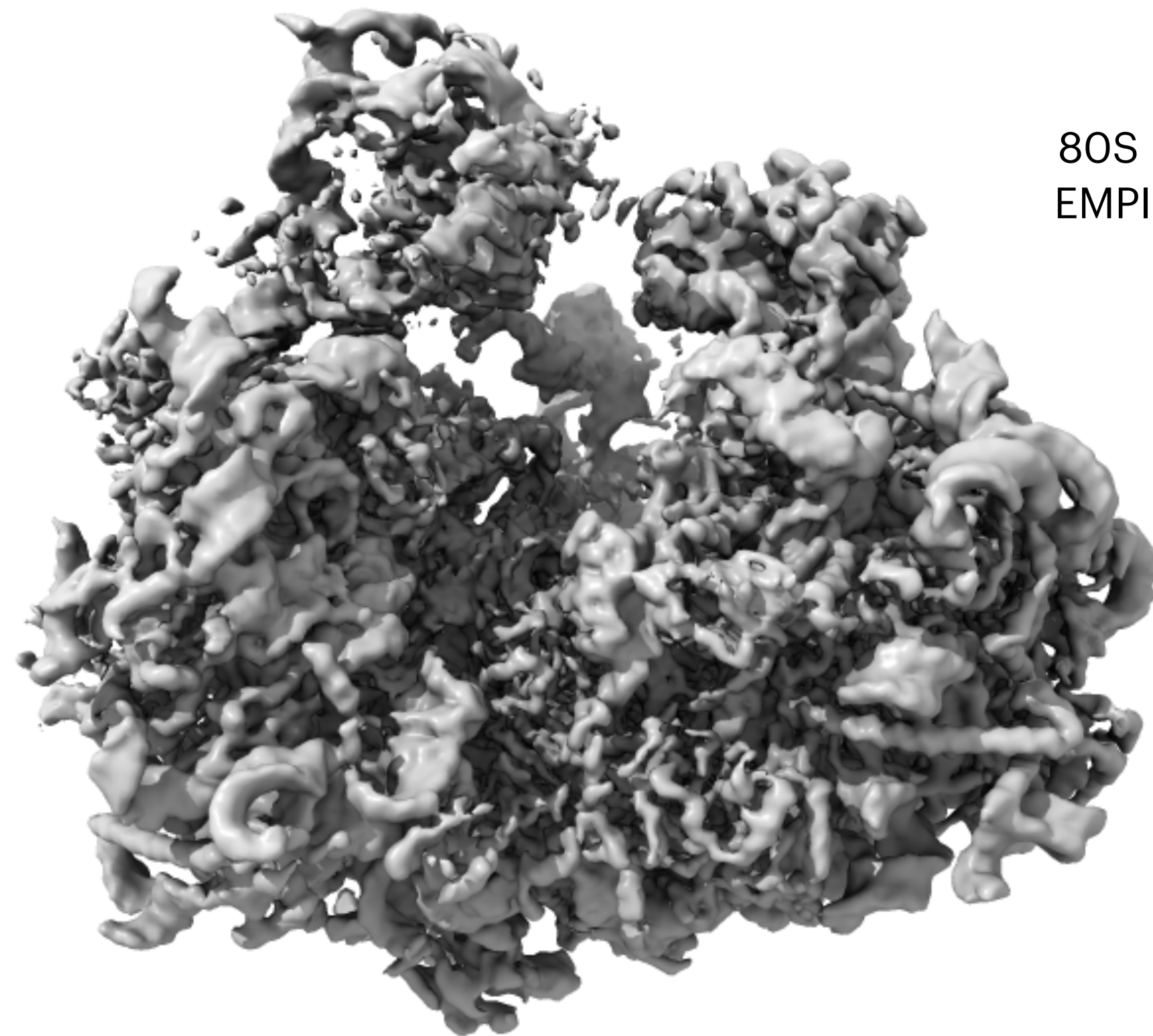
- 3.2 Å homogeneous reconstruction from 105k cryo-EM images (EMPIAR-10028)
- Difference map between structures with and without EME identified its binding site
- Lower local resolution in head group of small subunit and peripheral regions



CryoDRGN's neural model can learn high resolution cryo-EM density maps

- Train the cryoDRGN decoder (with no latent variable input) on images from EMPIAR-10028

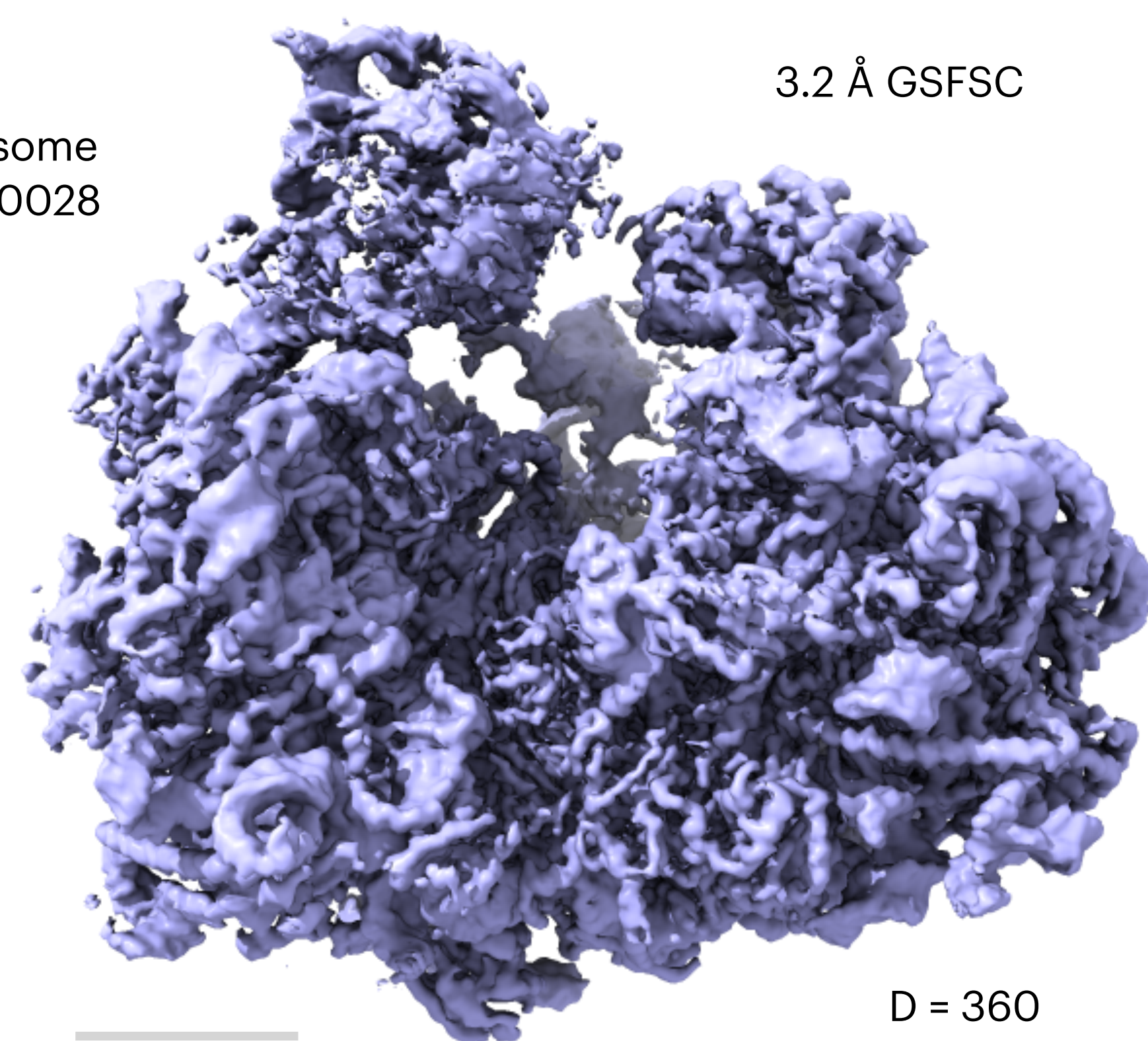
Neural network representation



1024x10 architecture
50 epochs

80S ribosome
EMPIAR 10028

Voxel-based representation



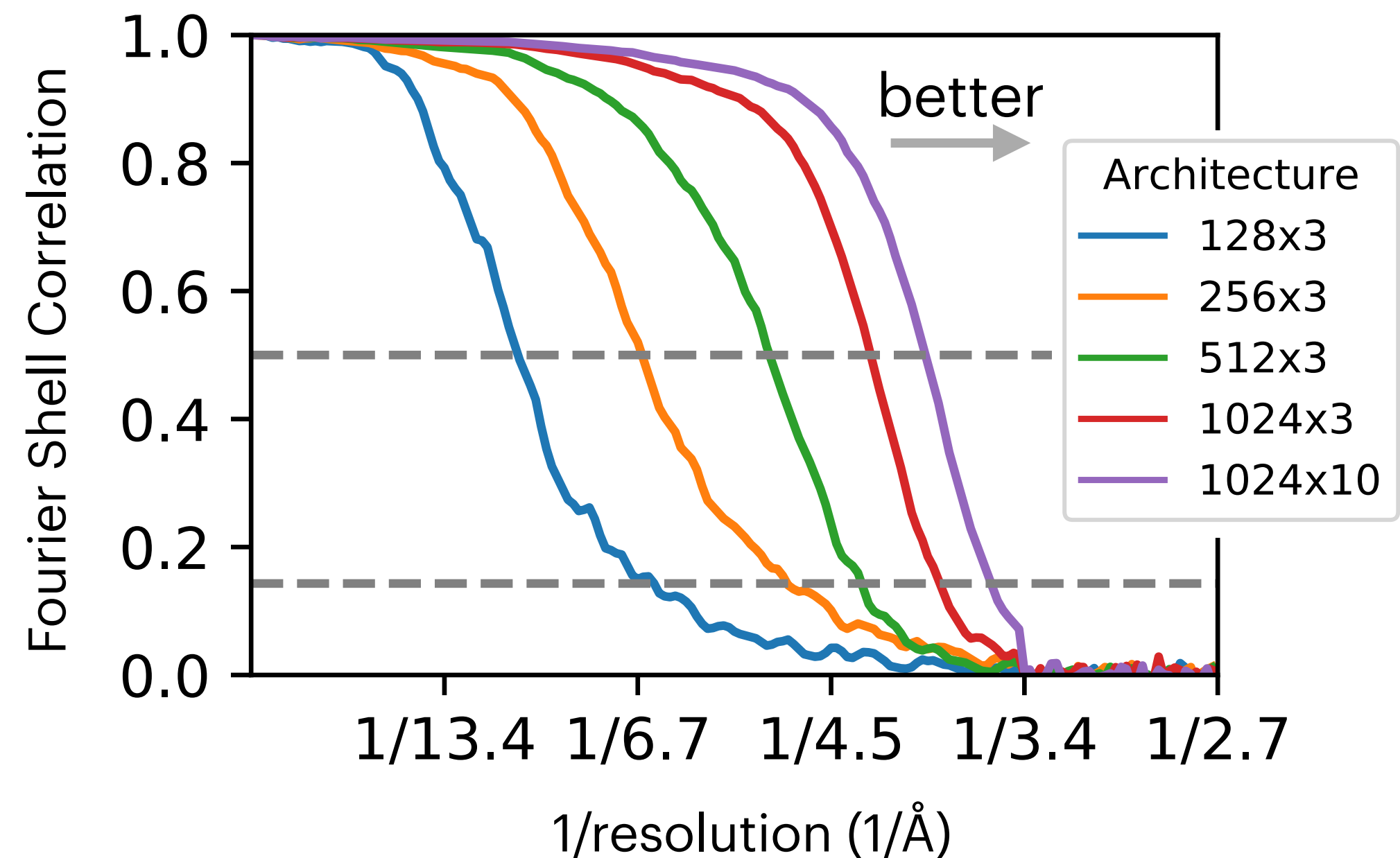
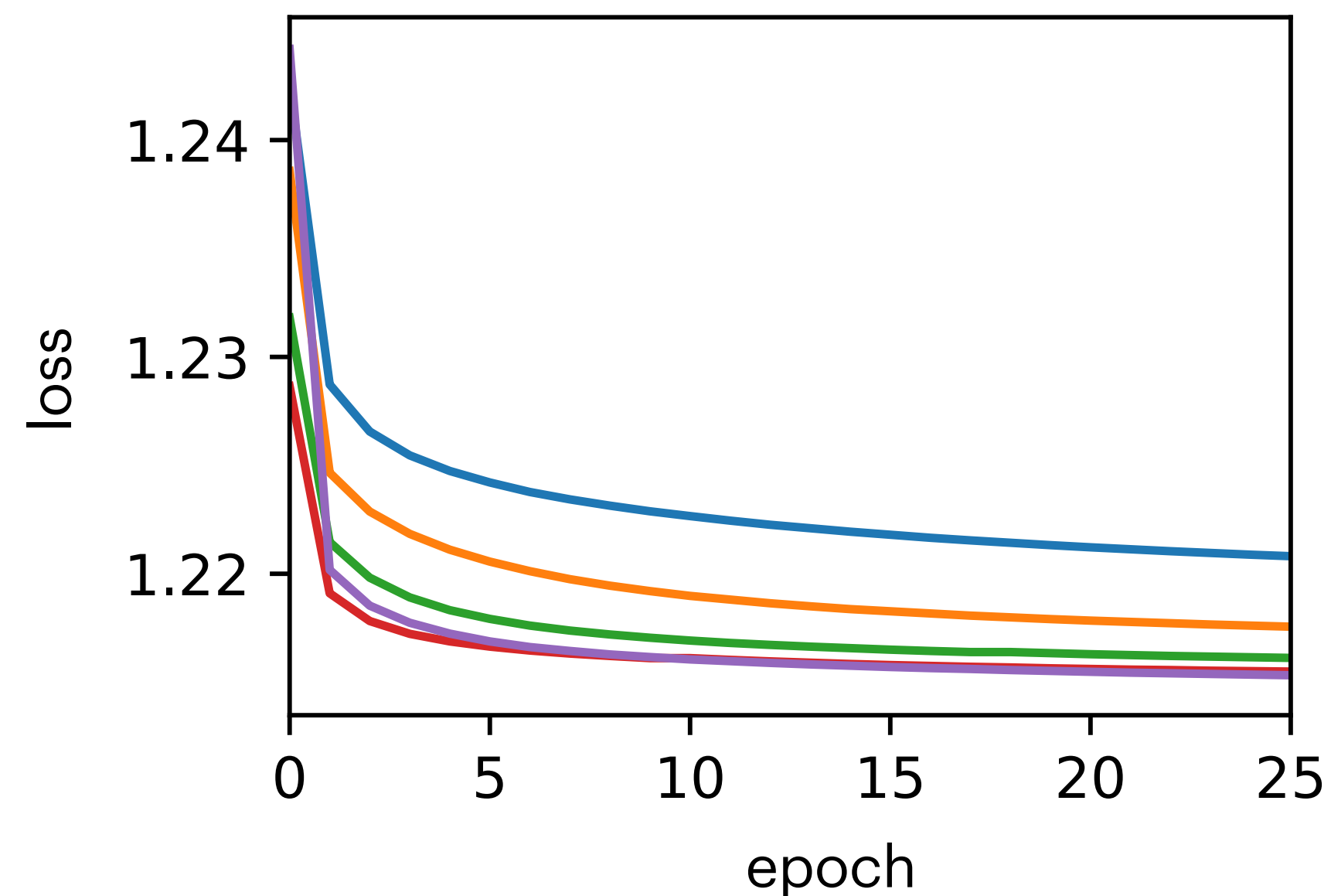
3.2 Å GSFSC

D = 360

50 Å

Achievable resolution is bounded by image size and model capacity

- The representation capacity of a cryoDRGN model is affected by:
 - Architecture
 - Fourier featurization, see Tancik et al. NeurIPS 2020
 - Latent variable dimension (for heterogeneous reconstruction)
- Inverse tradeoff between architecture size and training speed



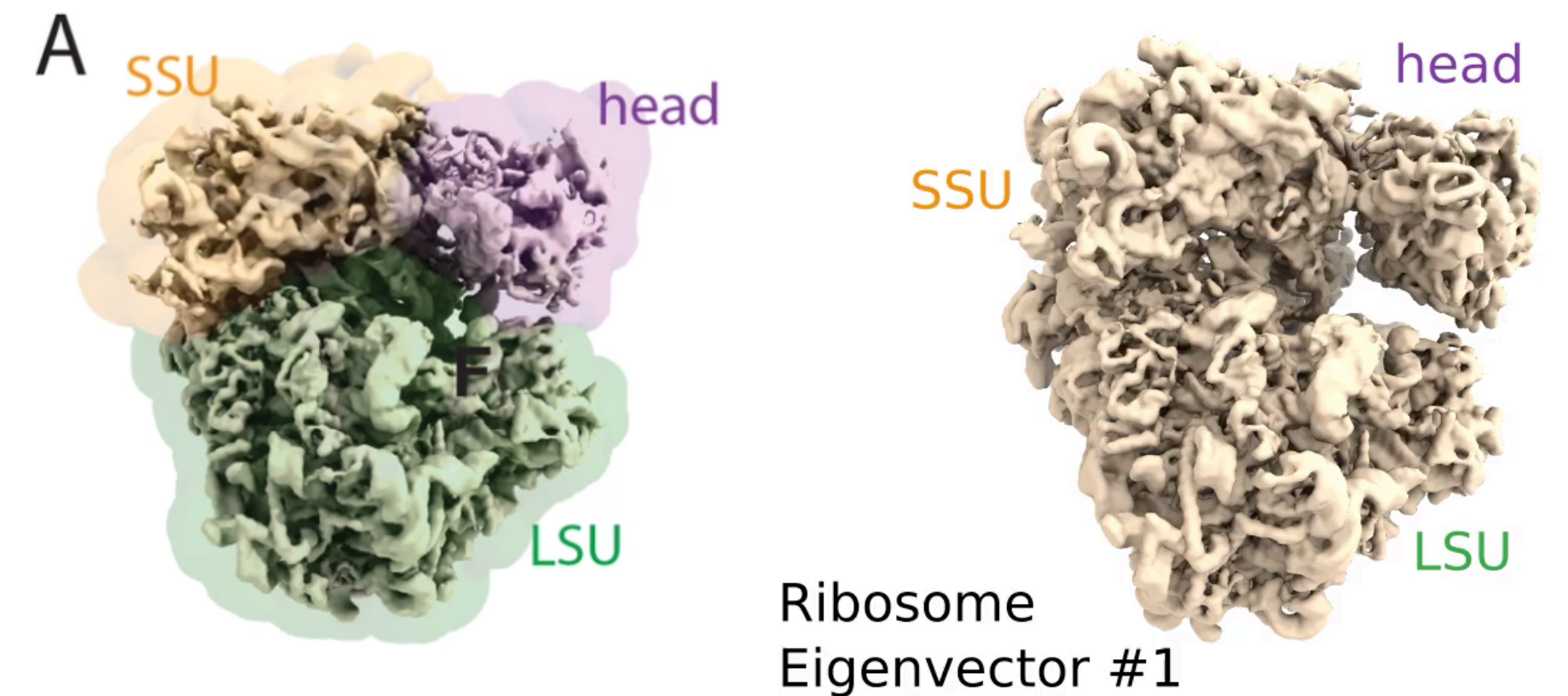
Advanced methods for heterogeneity analysis

- Multi-body analysis¹: Motions between B rigid bodies

$$Image X_i = CTF_i \left(\sum_{b=1}^B \mathbf{P}_{\phi_b} V_b \right) + N_i,$$

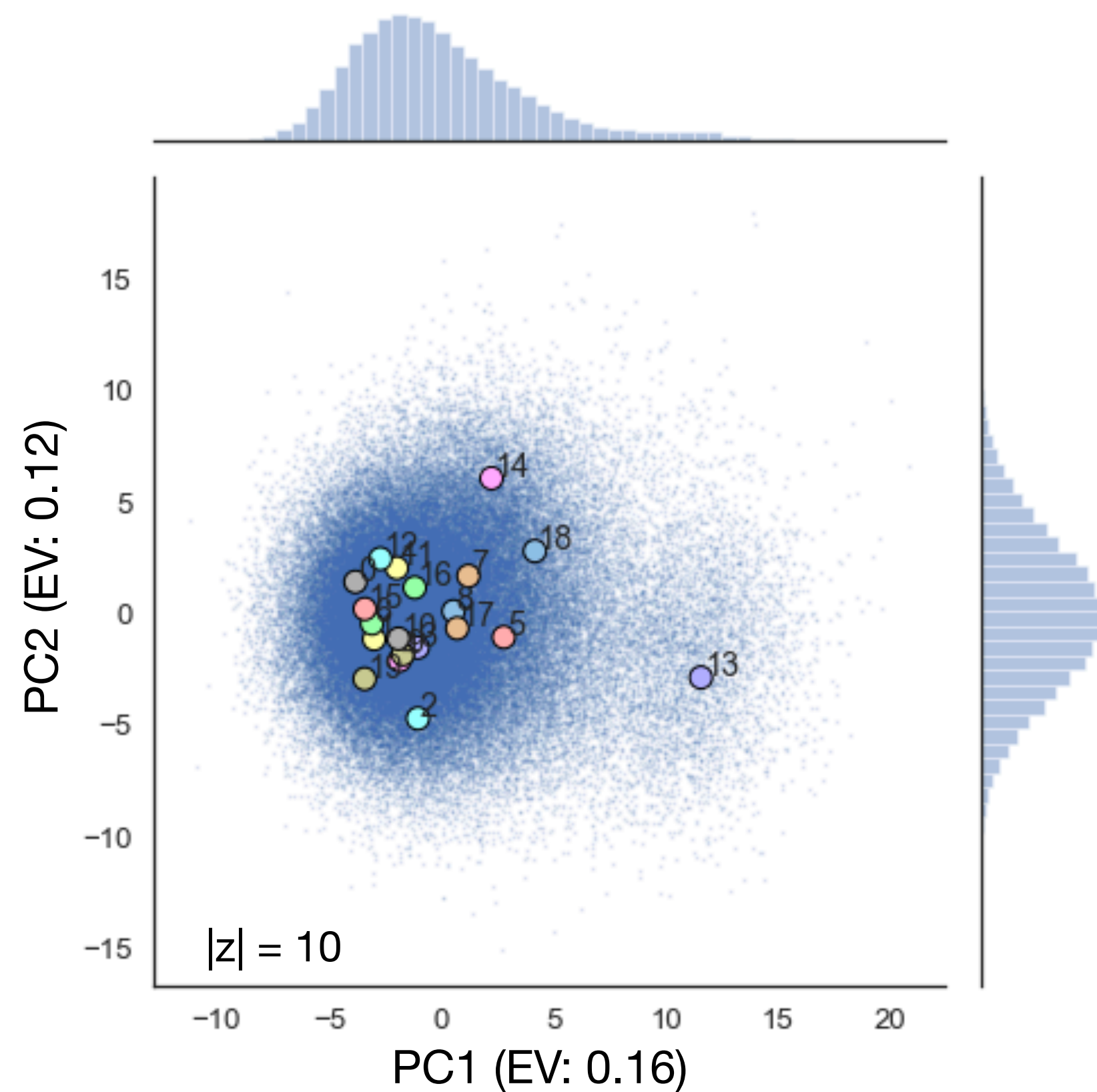
- cryoSPARC 3DVA²: Linear interpolations between “eigenvolumes”

$$\begin{aligned} Image X_i &= \alpha_i C_i P(\phi_i) \mathcal{V}(z_i) + \eta \\ &= \alpha_i C_i P(\phi_i) \left(V_0 + \sum_{k=1}^K z_{ik} V_k \right) + \eta \end{aligned}$$

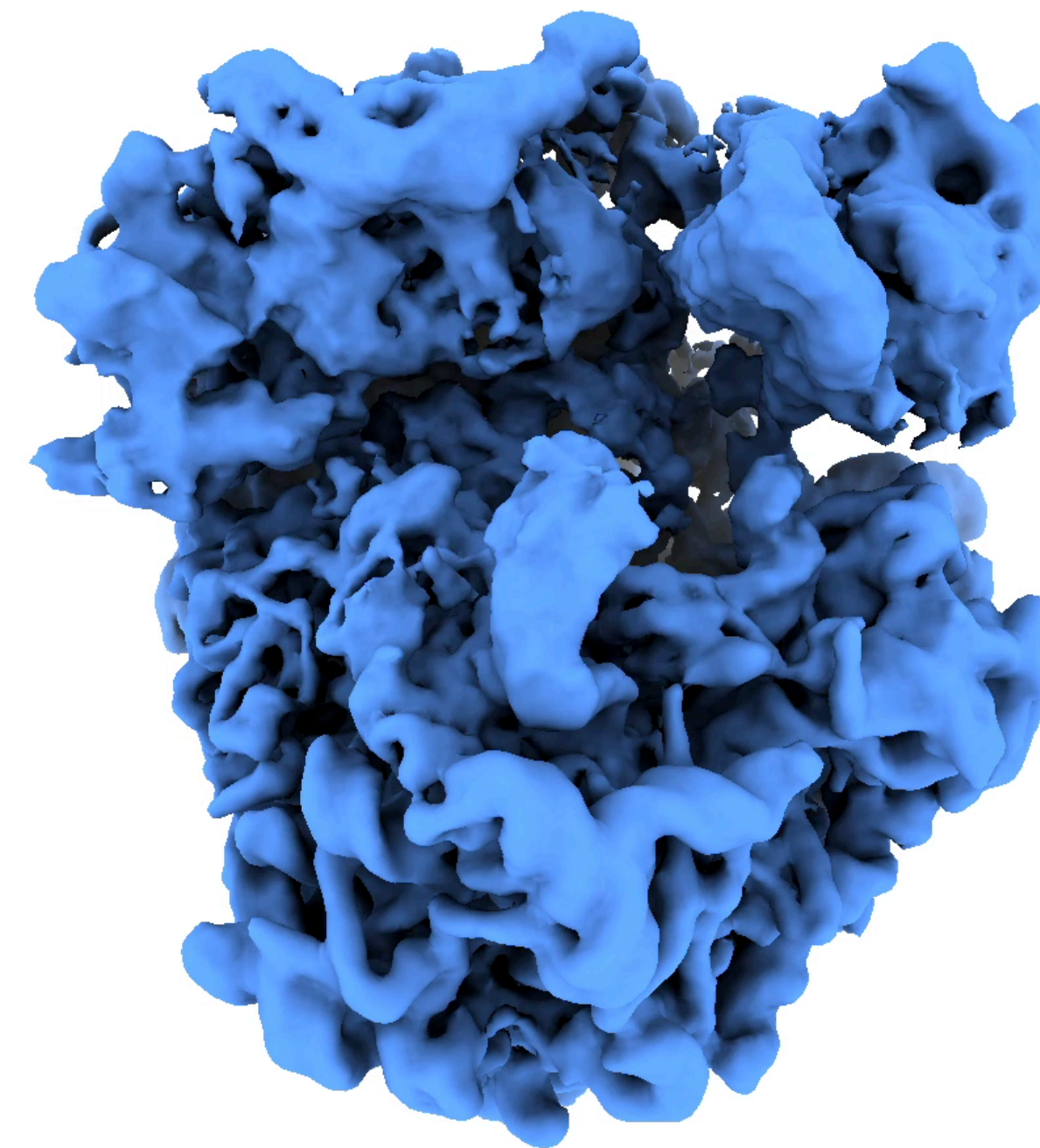


Discovering residual heterogeneity of the *Pf80S* ribosome [EMPIAR-10028]

CryoDRGN latent space:



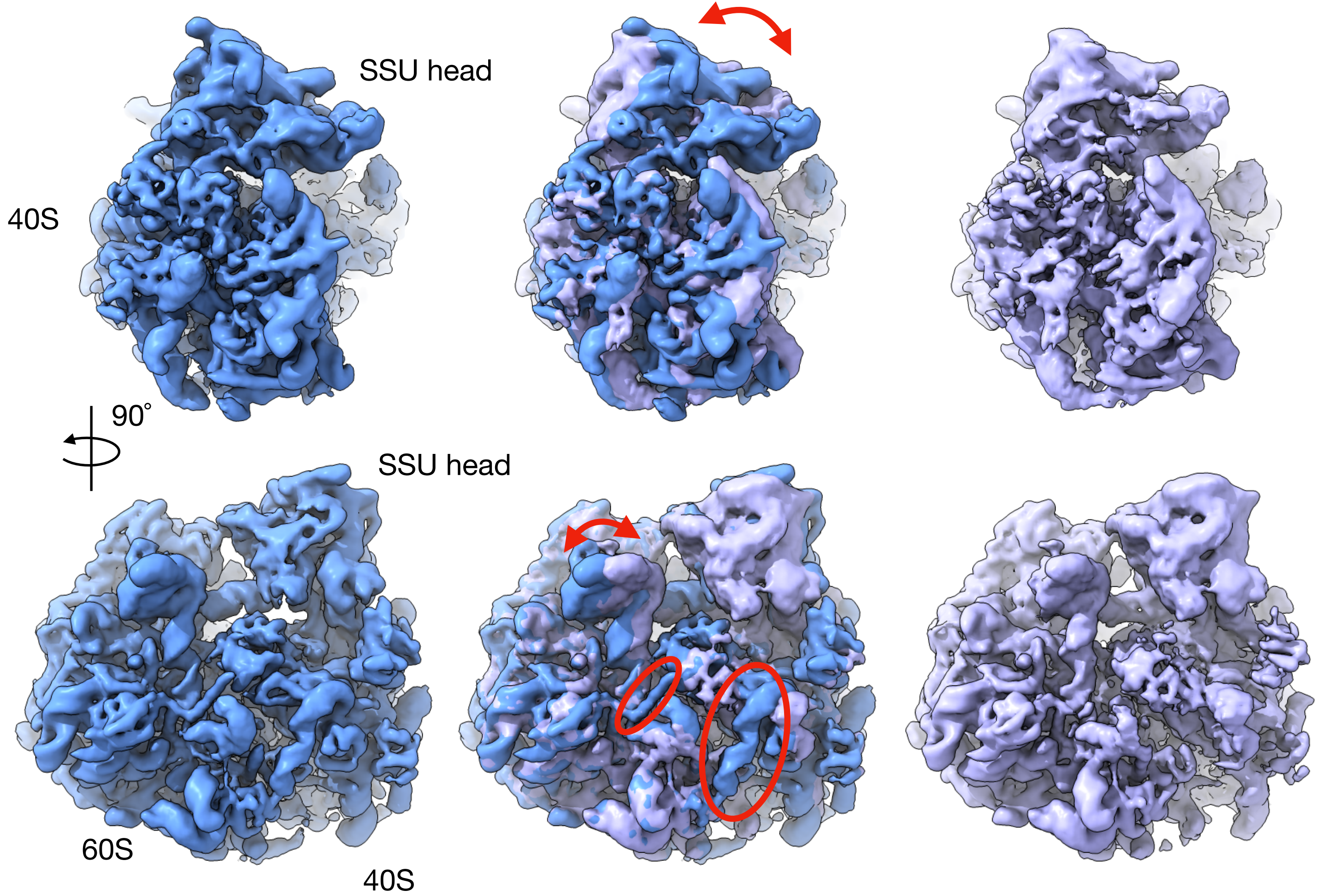
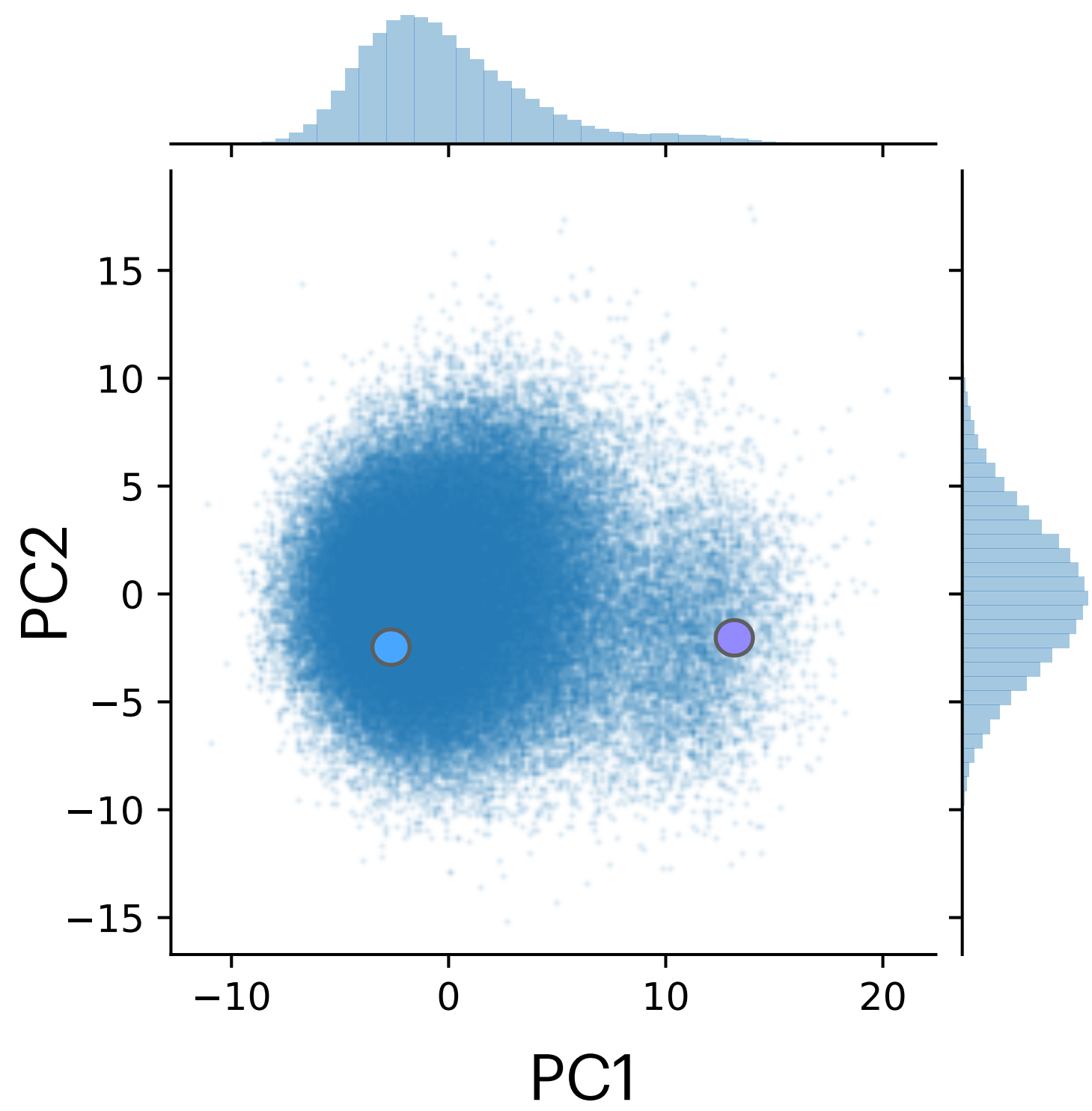
20 sampled structures:



- Subset of images separated by PC1 correspond to the 40S subunit in a rotated state
- Many heterogeneous elements in the large and small subunit

Variation in 40S SSU is consistent with other methods for heterogeneity analysis

CryoDRGN, comparison of 2 volumes



Roadmap

- Motivation and background
- CryoDRGN: Deep Reconstructing Generative Networks
- Validation on synthetic benchmarks
- **CryoDRGN reconstructions of real data**
 - Uncovering residual heterogeneity in high resolution “homogeneous” datasets
 - **Discovering new states of the assembling ribosome**
 - Reconstructing continuous motions of the pre-catalytic spliceosome
- Future vision

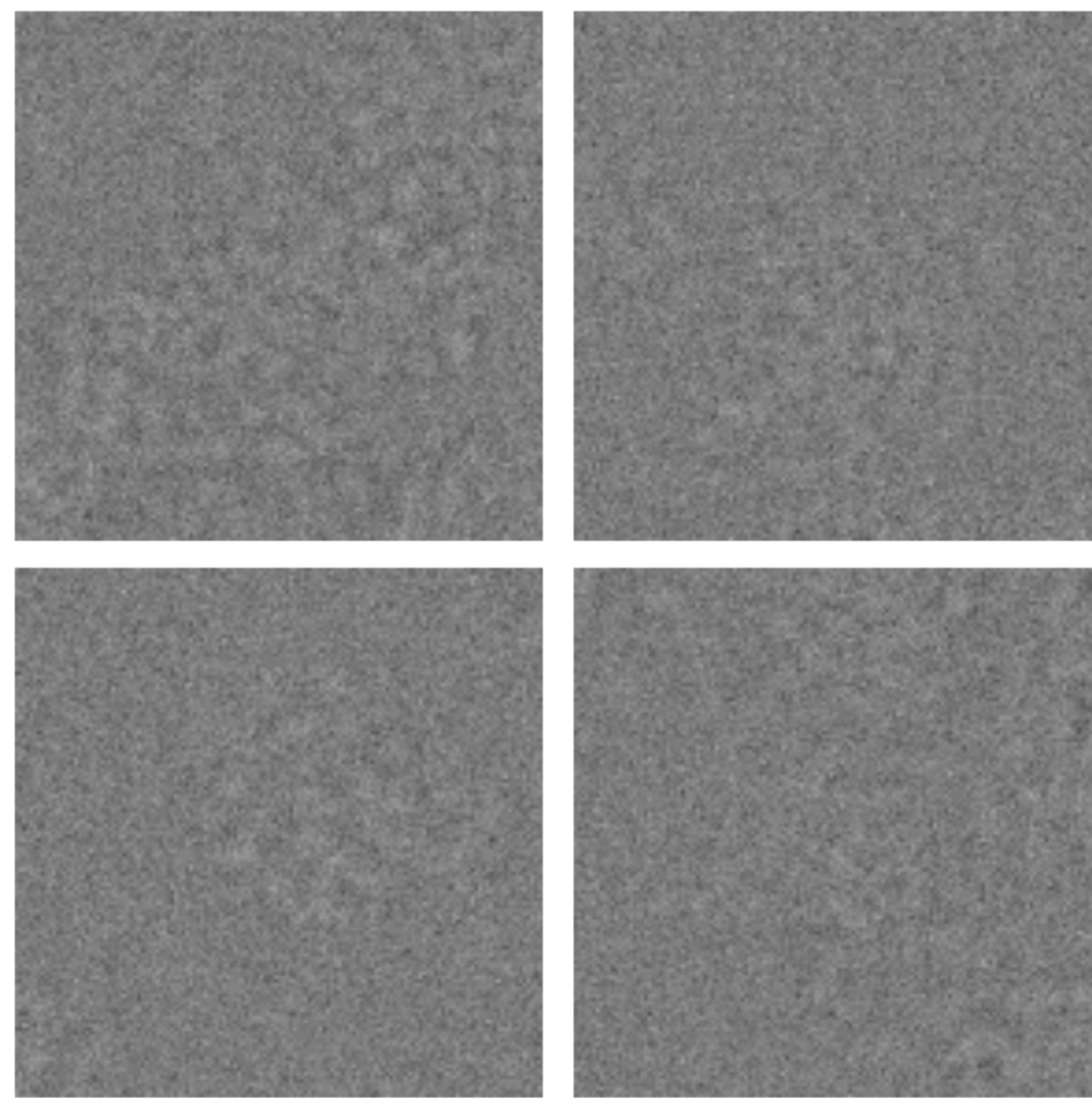
Learning ribosome assembly landscapes [EMPIAR-10076]

Modular assembly of the bacterial large ribosomal subunit (LSU)

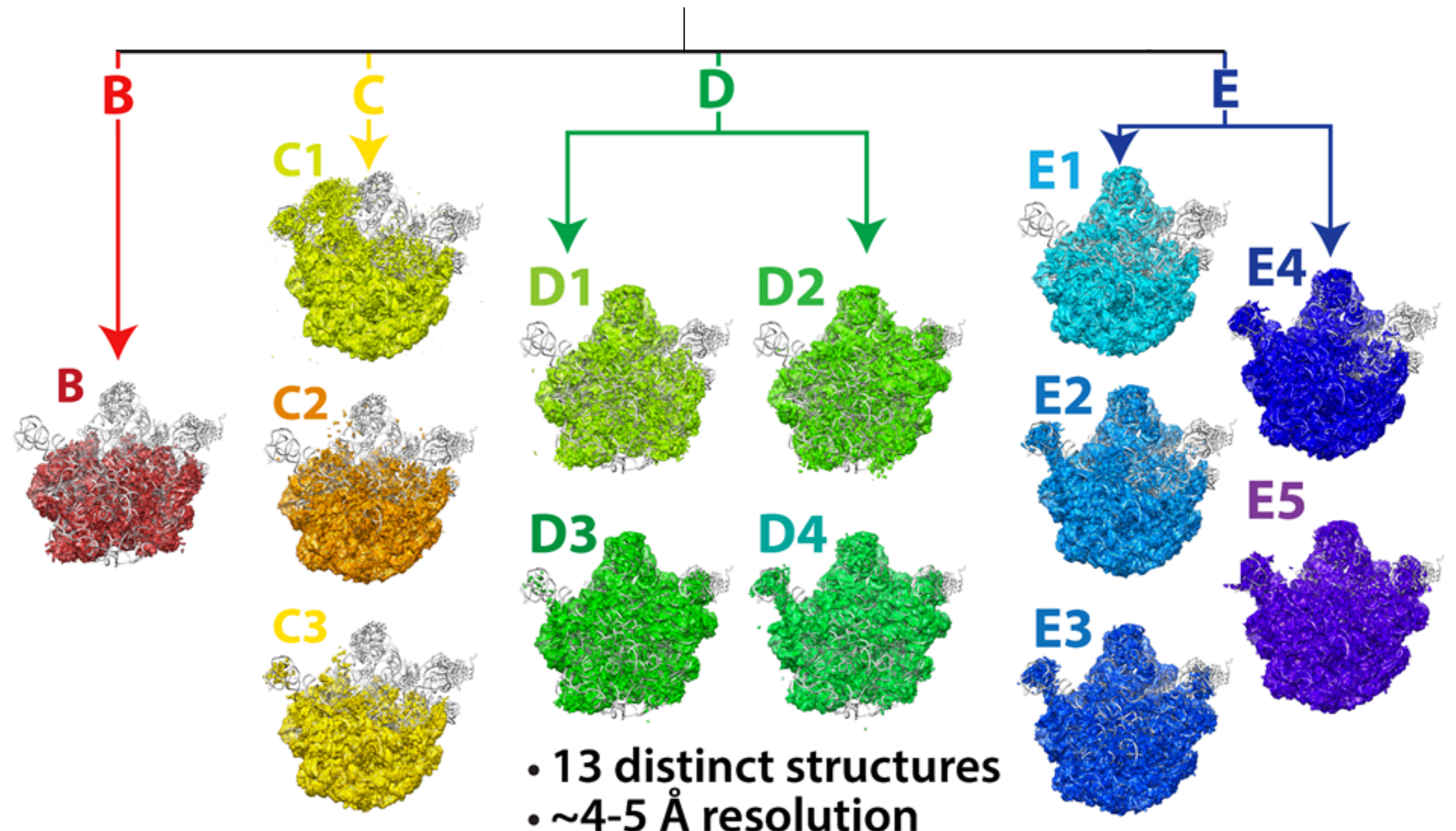
Dataset: 131k cryo-EM images of a mixture of LSU assembly intermediates

4 major and 13 minor states of the LSU identified from hierarchical multiclass reconstruction

Example images



[EMPIAR-10076]



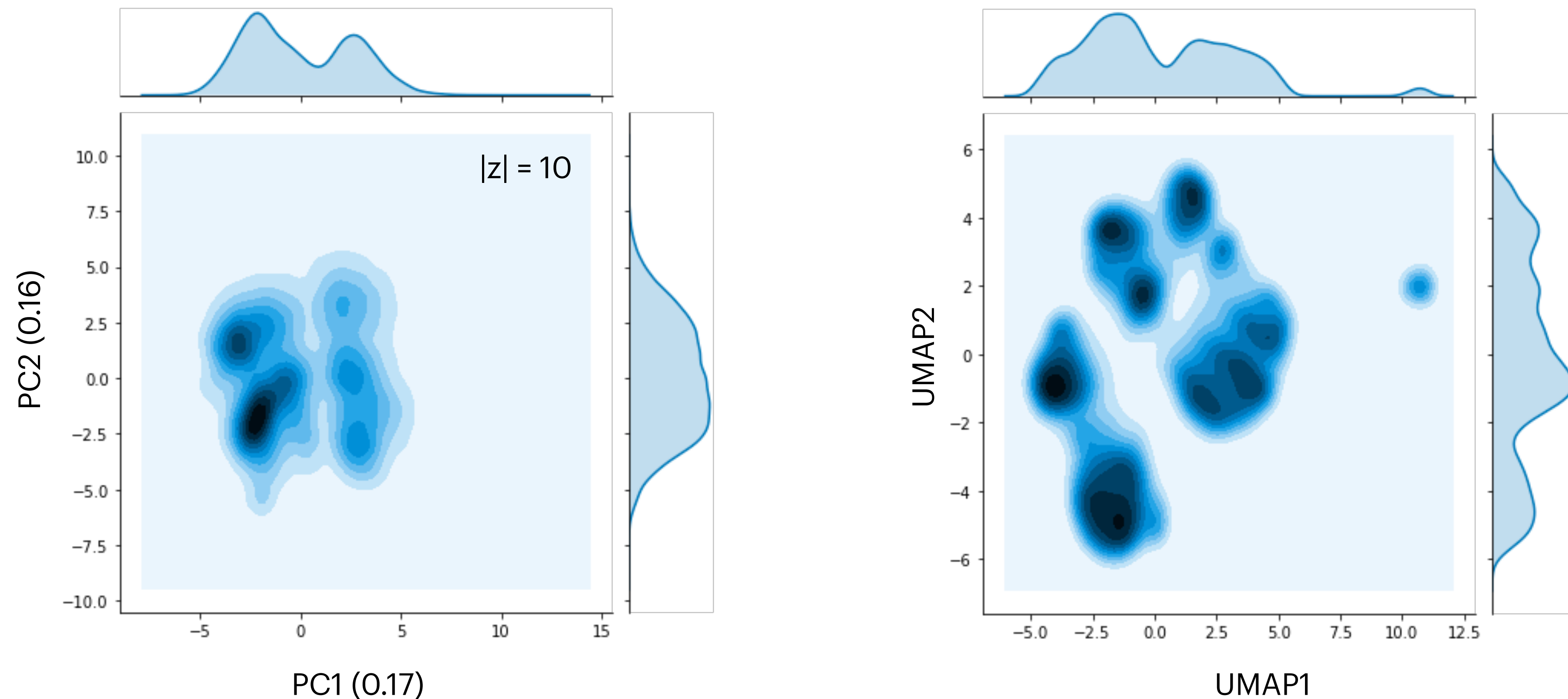
Learning ribosome assembly landscapes [EMPIAR-10076]

Modular assembly of the bacterial large ribosomal subunit (LSU)

Dataset: 131k cryo-EM images of a mixture of LSU assembly intermediates

4 major and 13 minor states of the LSU identified from hierarchical multiclass reconstruction

Latent embeddings from cryoDRGN

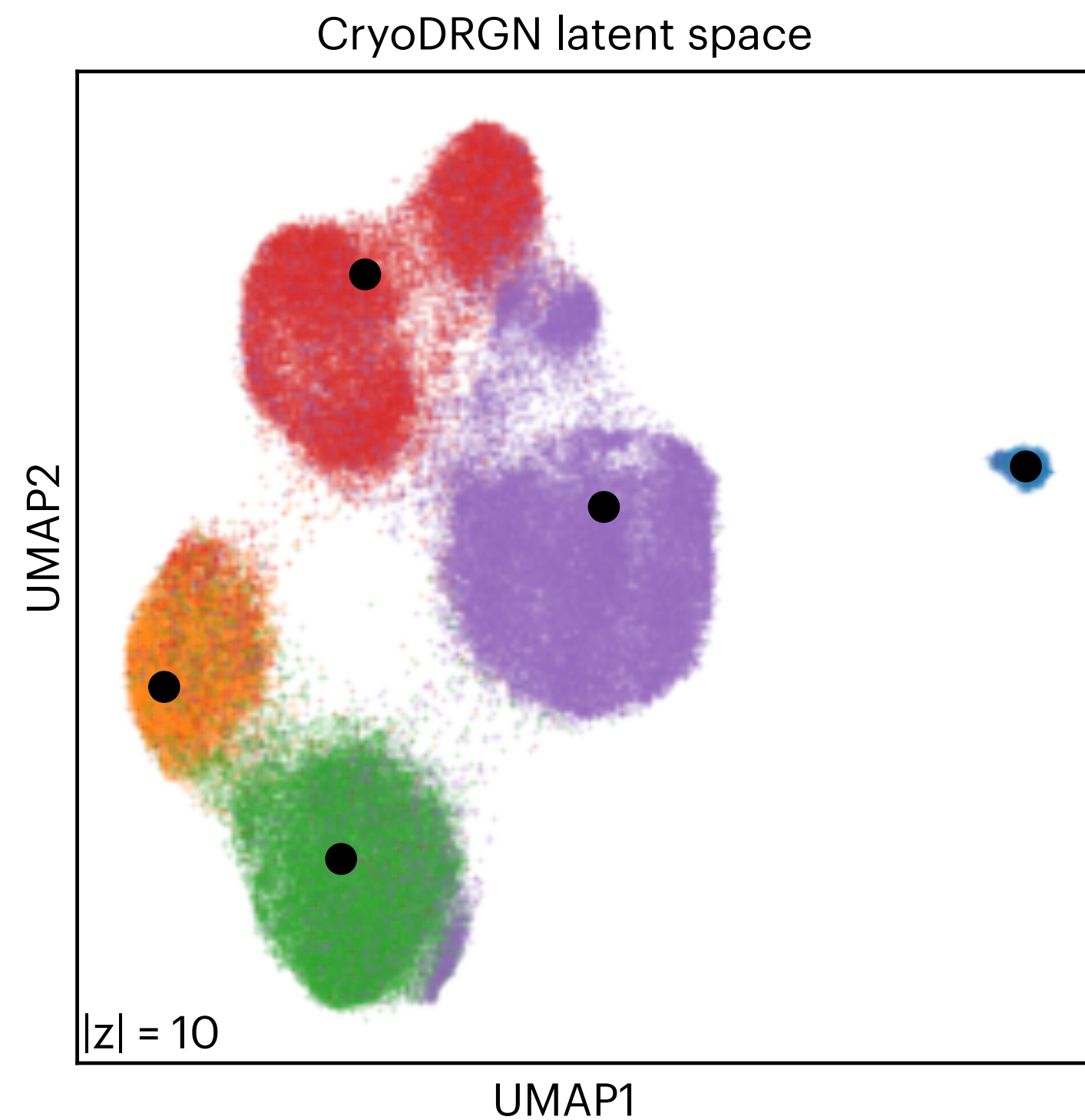


Learning ribosome assembly landscapes [EMPIAR-10076]

Modular assembly of the bacterial large ribosomal subunit (LSU)

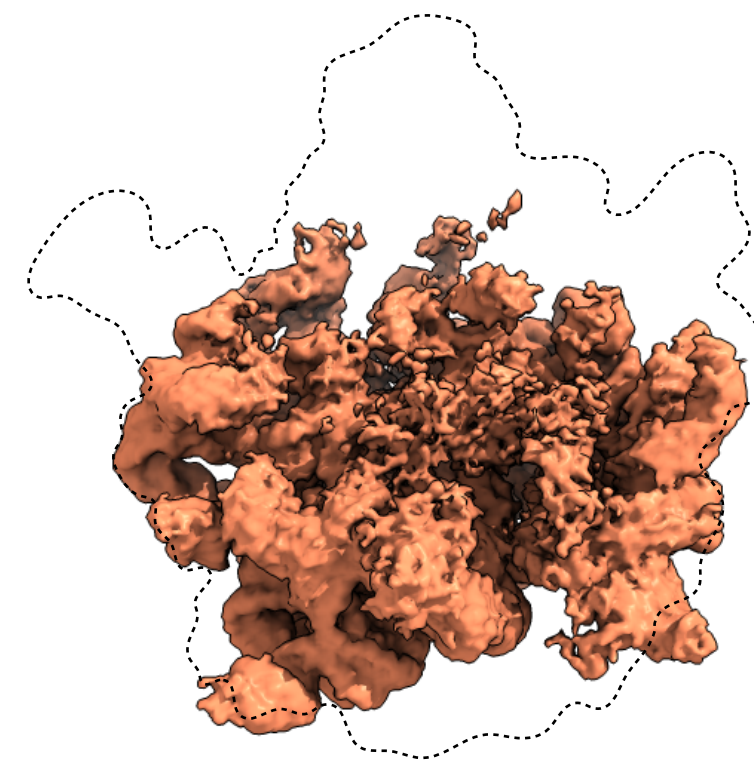
Dataset: 131k cryo-EM images of a mixture of LSU assembly intermediates

4 major and 13 minor states of the LSU identified from hierarchical multiclass reconstruction

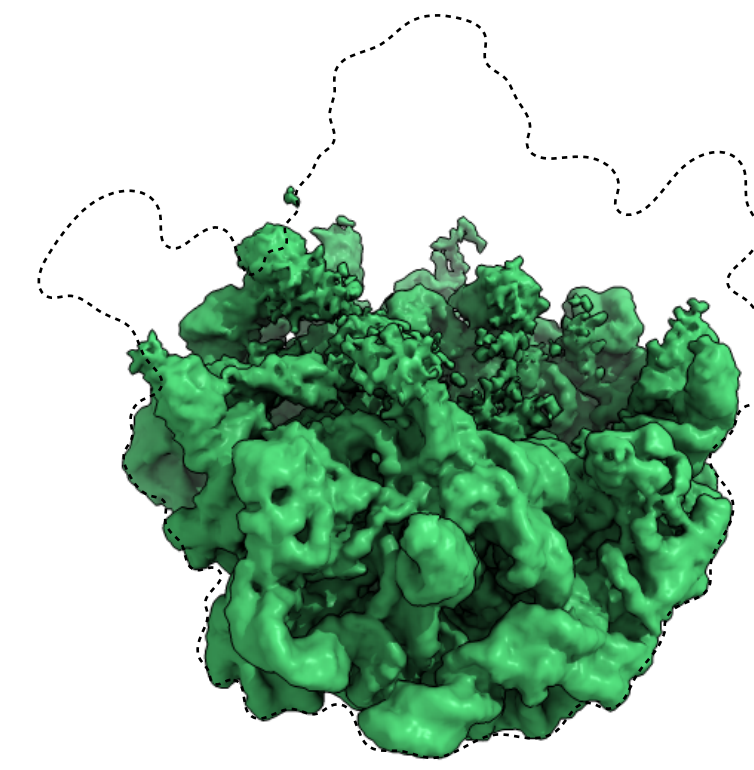


Published class assignment, major states

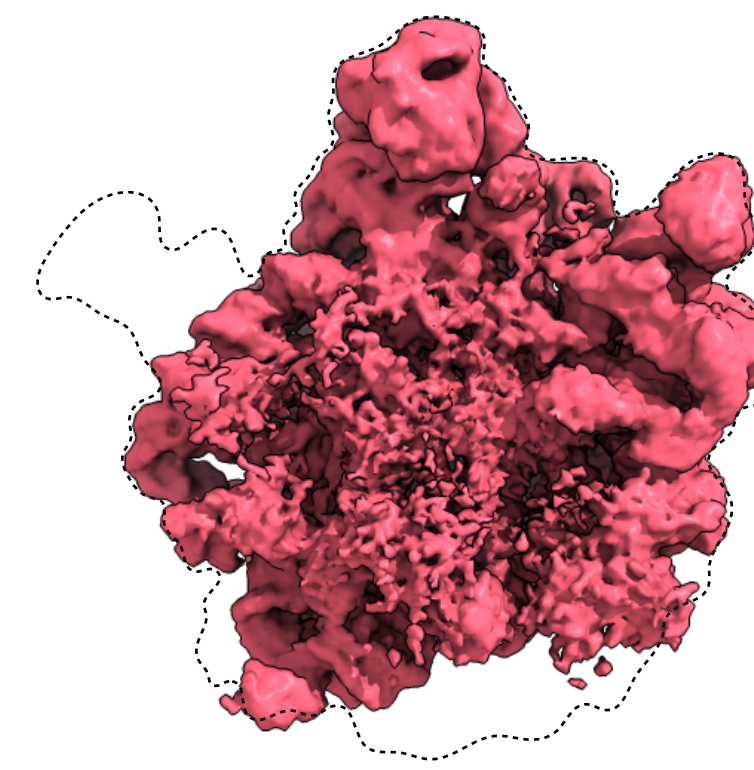
● A ● B ● C ● D ● E



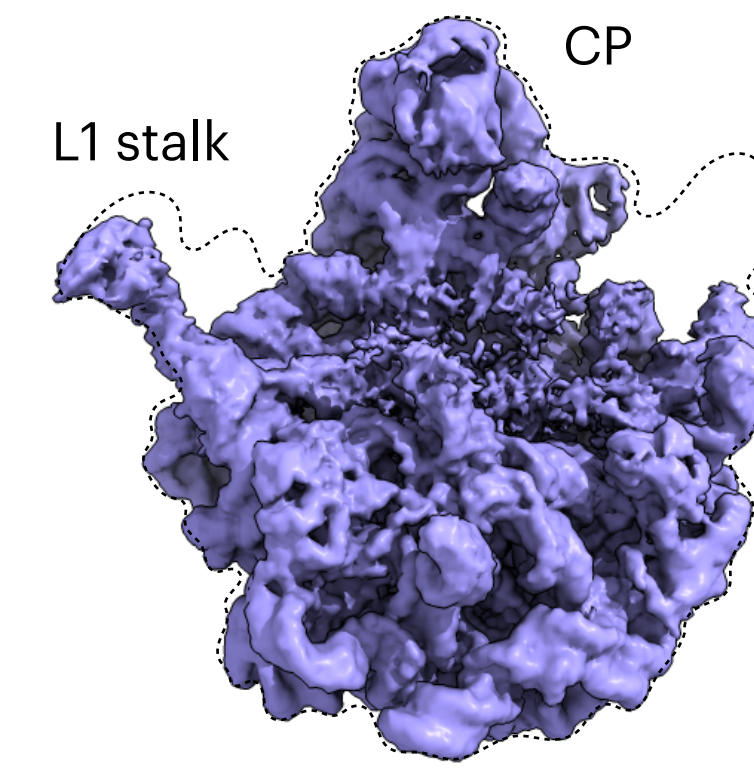
LSU assembly class B



LSU assembly class C

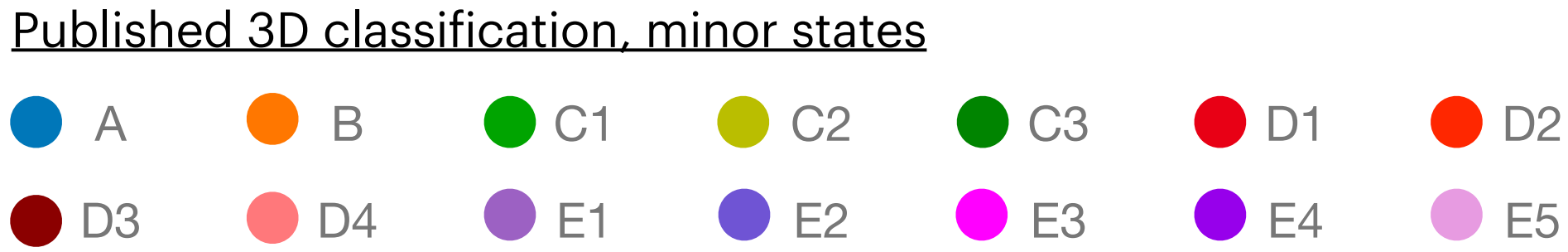
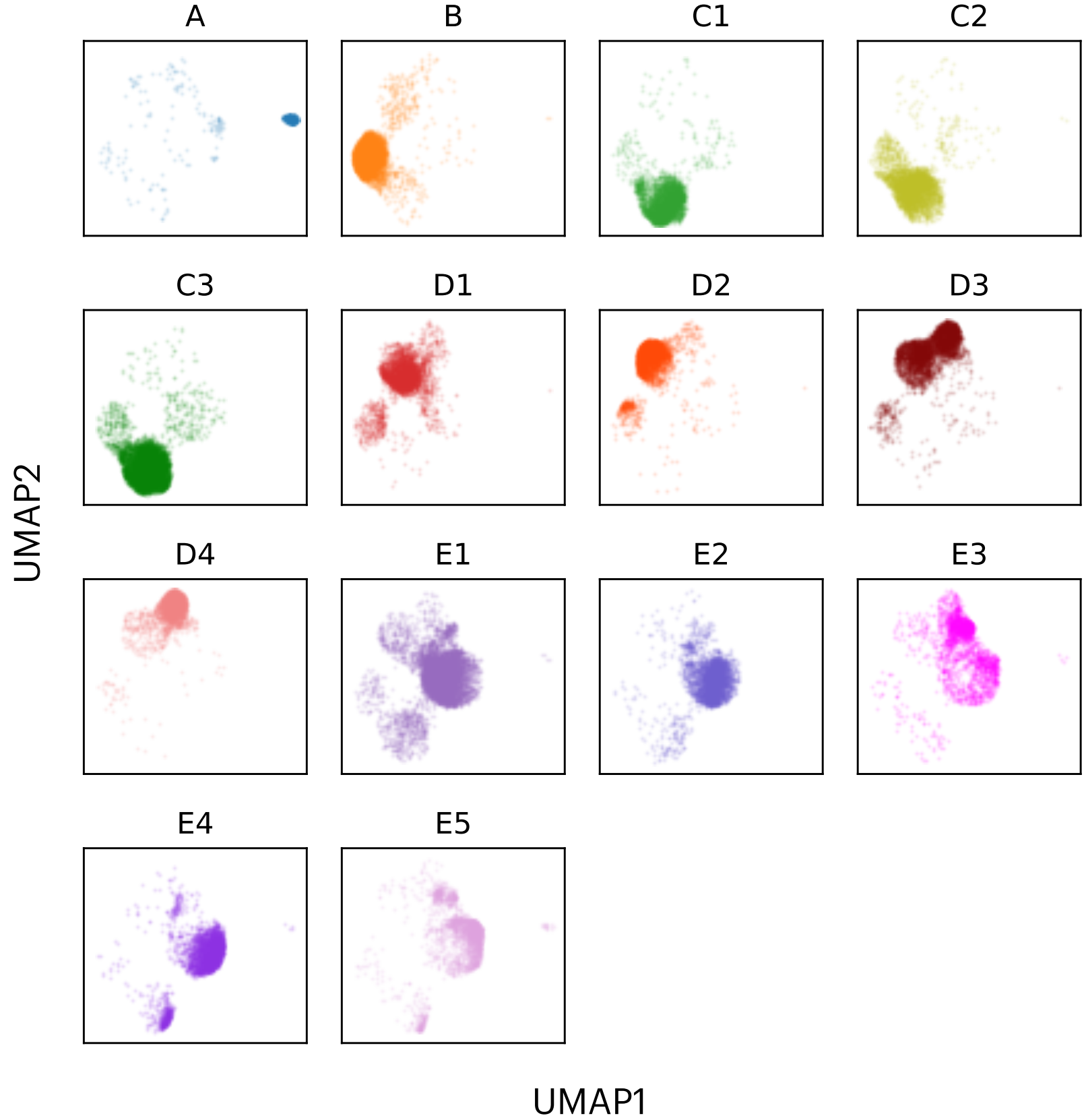
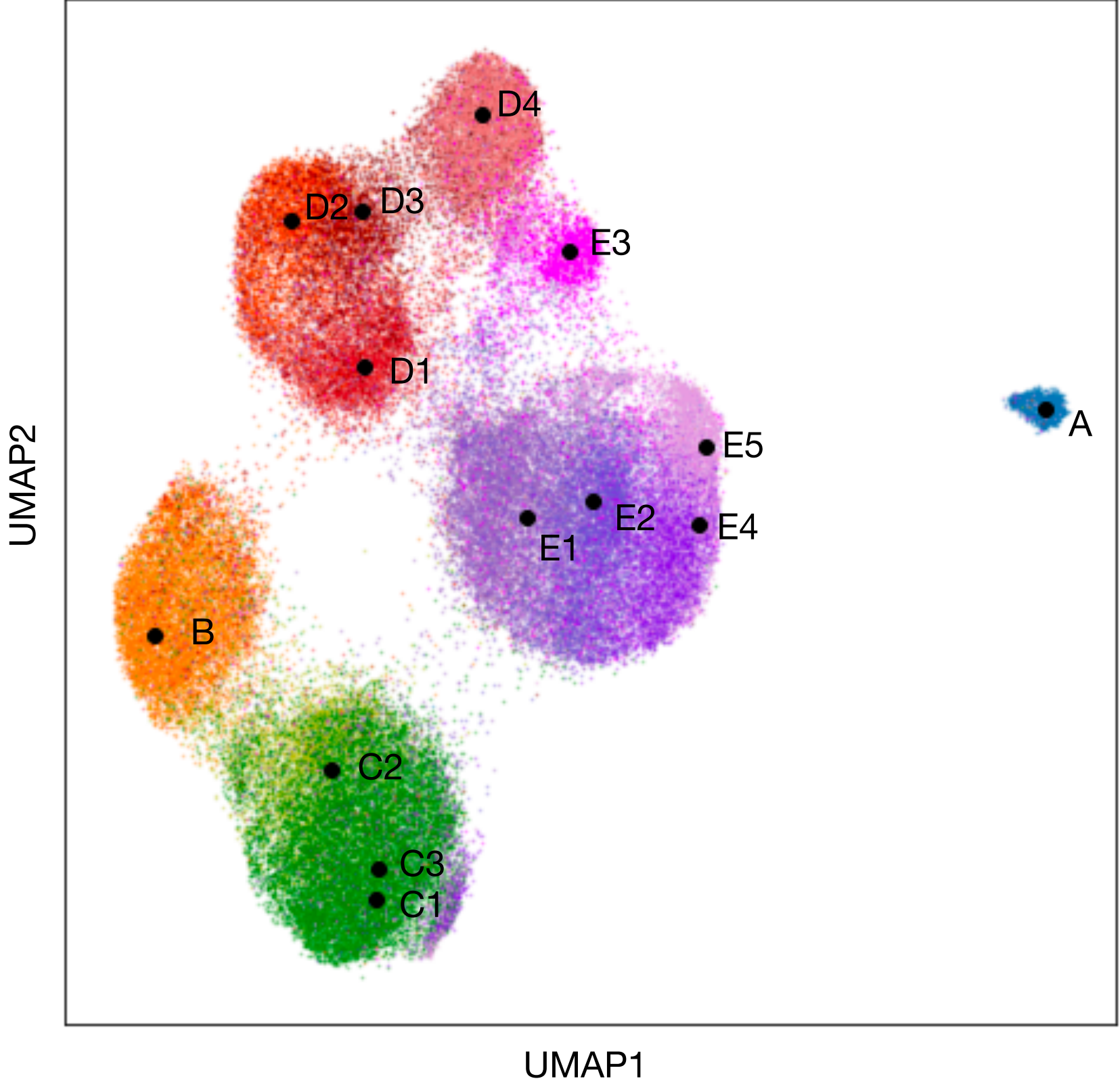


LSU assembly class D

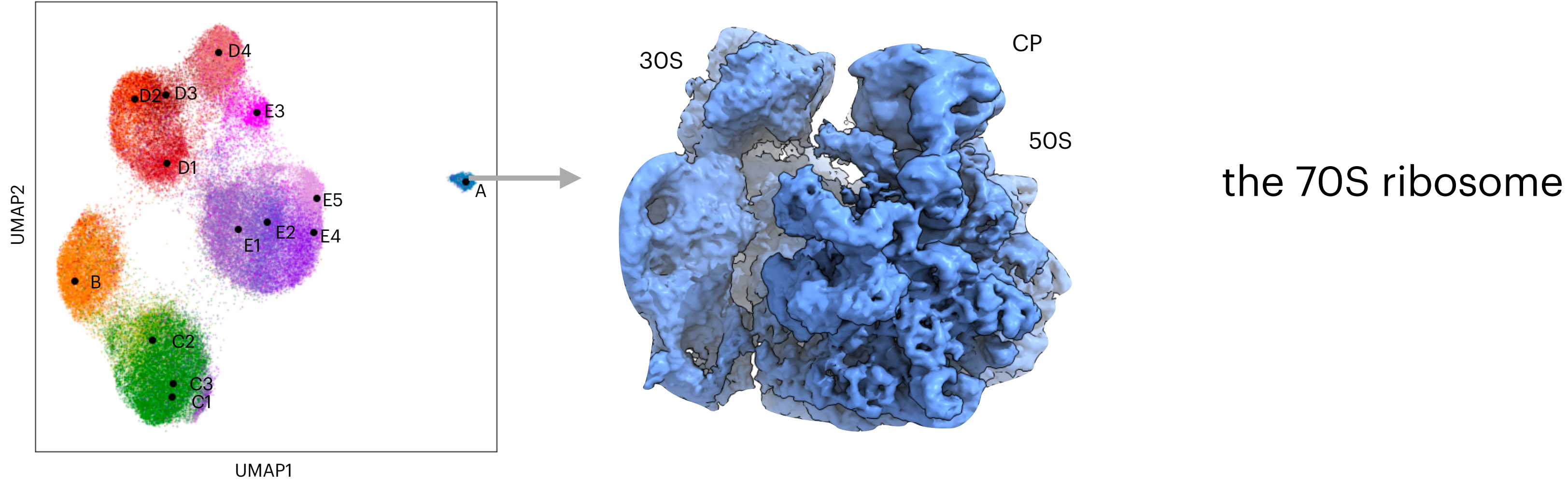


LSU assembly class E

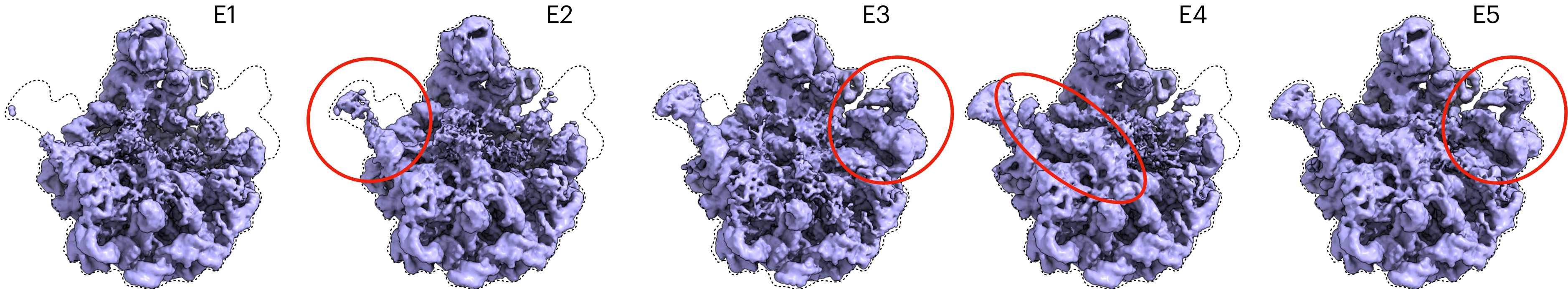
Clusters in the latent space vs. expert-driven hierarchical classification



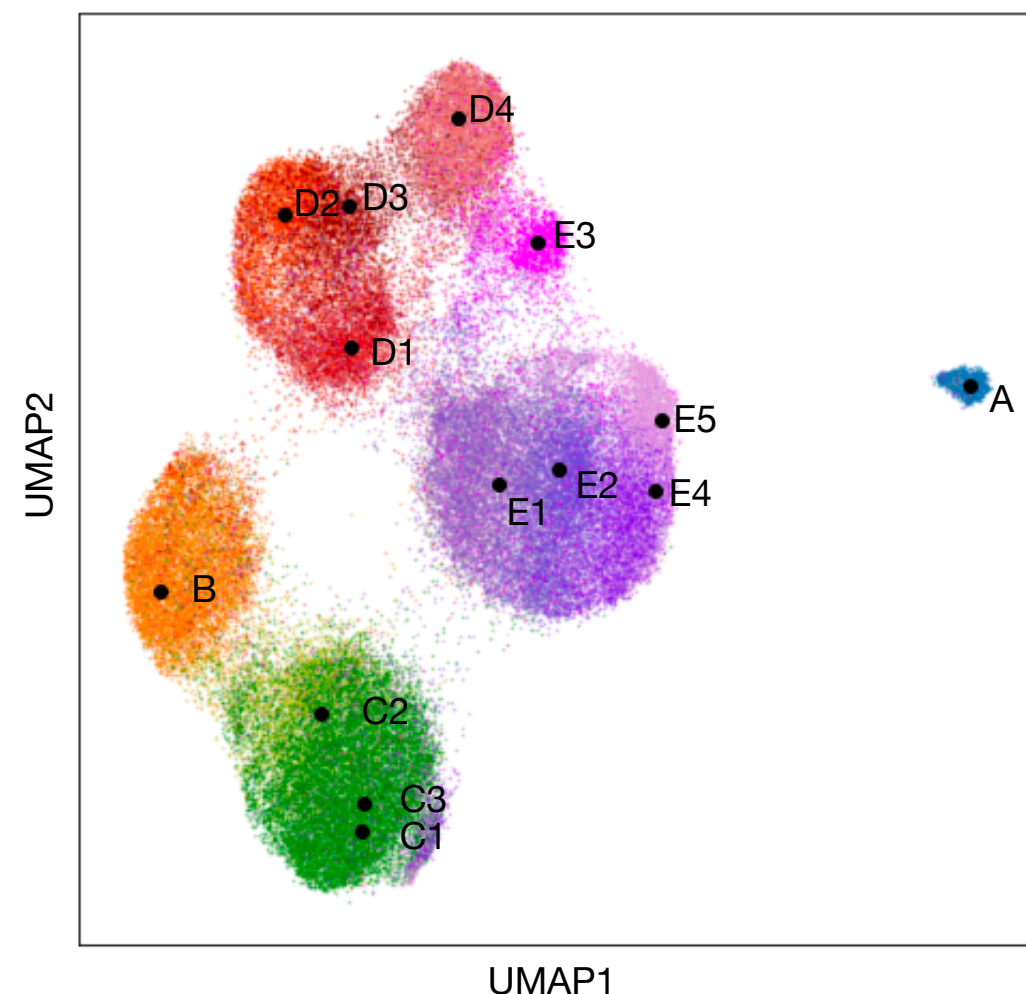
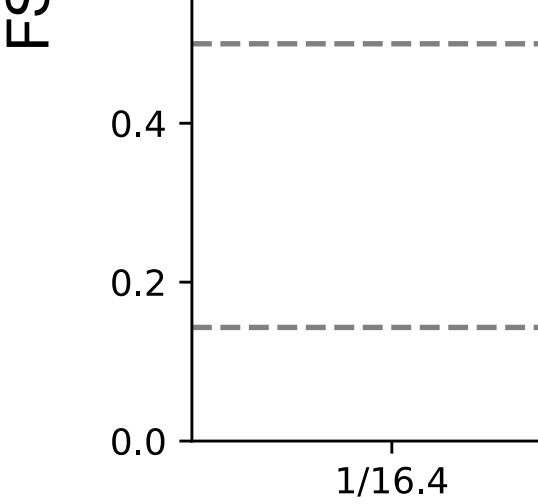
Additional samples from the latent space



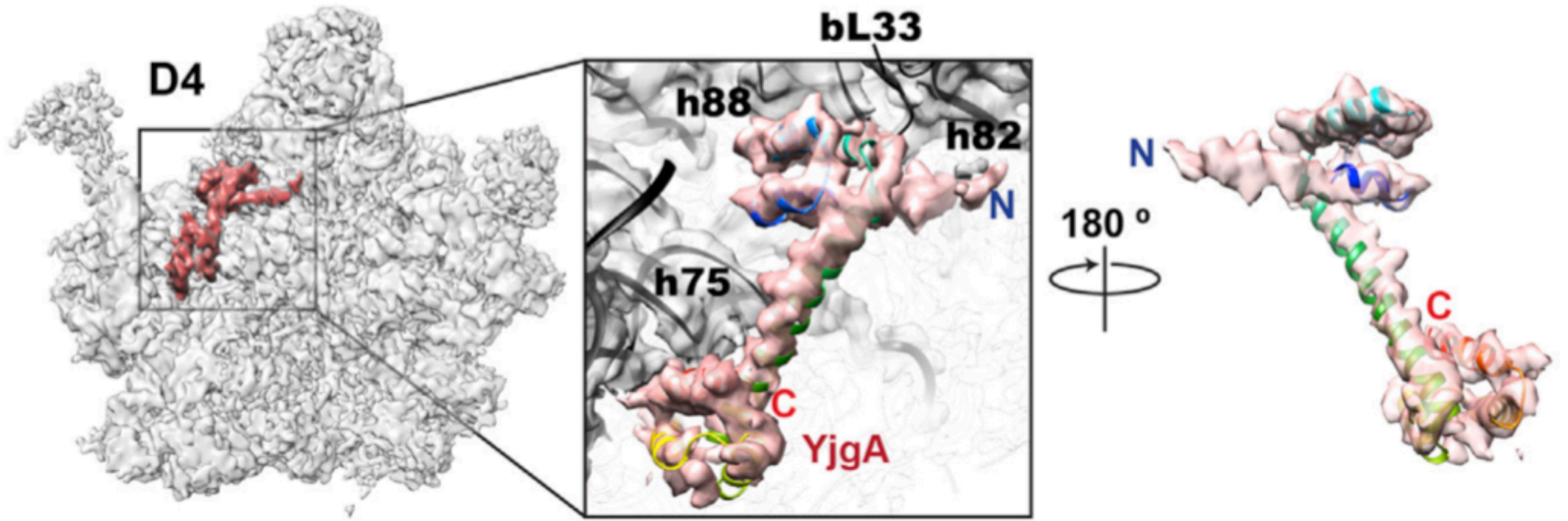
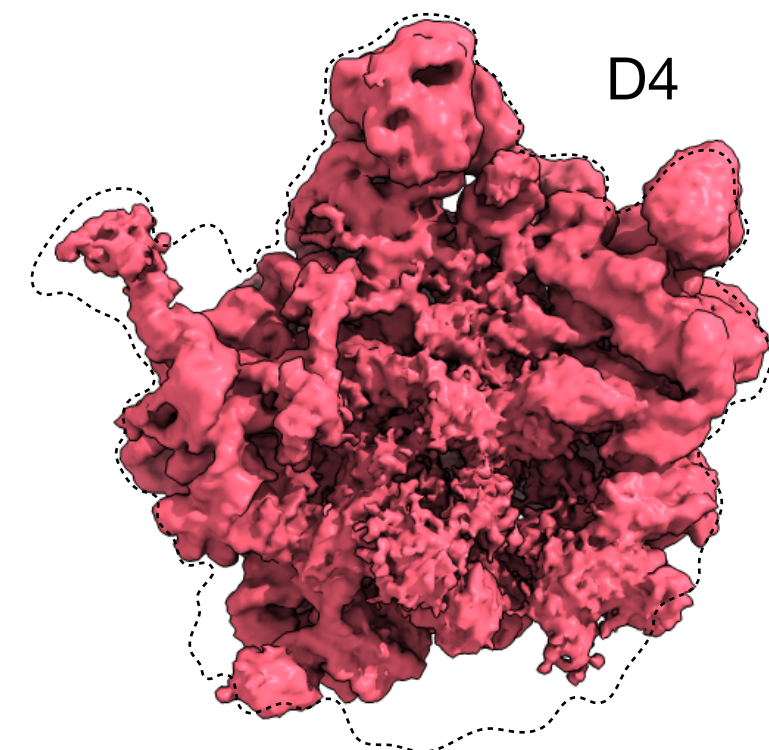
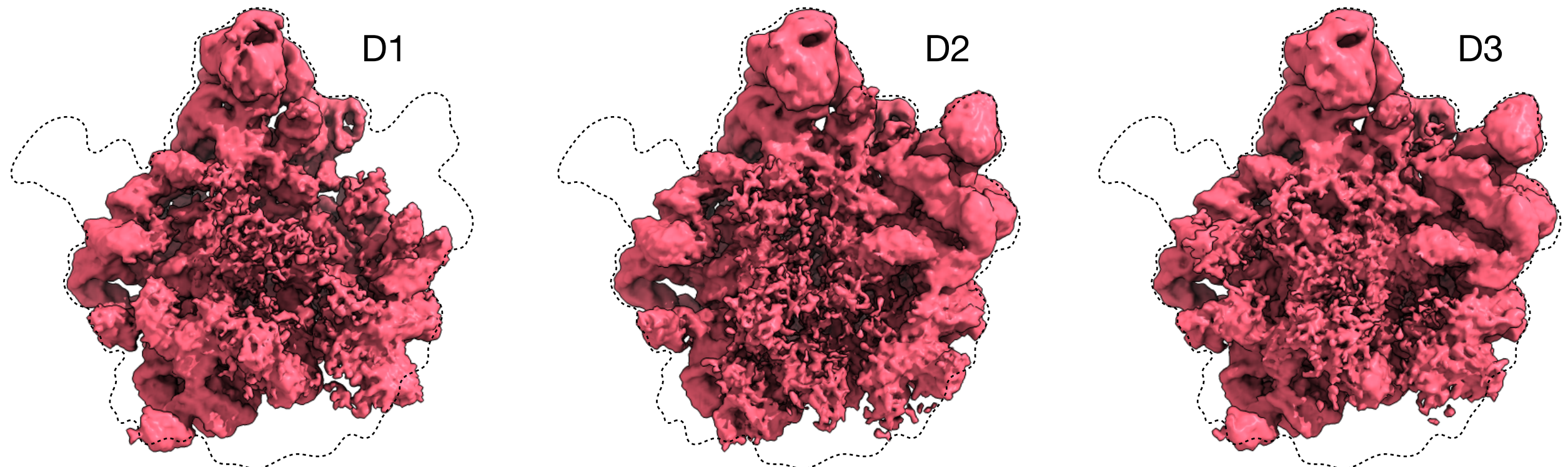
E class minor assembly states



Additional samples from the latent space

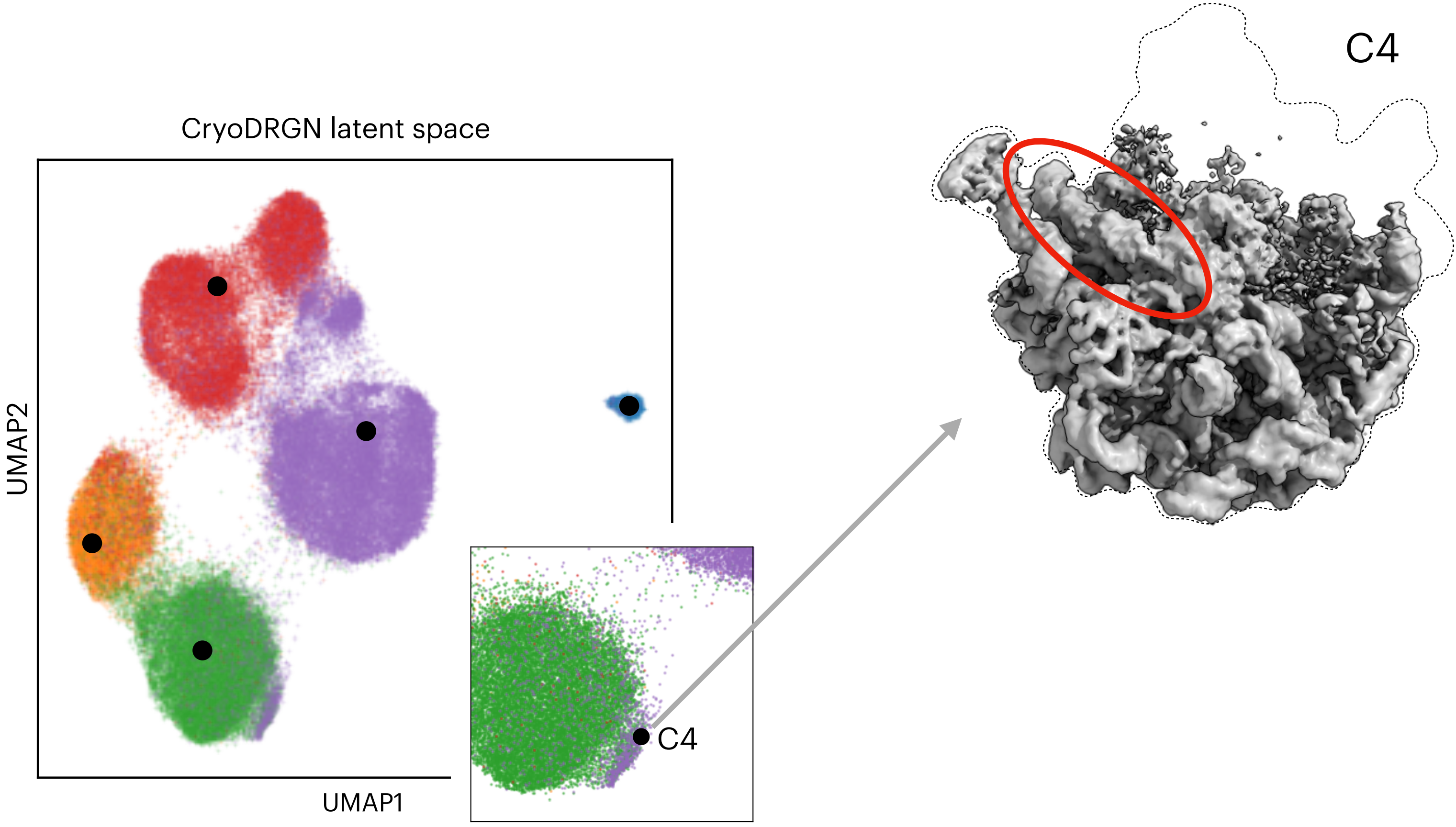


D class minor assembly states



Adapted from Figure 5, Davis et al 2016

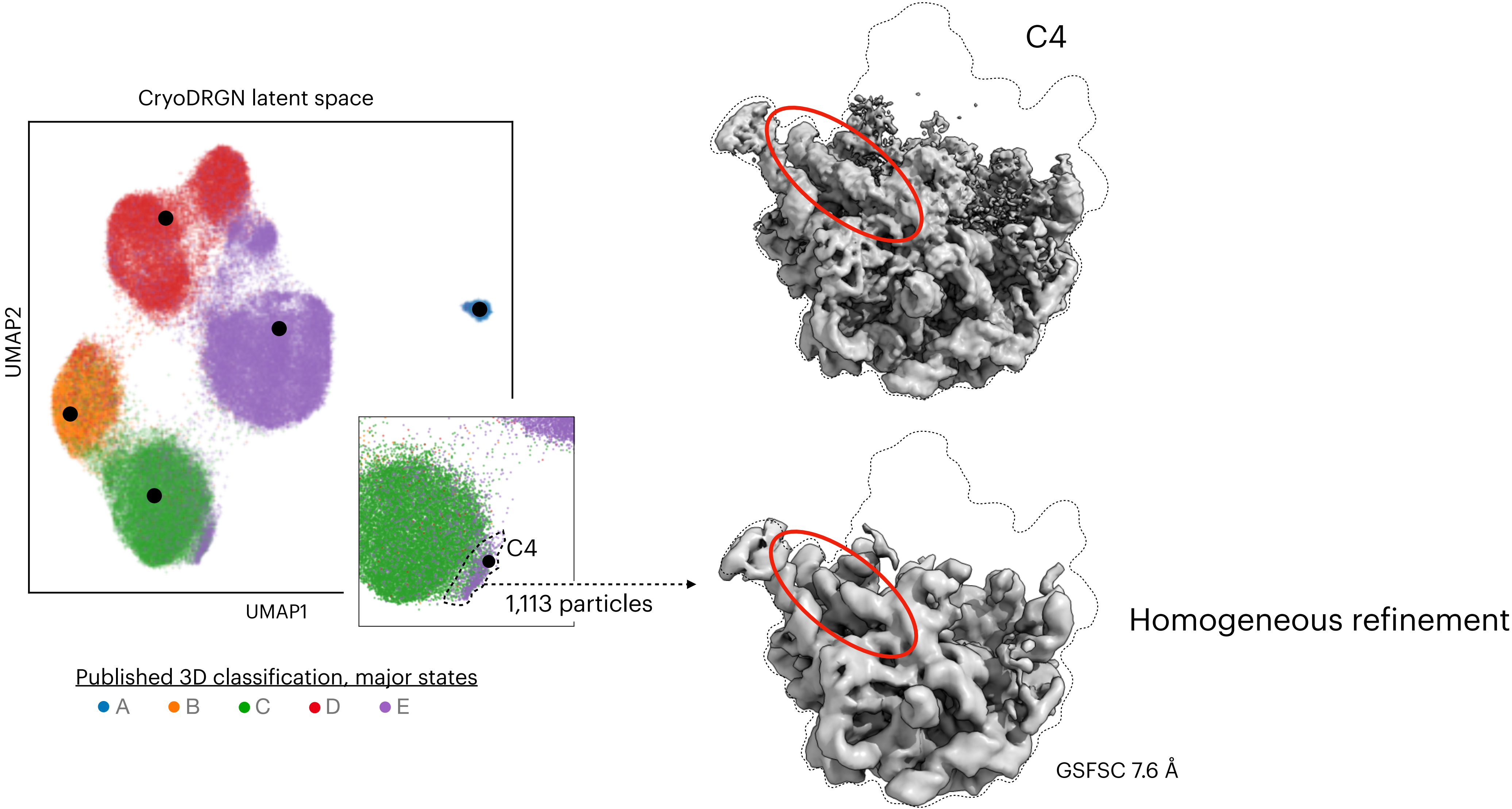
Discovery of a new assembly state, C4



Published 3D classification, major states

- A
- B
- C
- D
- E

Discovery of a new assembly state, C4

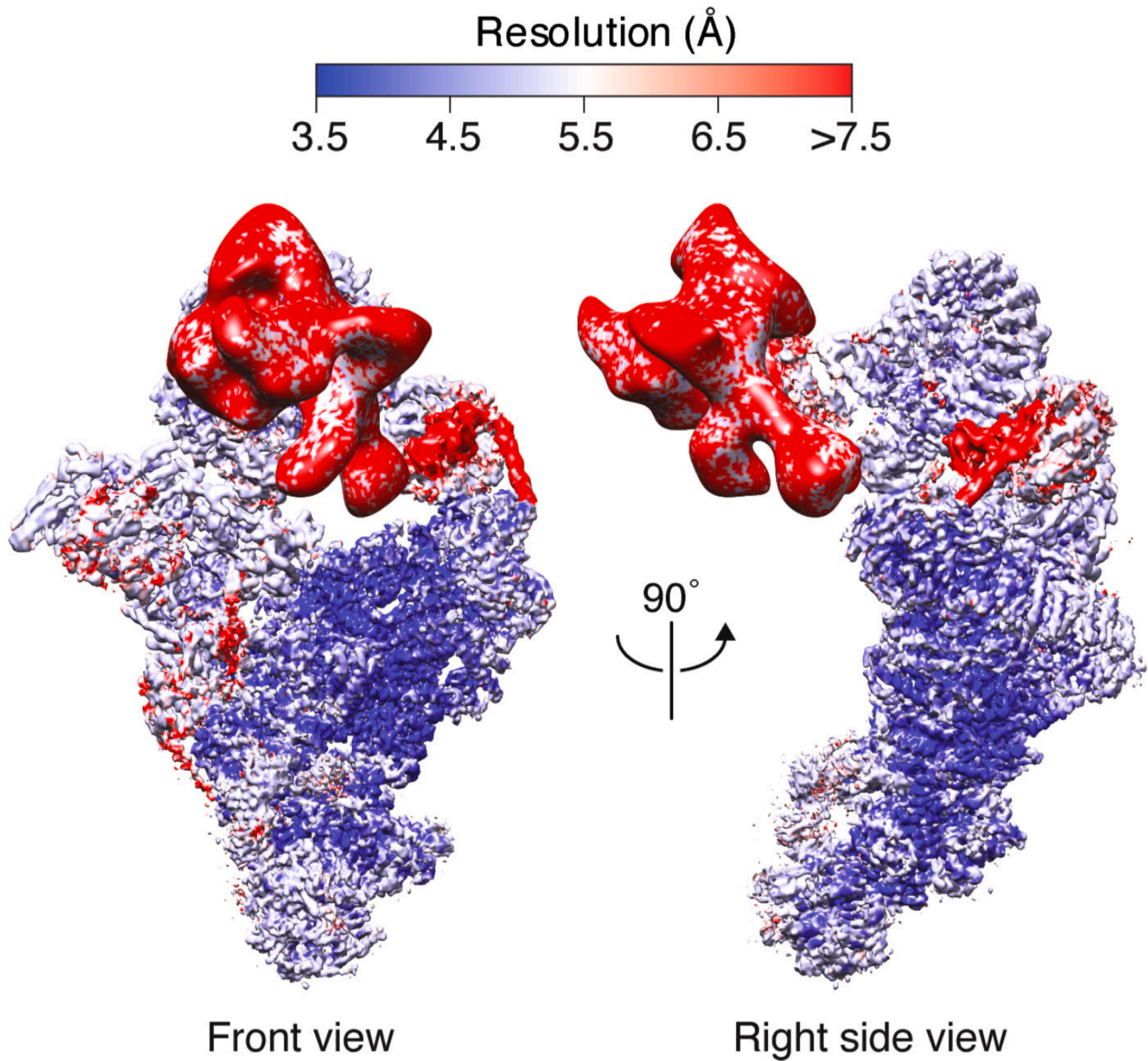


Roadmap

- Motivation and background
- CryoDRGN: Deep Reconstructing Generative Networks
- Validation on synthetic benchmarks
- **CryoDRGN reconstructions of real data**
 - Uncovering residual heterogeneity in high resolution “homogeneous” datasets
 - Discovering new states of the assembling ribosome
 - **Reconstructing continuous motions of the pre-catalytic spliceosome**
- Future vision

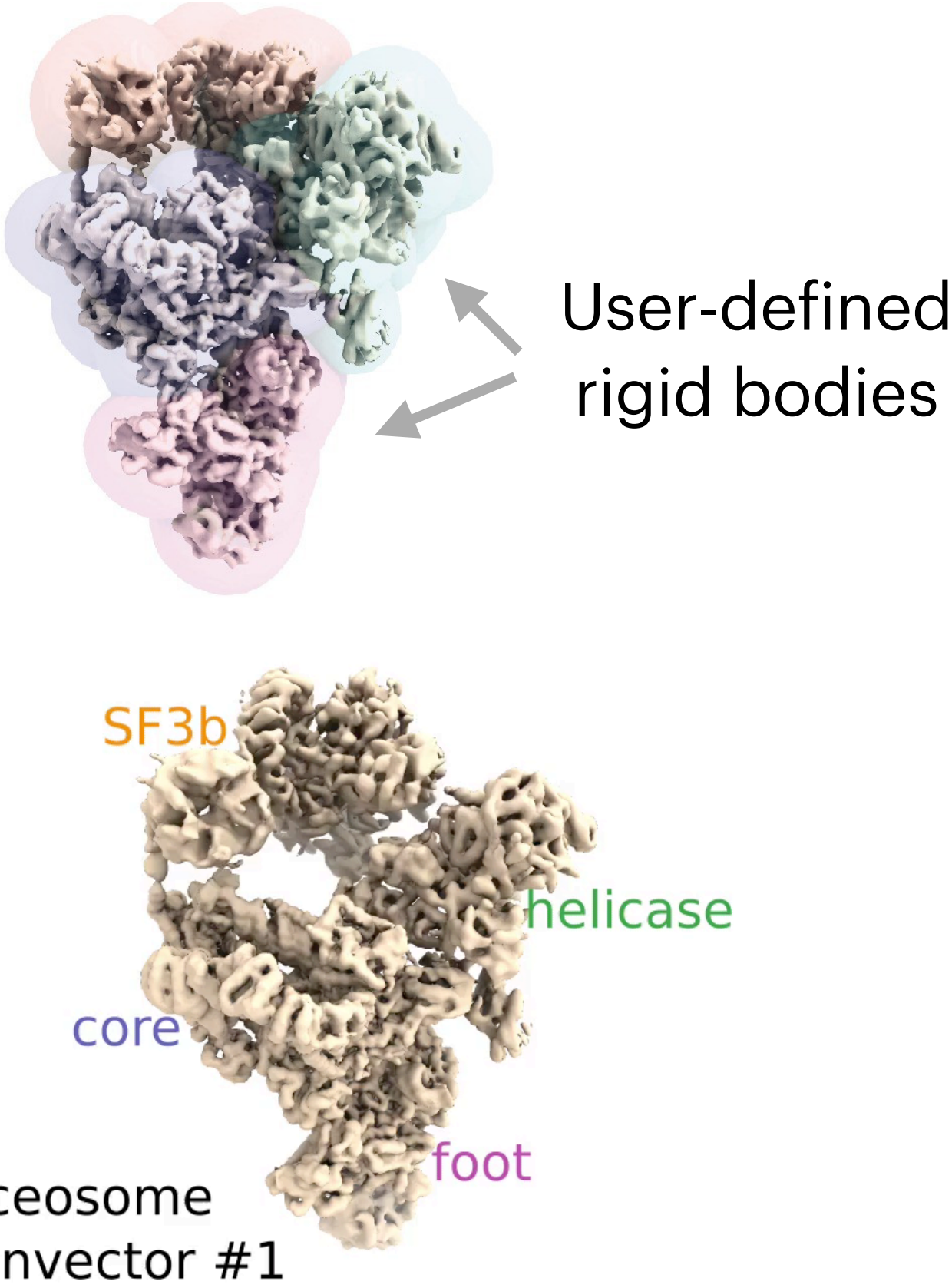
Structure of the pre-catalytic spliceosome

Sub-complexes resolved separately through many rounds of focused classification



Plaschka, Lin, & Nagai. Nature 2017

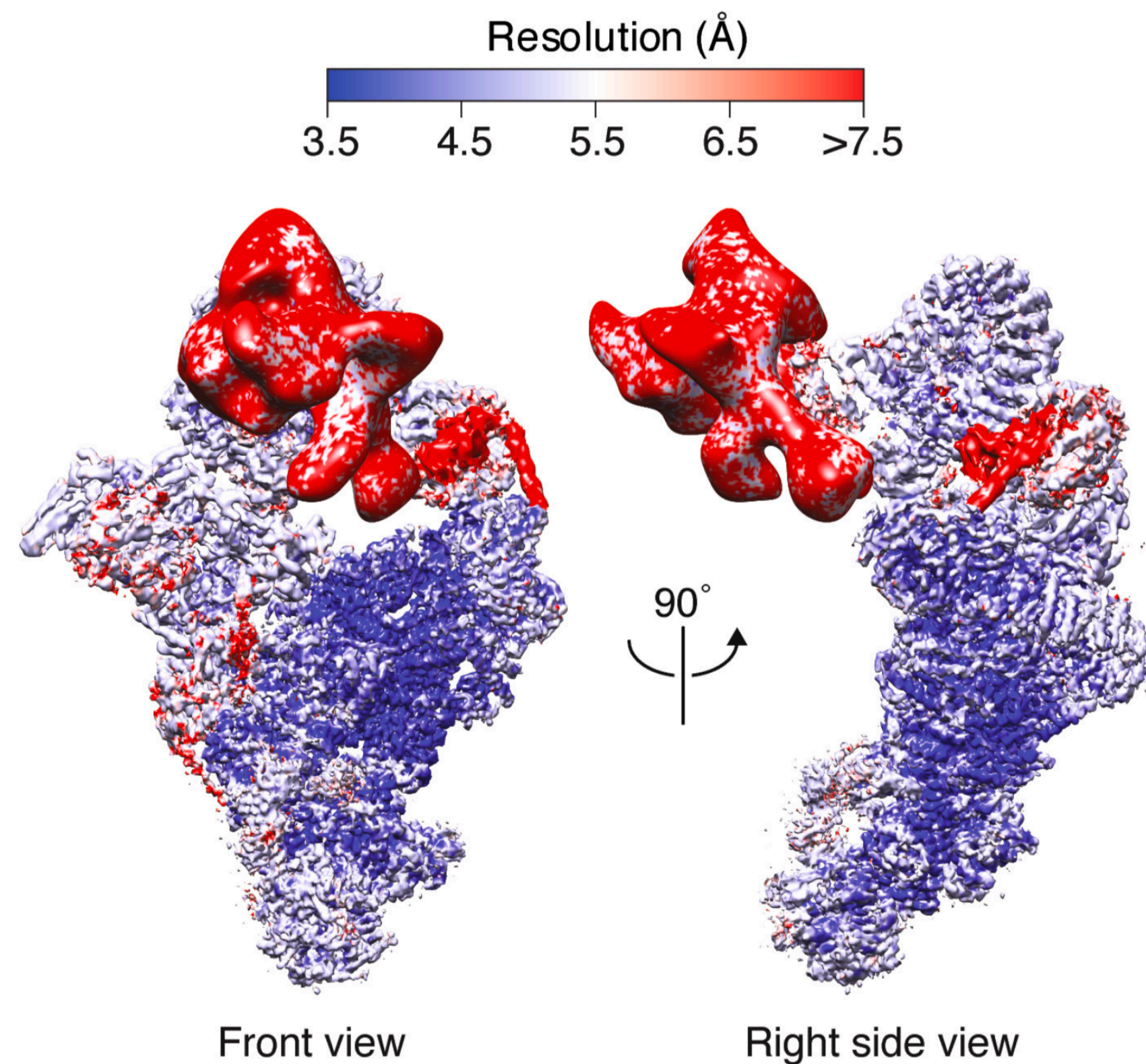
Multibody refinement



Nakane et al. eLife 2018

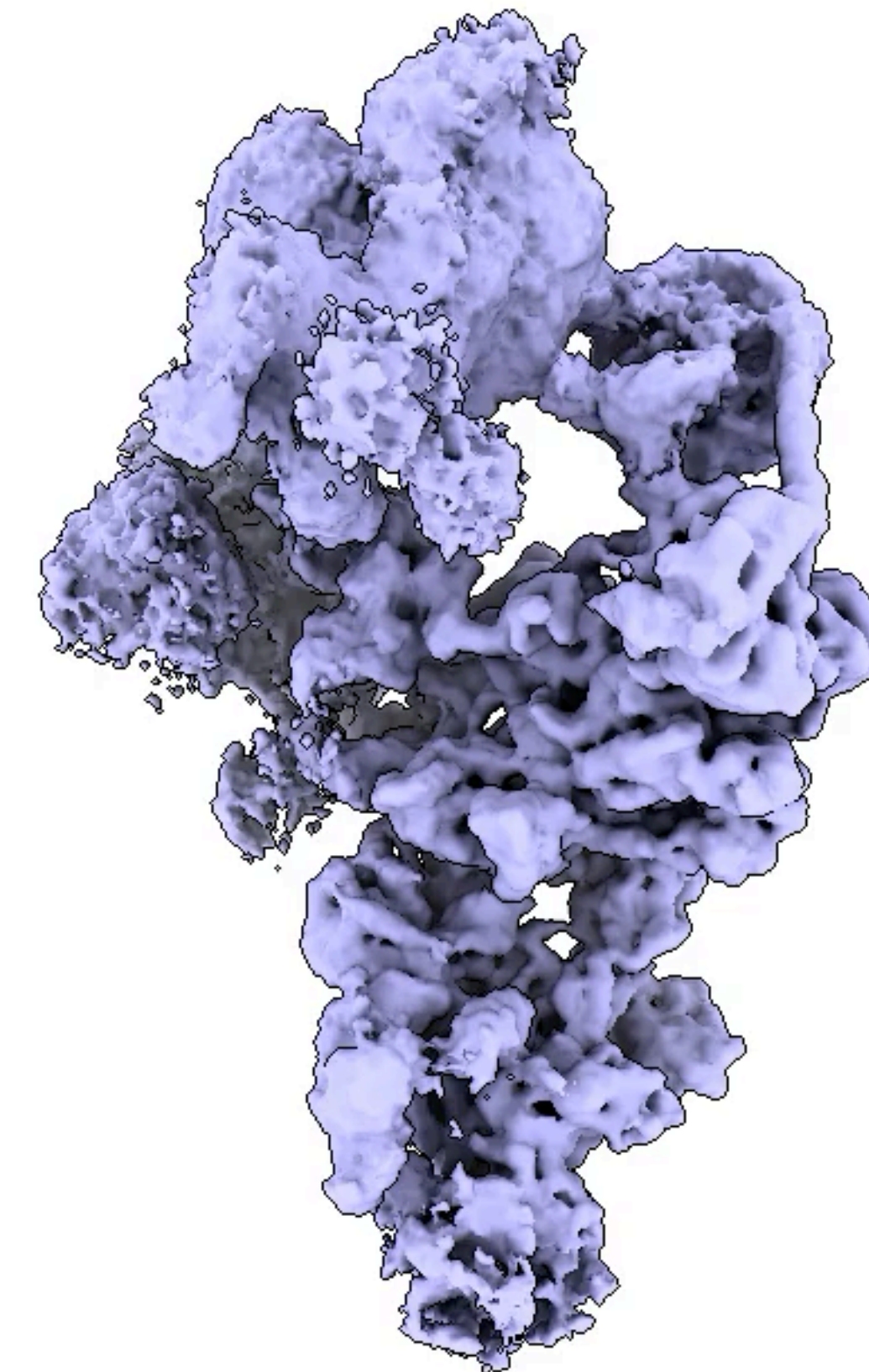
Structure of the pre-catalytic spliceosome

Sub-complexes resolved separately through many rounds of focused classification



Plaschka, Lin, & Nagai. Nature 2017

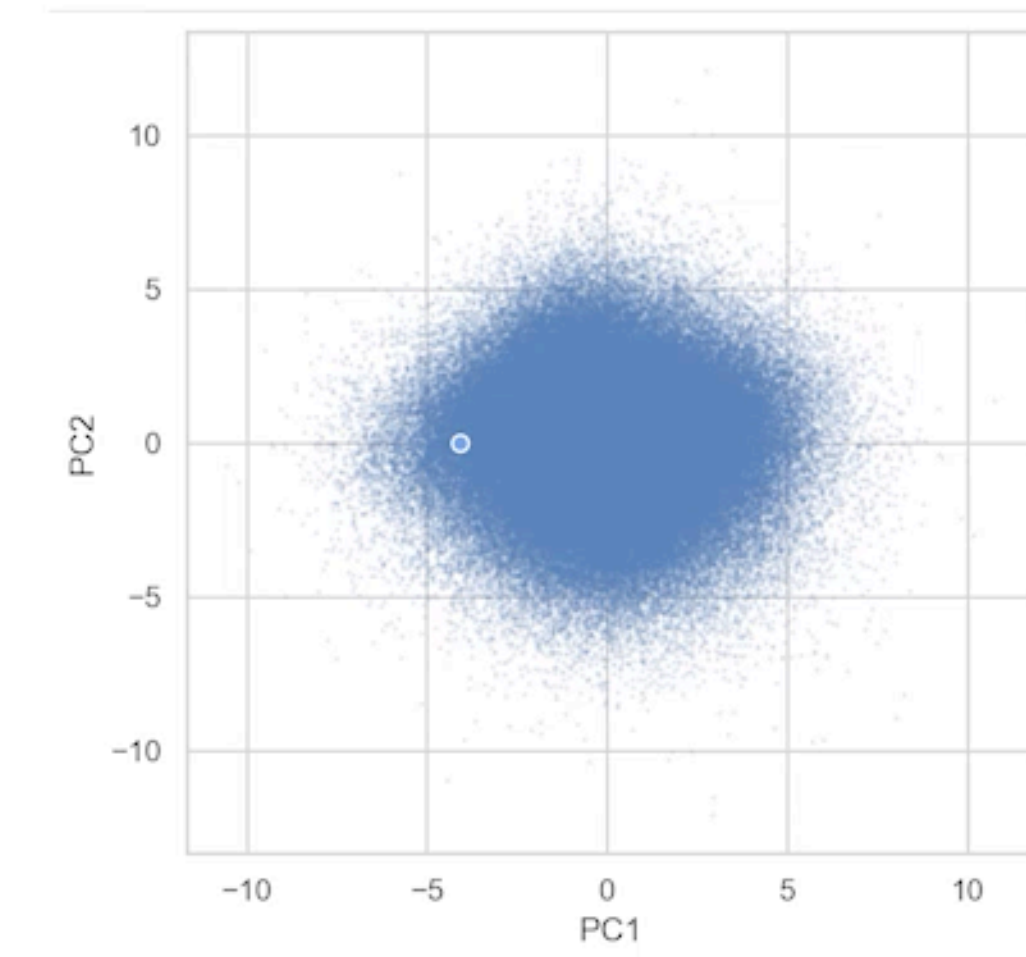
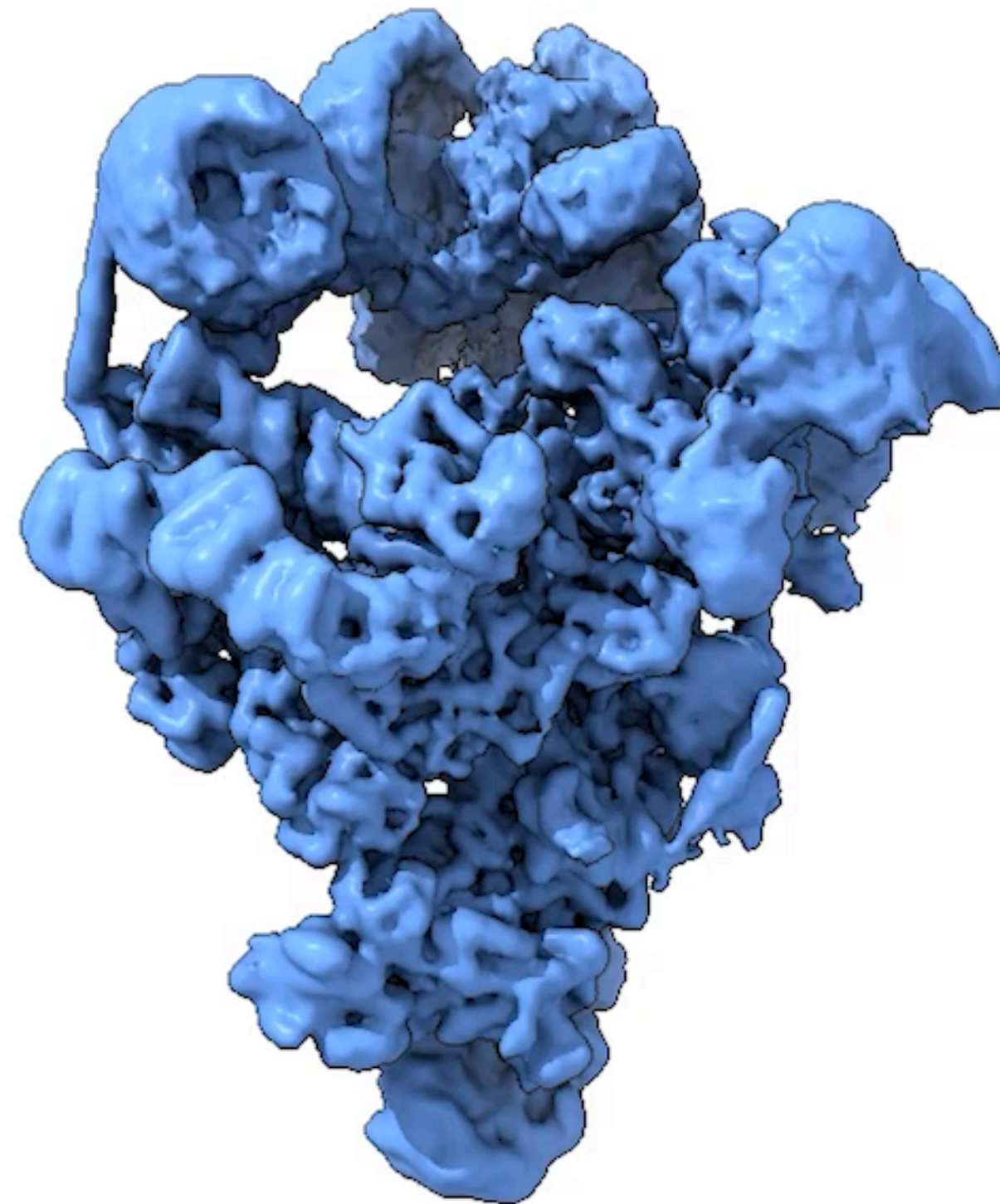
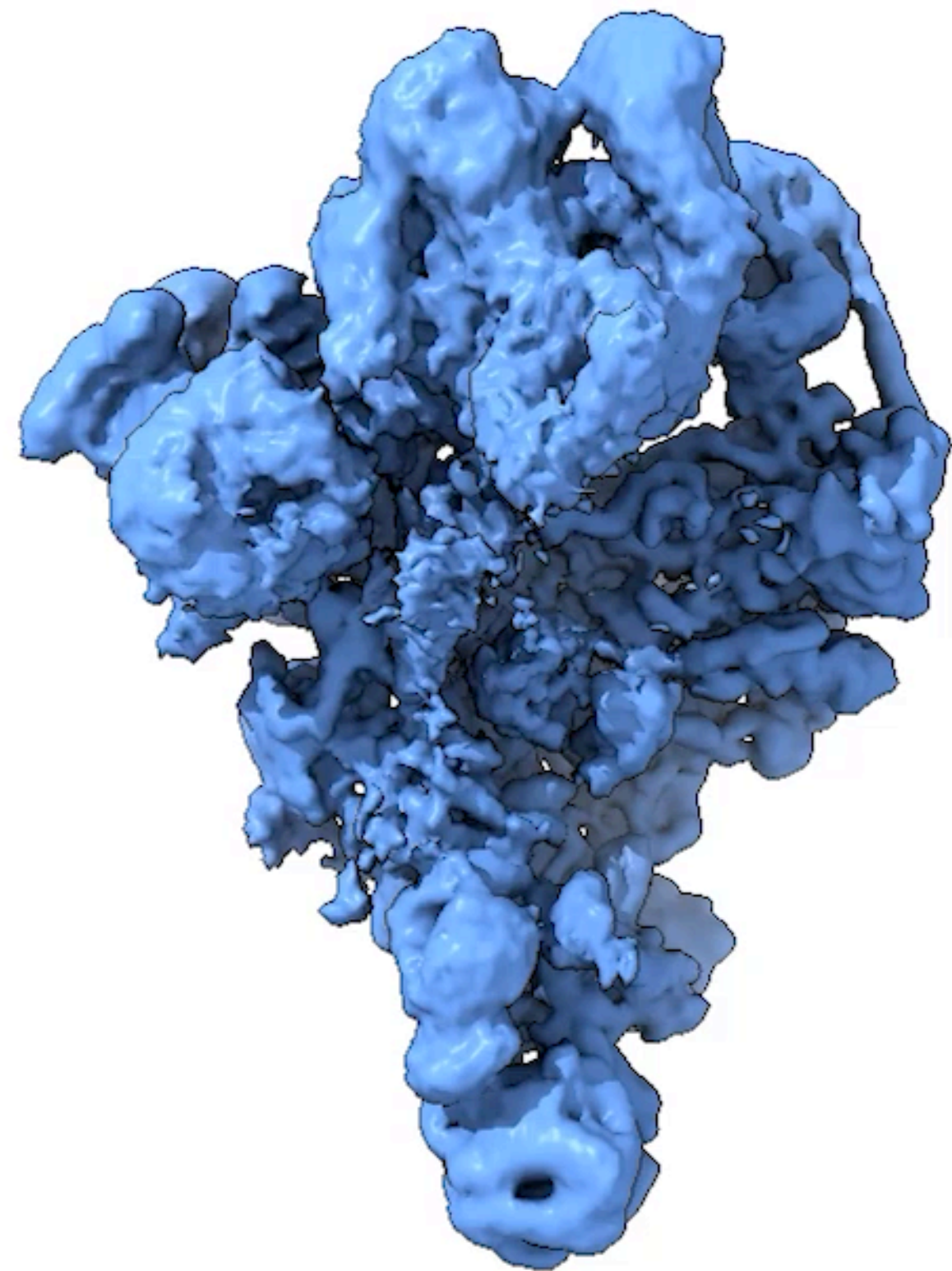
cryoSPARC 3DVA



Punjani & Fleet. JSB 2021

Reconstructing continuous motions of the pre-catalytic spliceosome [EMPIAR-10180]

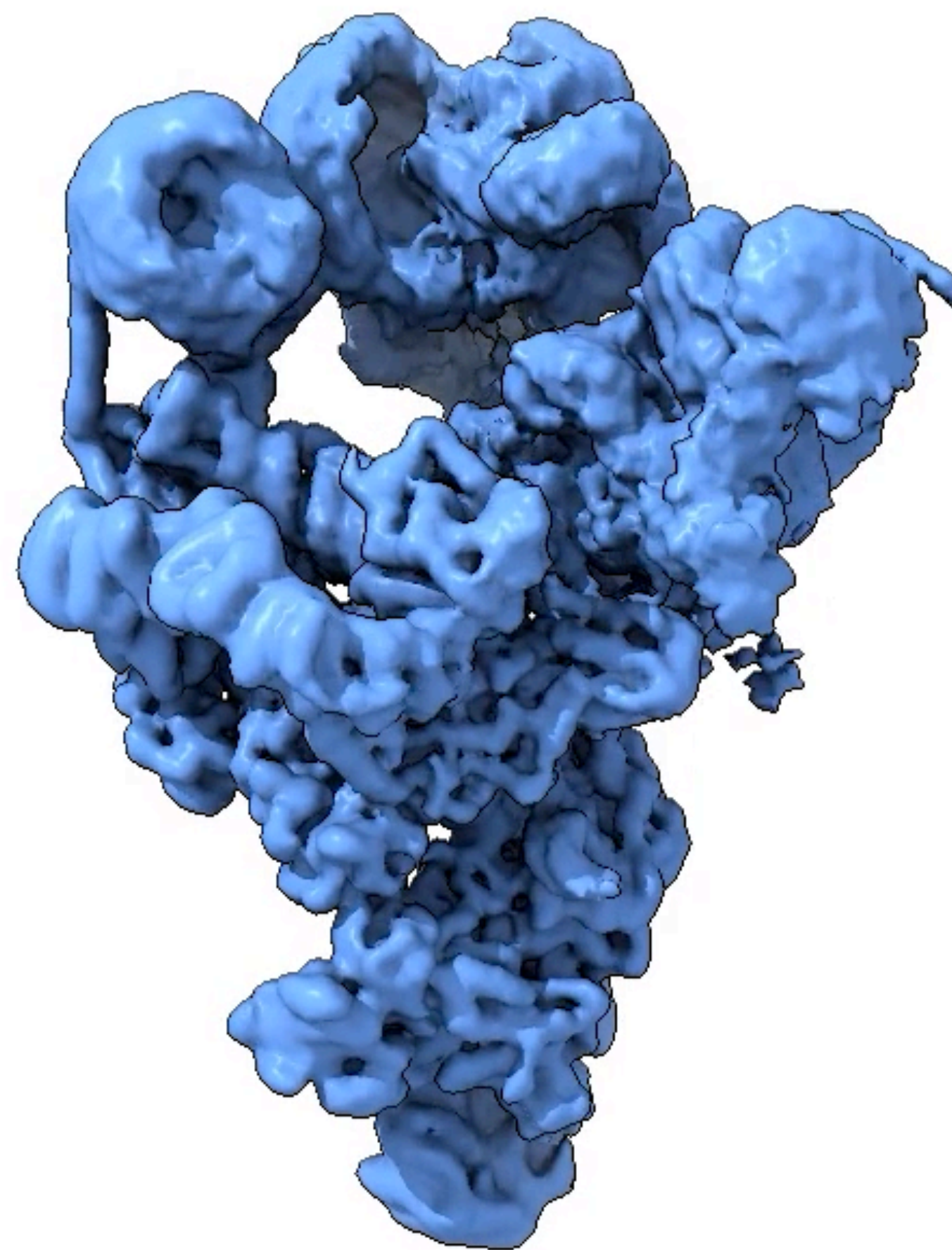
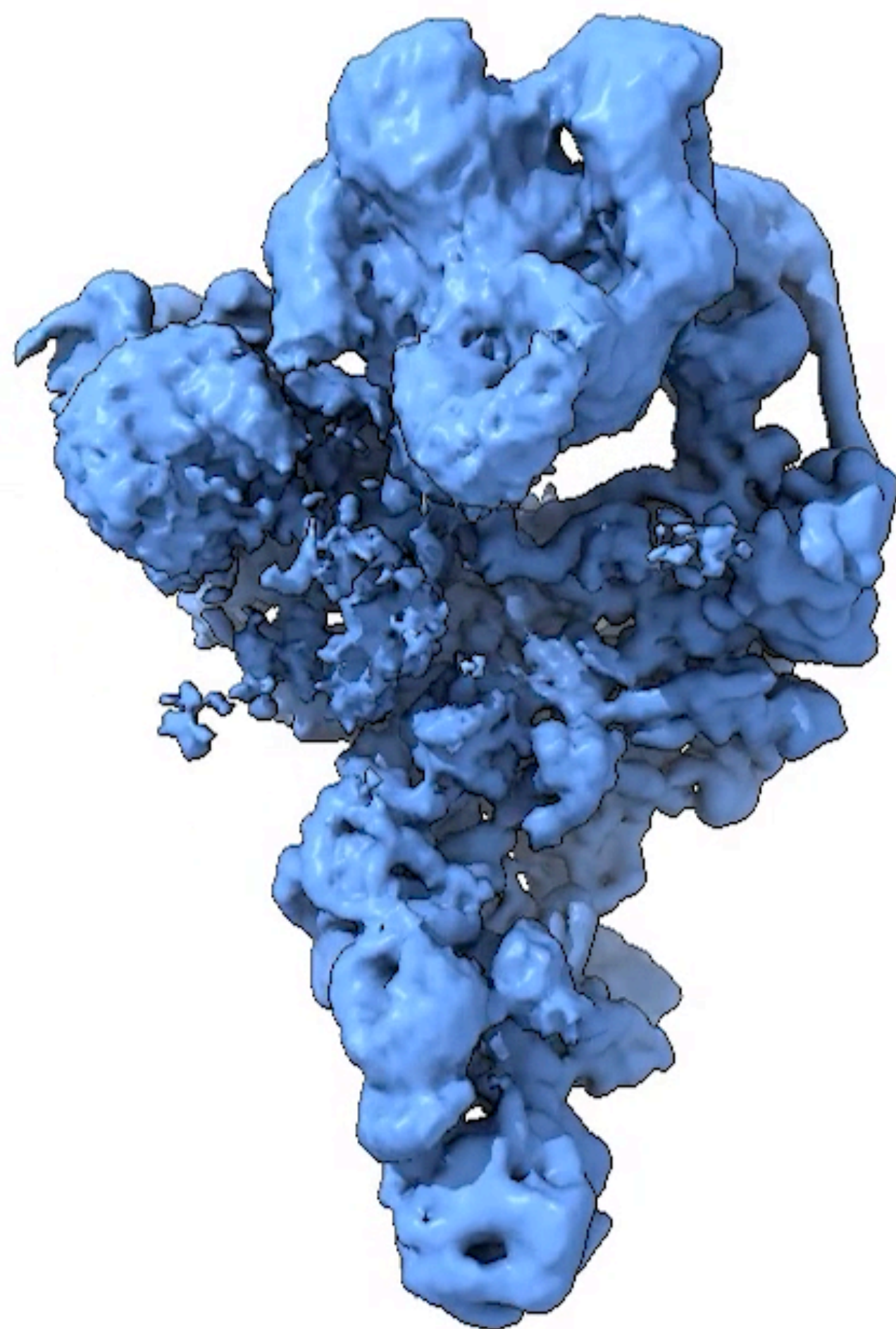
Trajectories along principle component axis of the latent space show variability within dataset



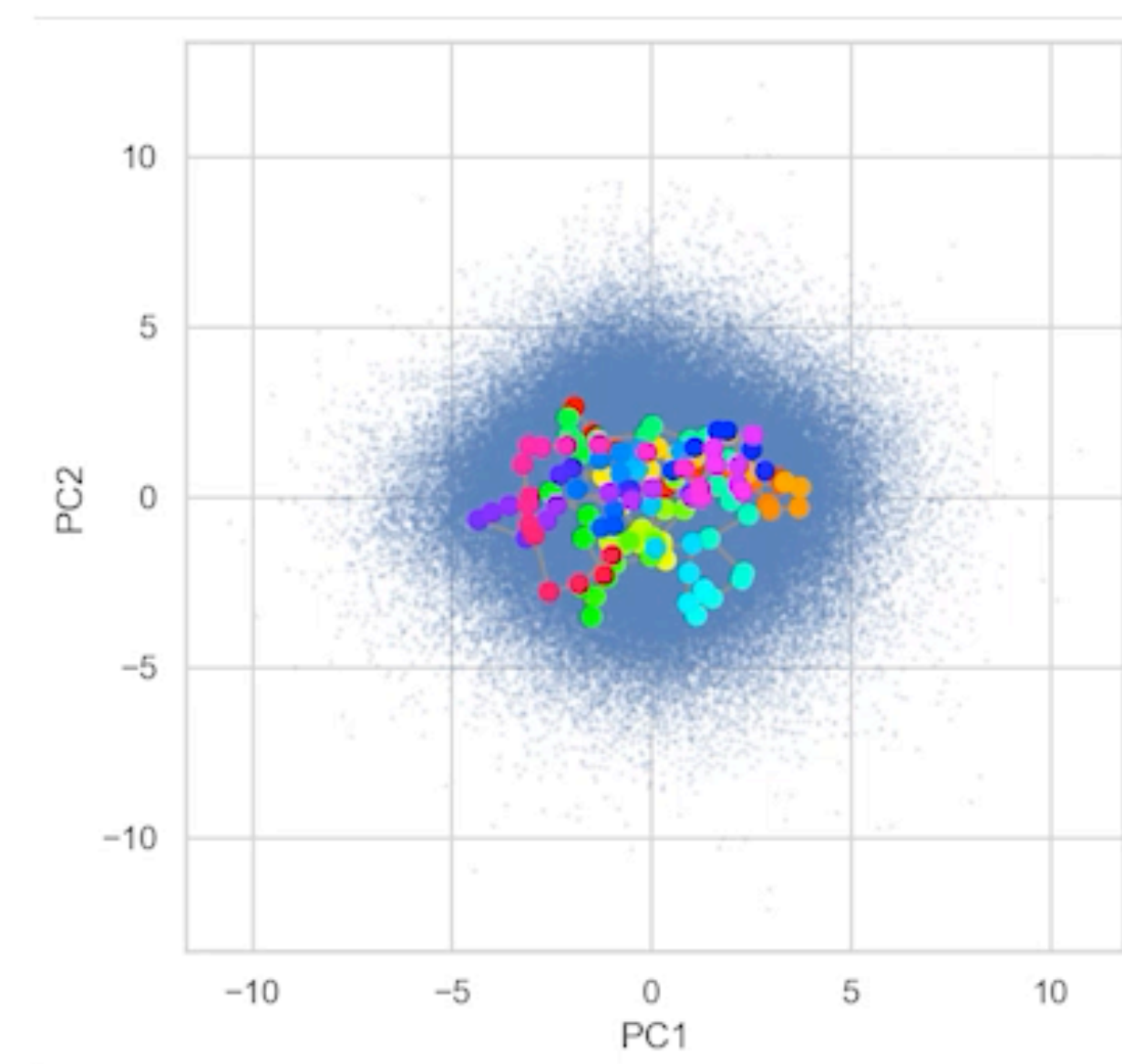
`cryodrgn pc_traversal`

Caveat: Interpolation along PCs can produce nonphysical motions e.g. under compositional heterogeneity and in general when the data distribution is not supported along the interpolation path

Generating trajectories with a graph traversal algorithm



cryodrgn graph_traversal

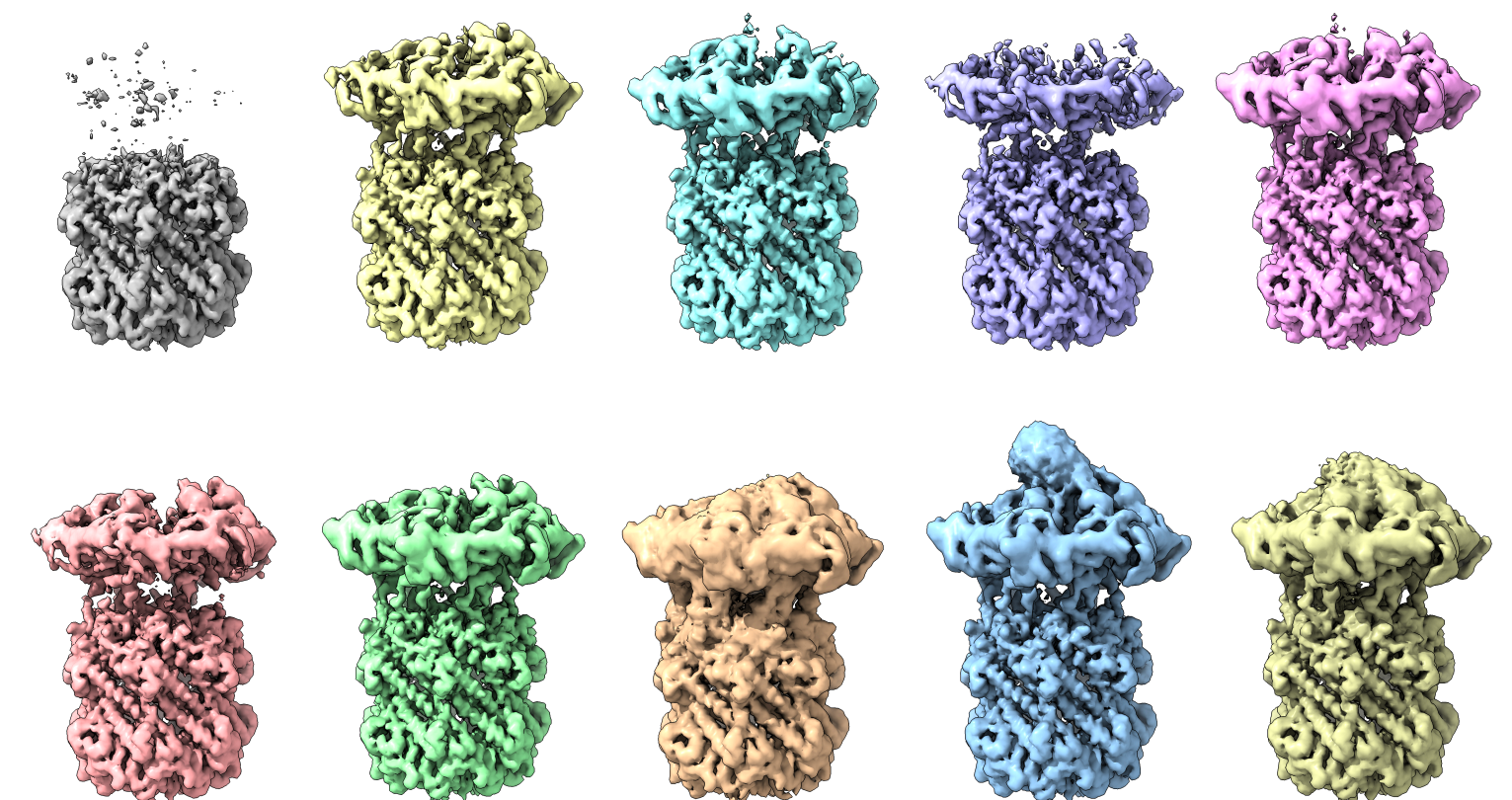
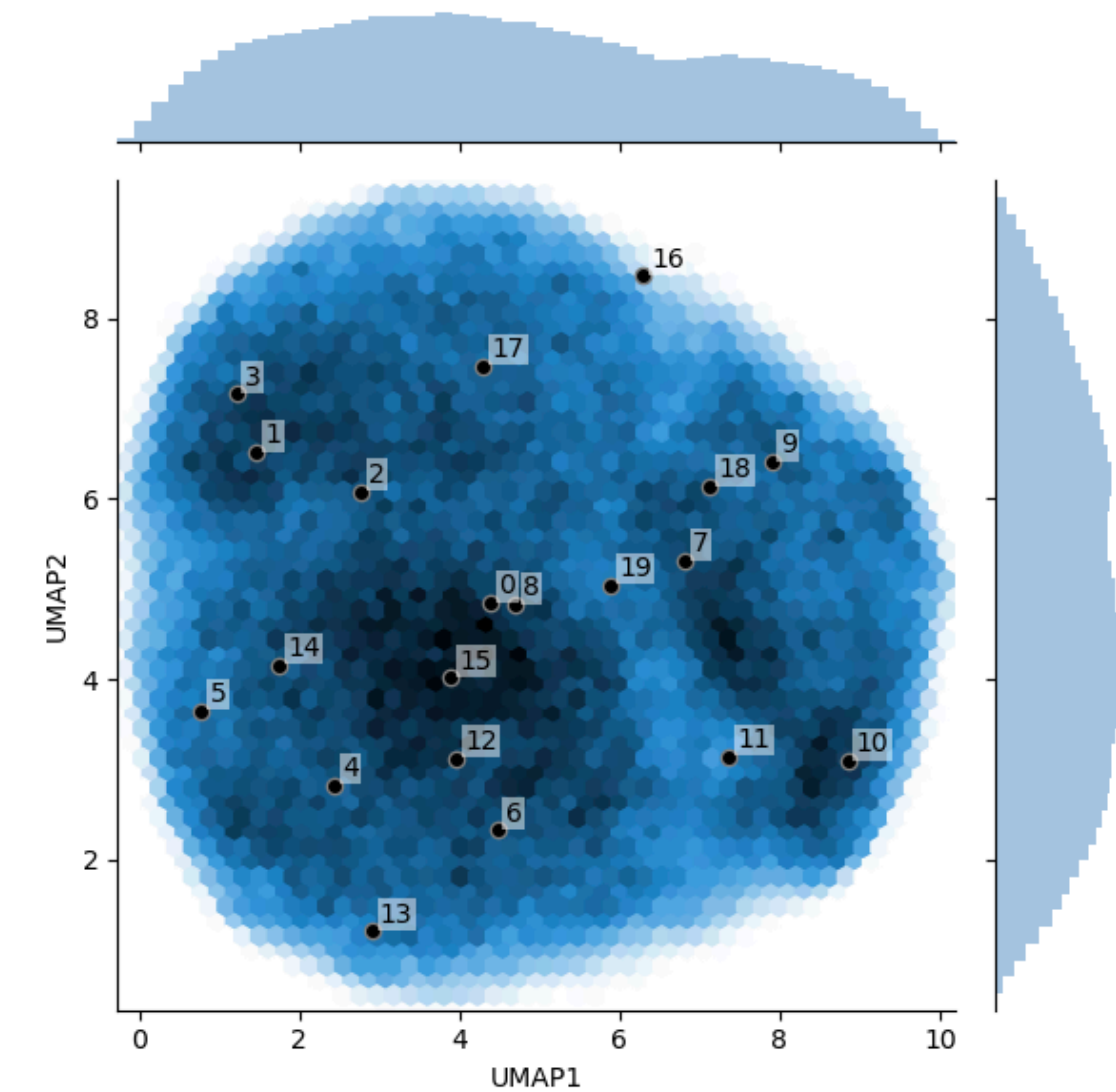


- Graph traversal algorithm along latent embedding nearest-neighbor graph
- Explore the learned distribution

CryoDRGN interactive analysis

In the generative modeling paradigm, cryoDRGN can reconstruct an arbitrary number of cryo-EM volumes.

How do we analyze the resulting ensemble of structures?



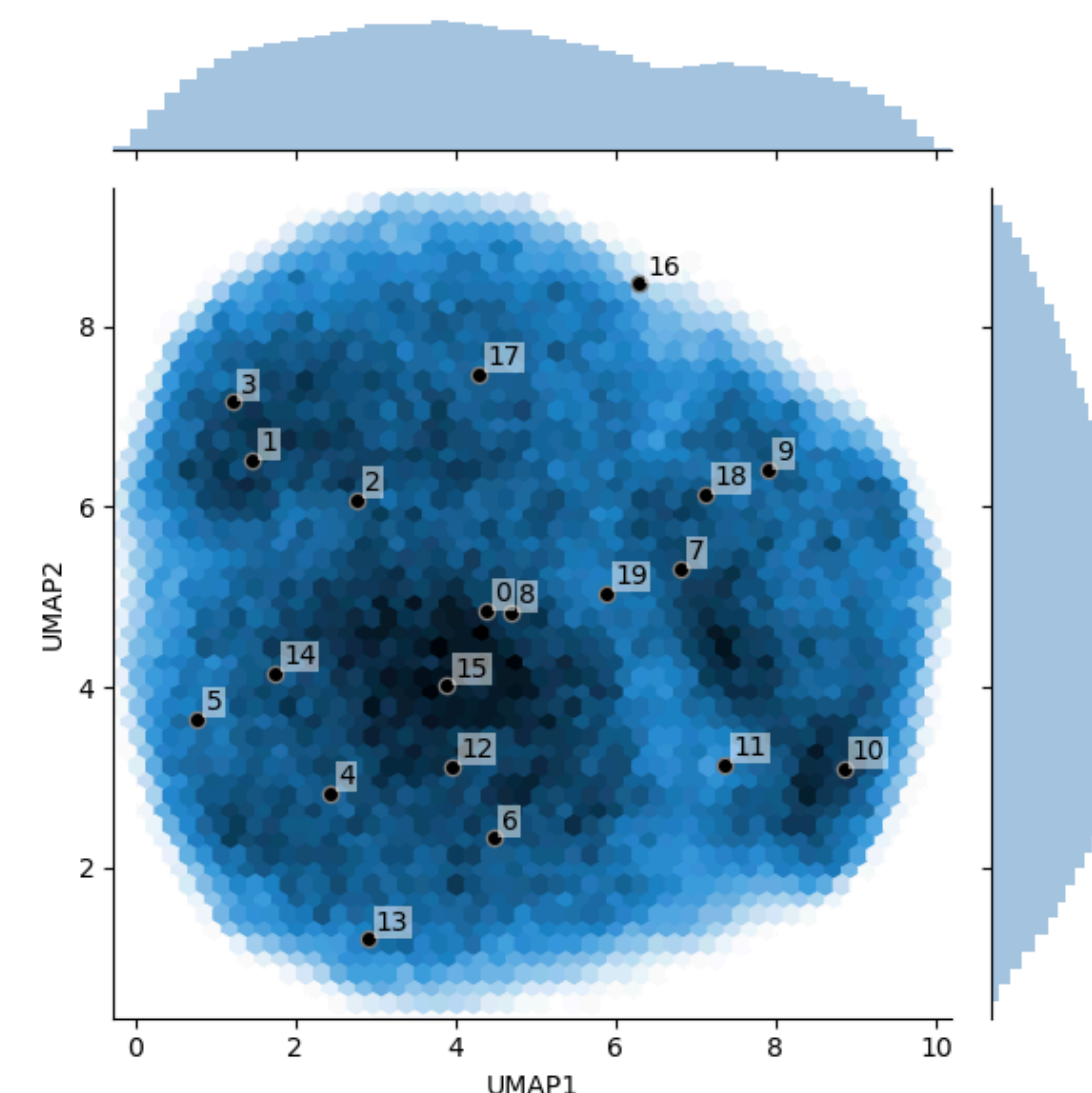
CryoDRGN interactive analysis

In the generative modeling paradigm, cryoDRGN can reconstruct an arbitrary number of cryo-EM volumes.

How do we analyze the resulting ensemble of structures?

The cryoDRGN jupyter notebook is a web application that allows **exploratory** data analysis:

- visualization of the latent space embeddings
- visualization of images
- generation of new volumes



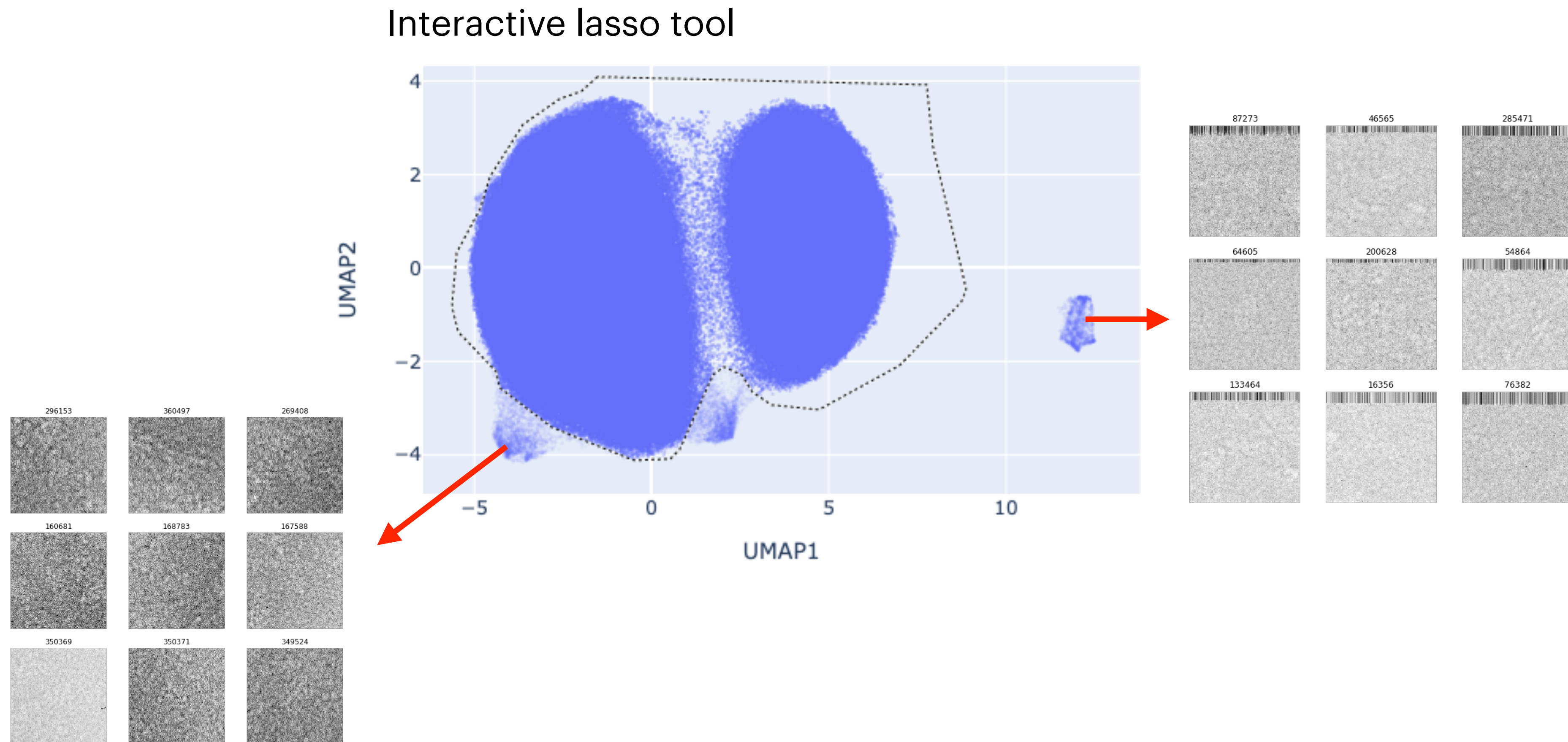
```
In [ ]: import pandas as pd
import numpy as np
import pickle
import subprocess
import os, sys

from cryodrgn import analysis
from cryodrgn import utils
from cryodrgn import dataset
from cryodrgn import ctf

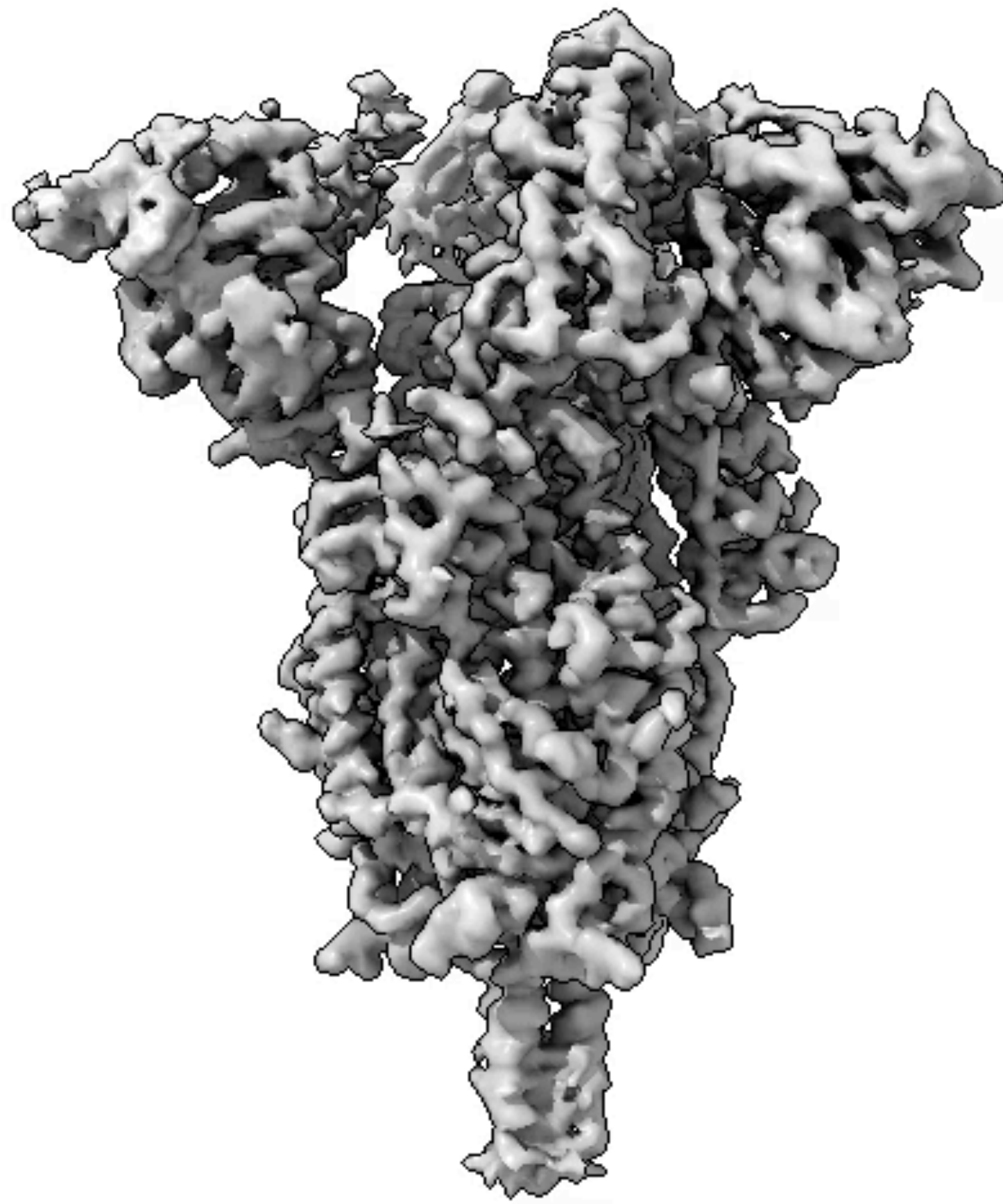
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.graph_objs as go
import plotly.offline as py
```

Interactive filtering of non-structural imaging variability

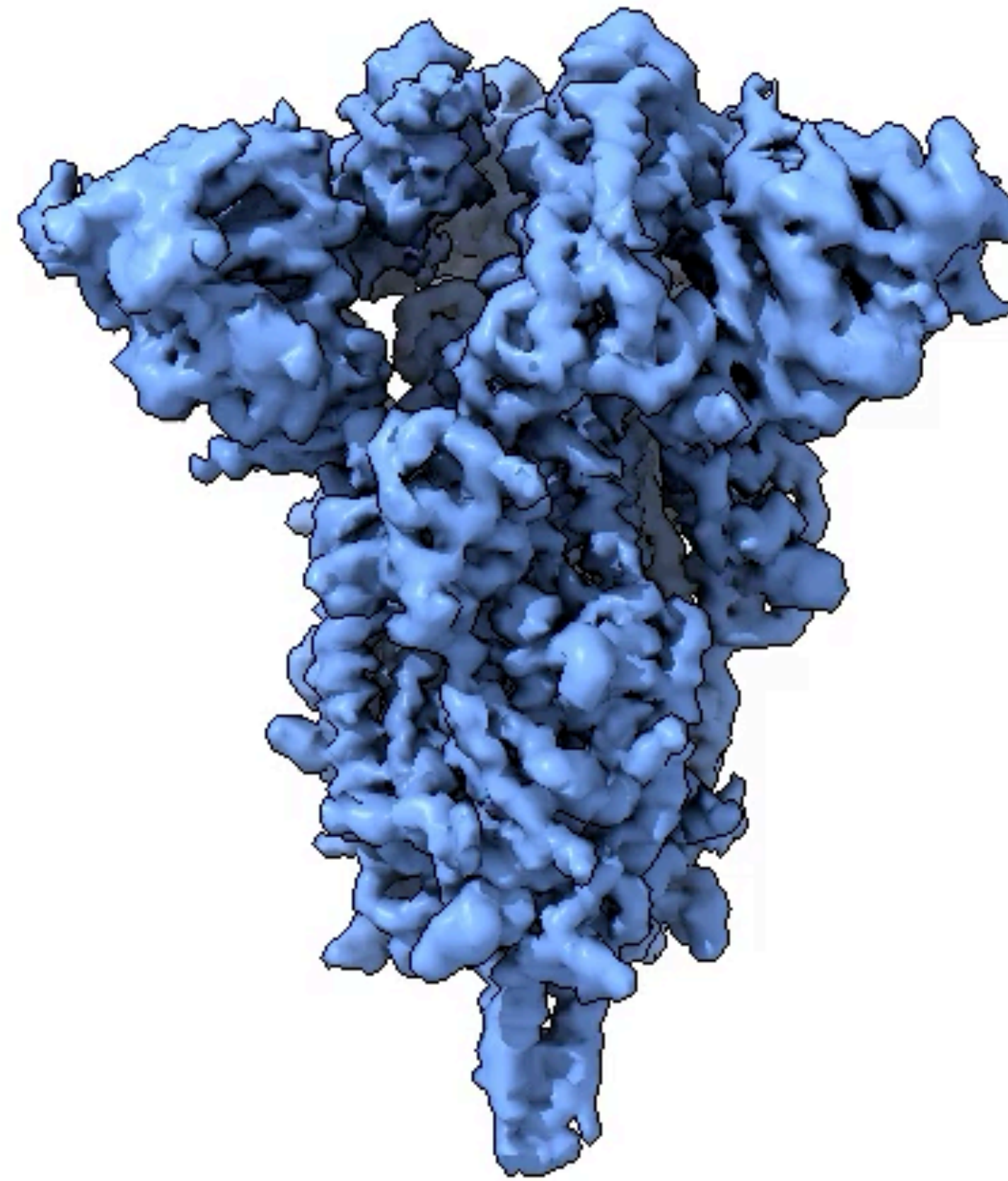
- Non-structural imaging variability (e.g. junk particles, ice artifacts, peripheral particles) may interfere with optimization and representation learning



CryoDRGN reconstruction of the SARS CoV-2 spike protein



20 sampled structures



Graph traversal trajectory

Dataset and training details:

- * Walls et al 2020
- * 276k particles
- * D=256, large architecture
- * 25 epochs
- * 4 GPU training, 10 hr total

Towards automated analysis of the structure distribution

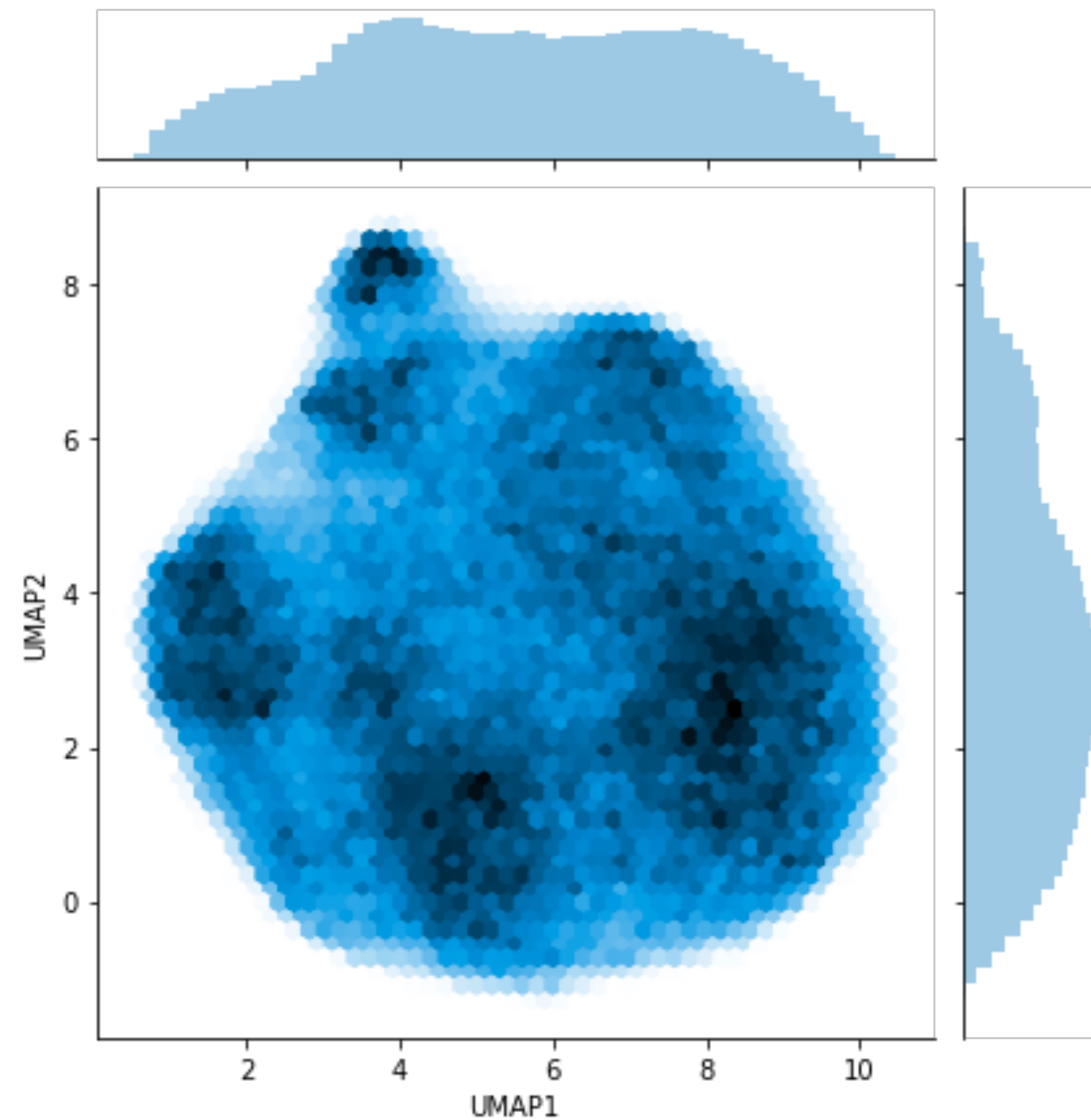
Dataset: Structural basis of ClpXP recognition and unfolding of ssrA-tagged substrates
Fei et al. 2020, eLife



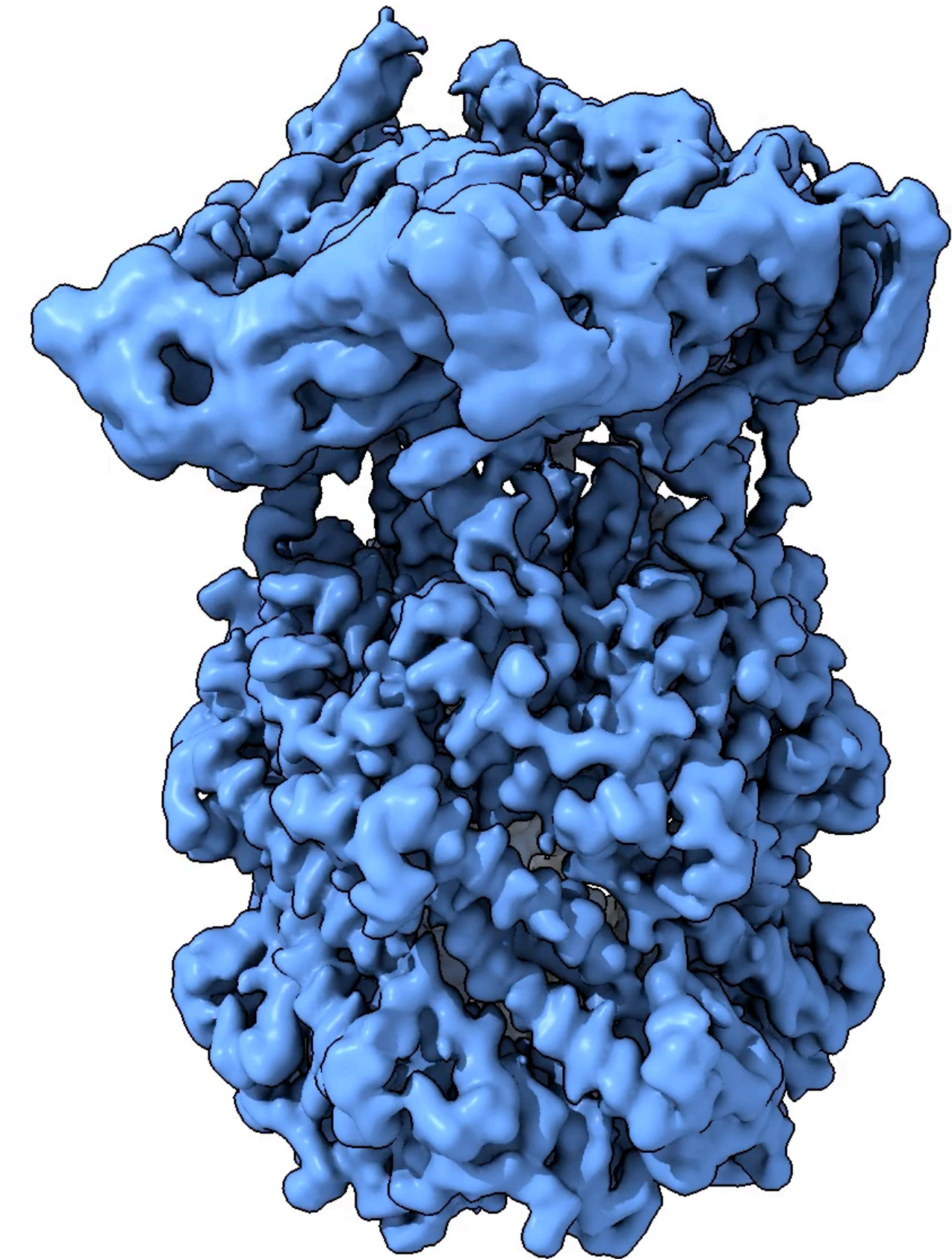
Xue Fei

Bob Sauer

cryoDRGN latent space



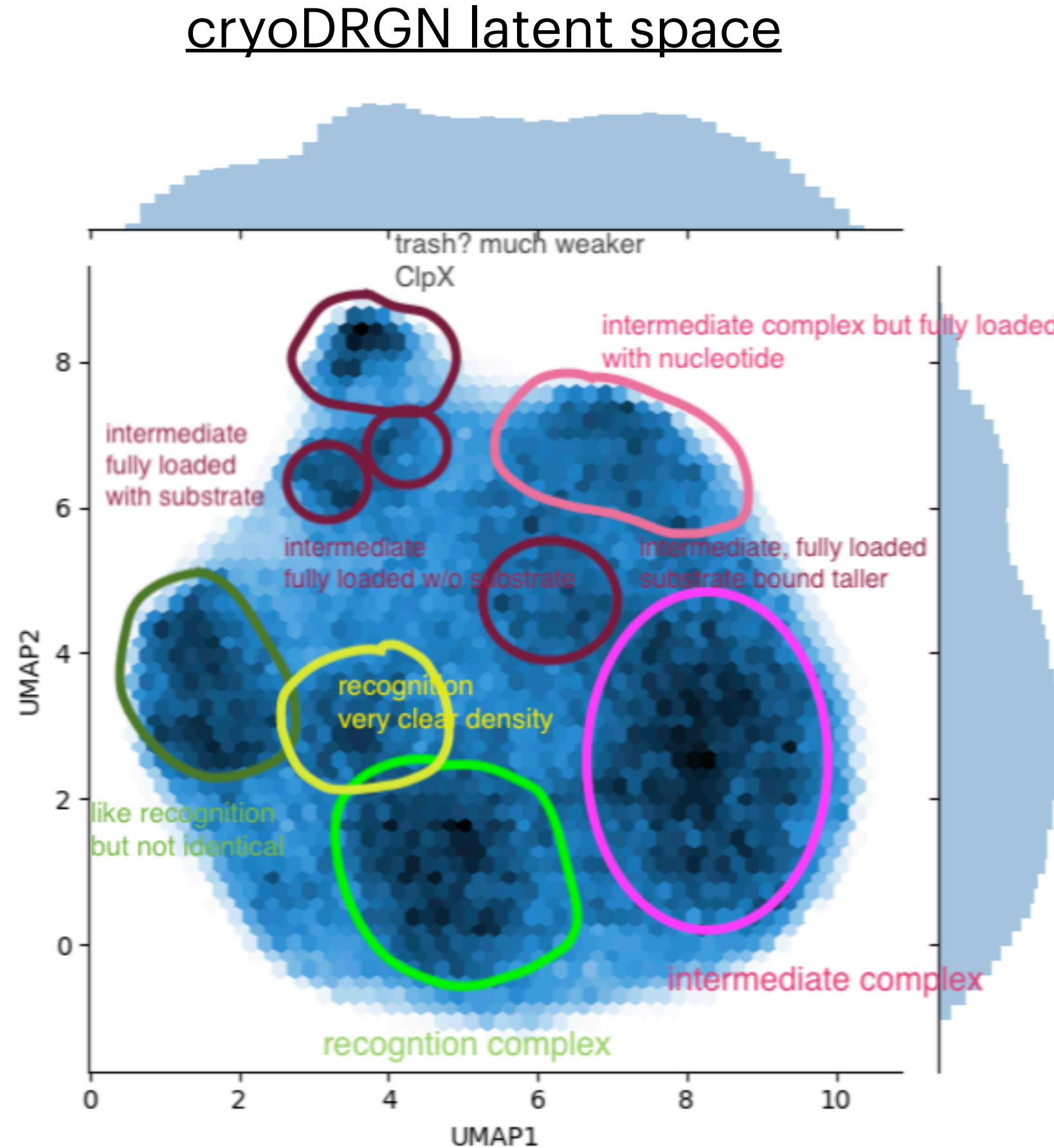
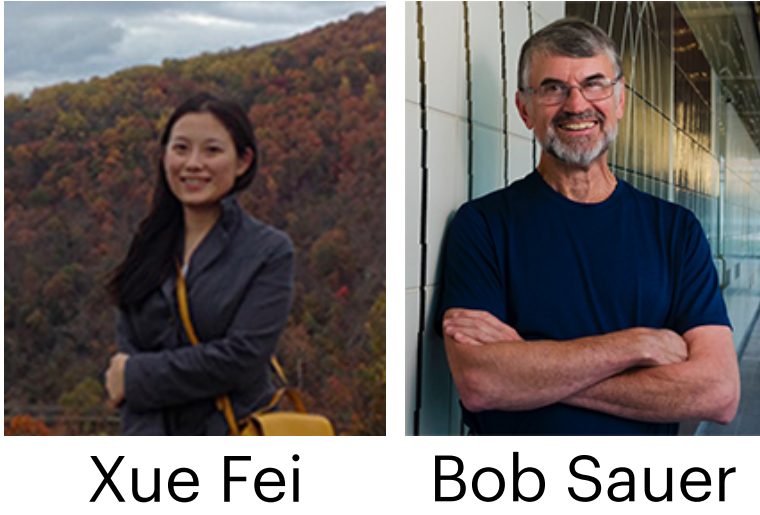
cryodrgn pc_traversal



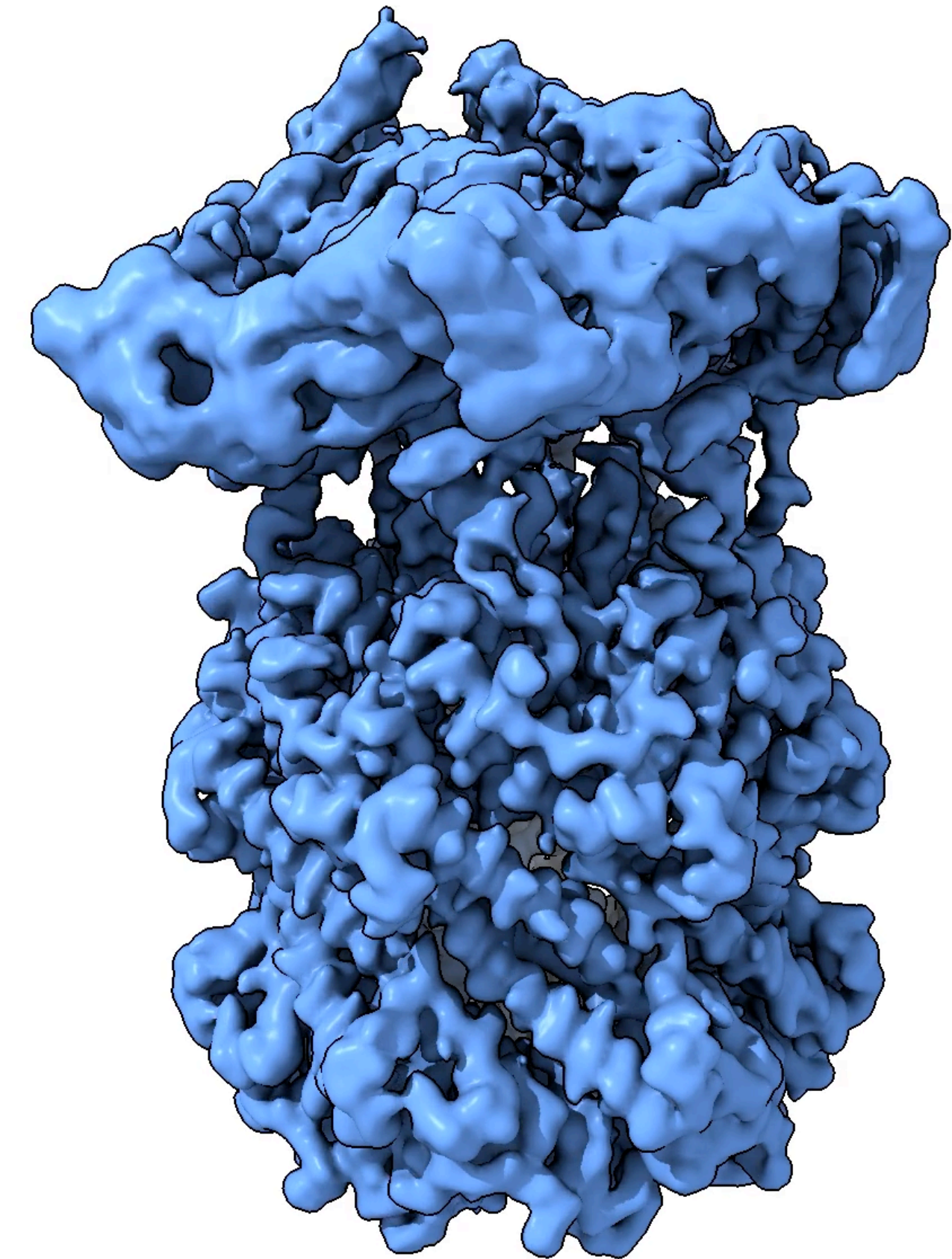
Intermediate \leftrightarrow Recognition complex

Towards automated analysis of the structure distribution

Dataset: Structural basis of ClpXP recognition and unfolding of ssrA-tagged substrates
Fei et al. 2020, eLife



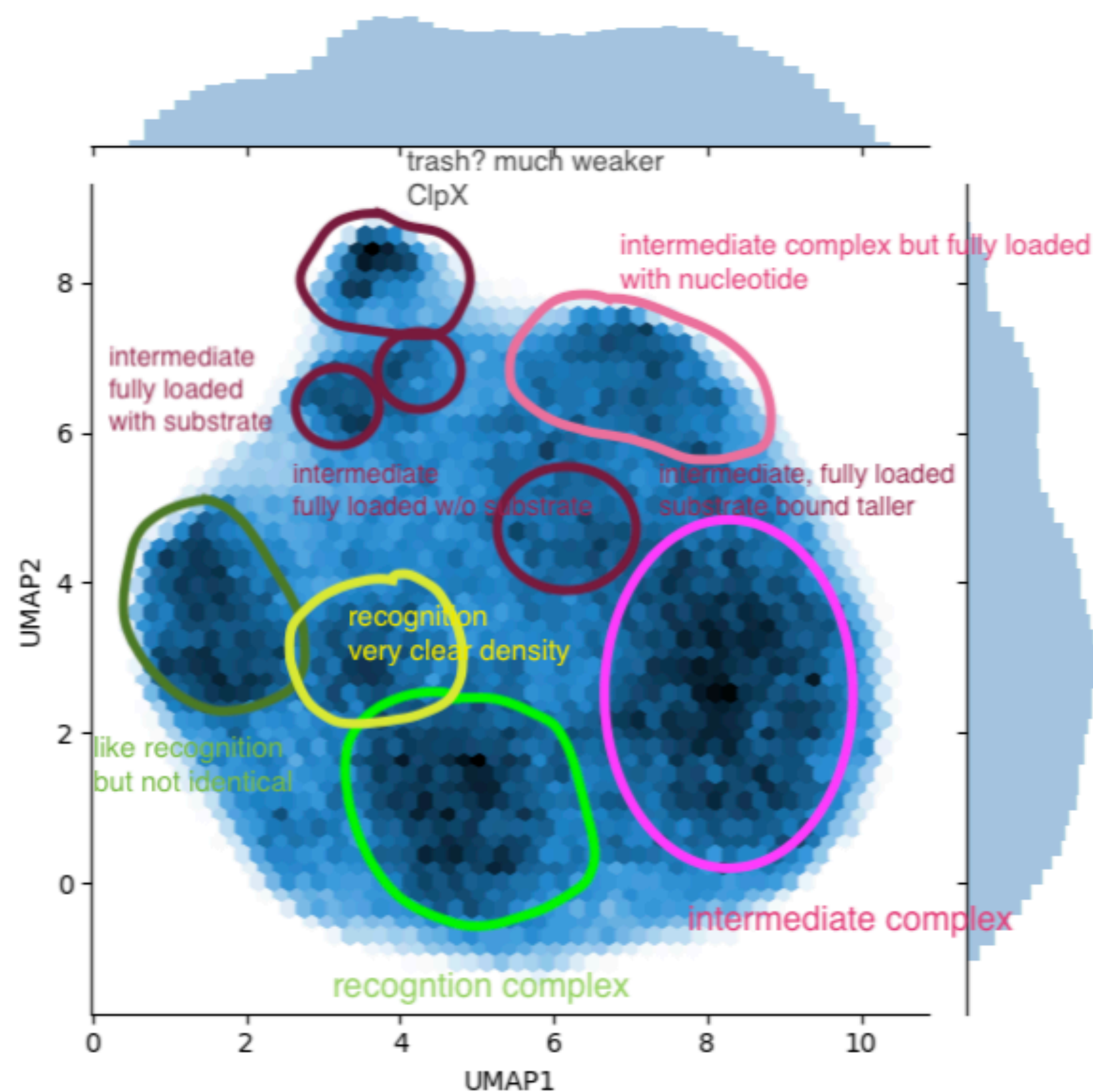
cryodrgn pc_traversal



Intermediate <-> Recognition complex

Extending the cryoDRGN toolkit with a scalable structural landscape analysis: `cryodrnn analyze_landscape`

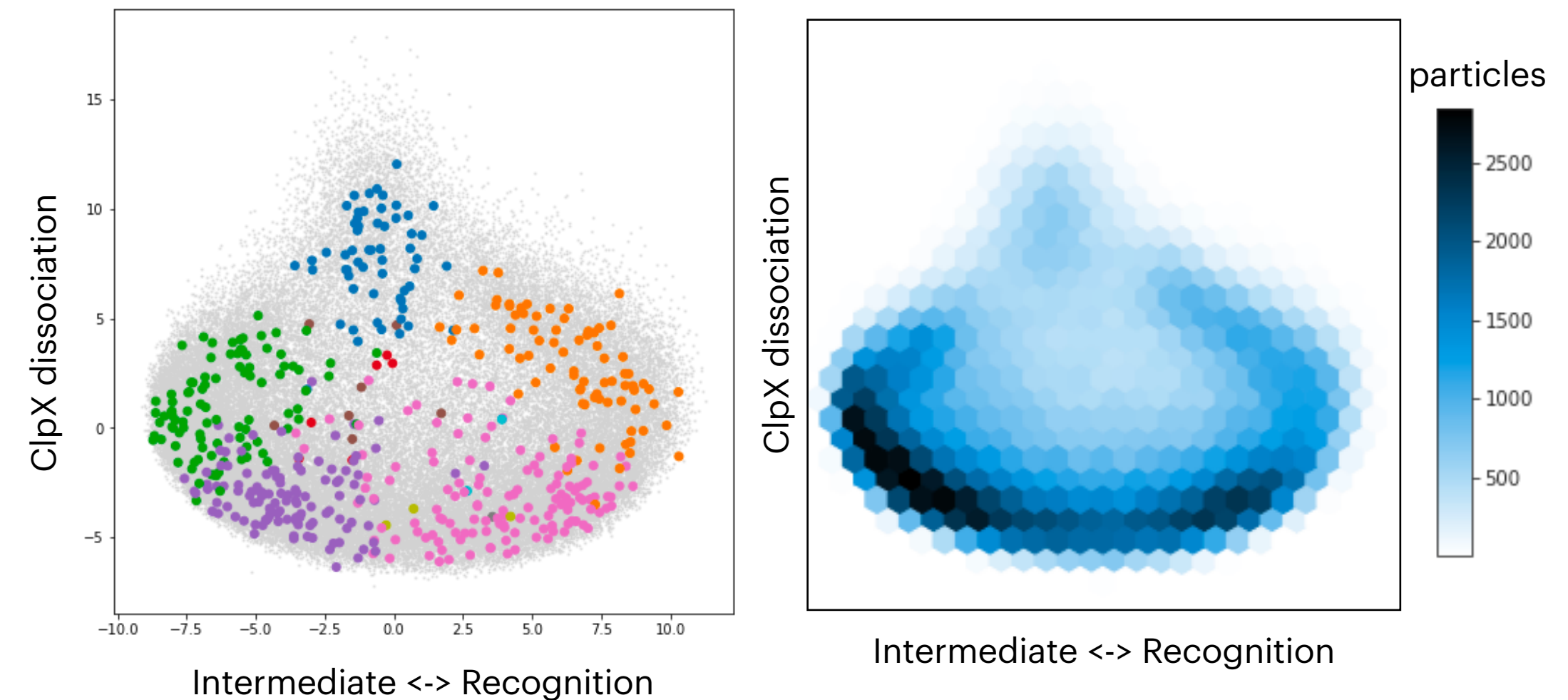
How can we gain insight from high-dimensional biological datasets?



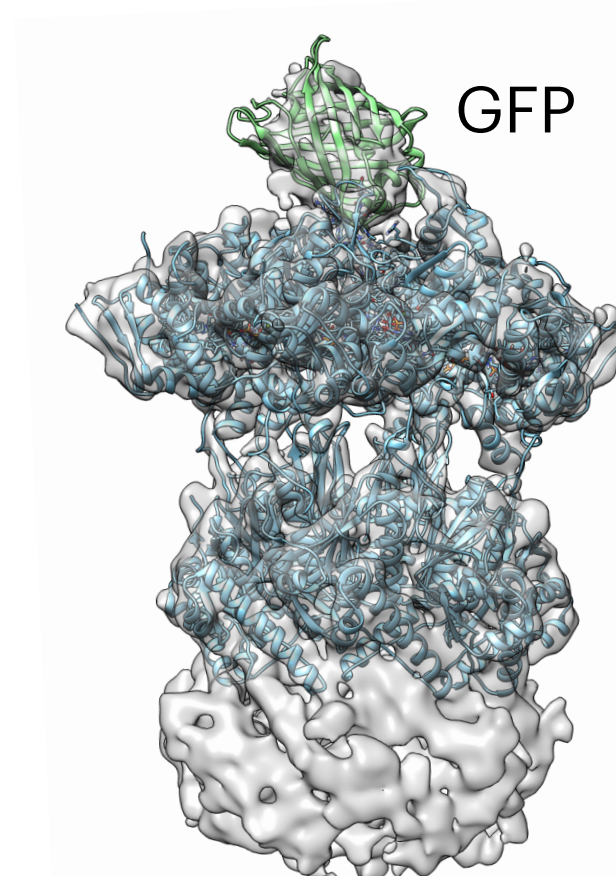
continuous

discrete

Mapping interpretable reaction coordinates



Identification of rare states

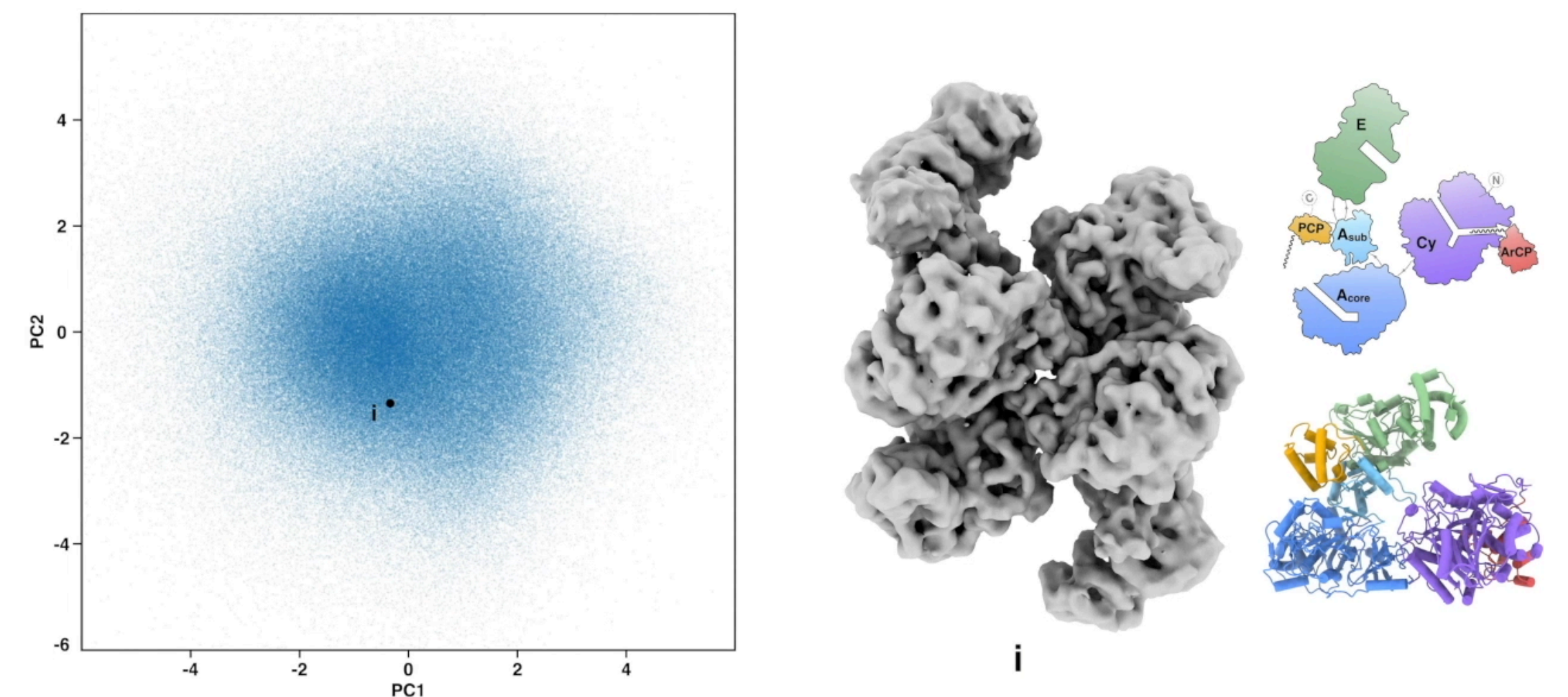
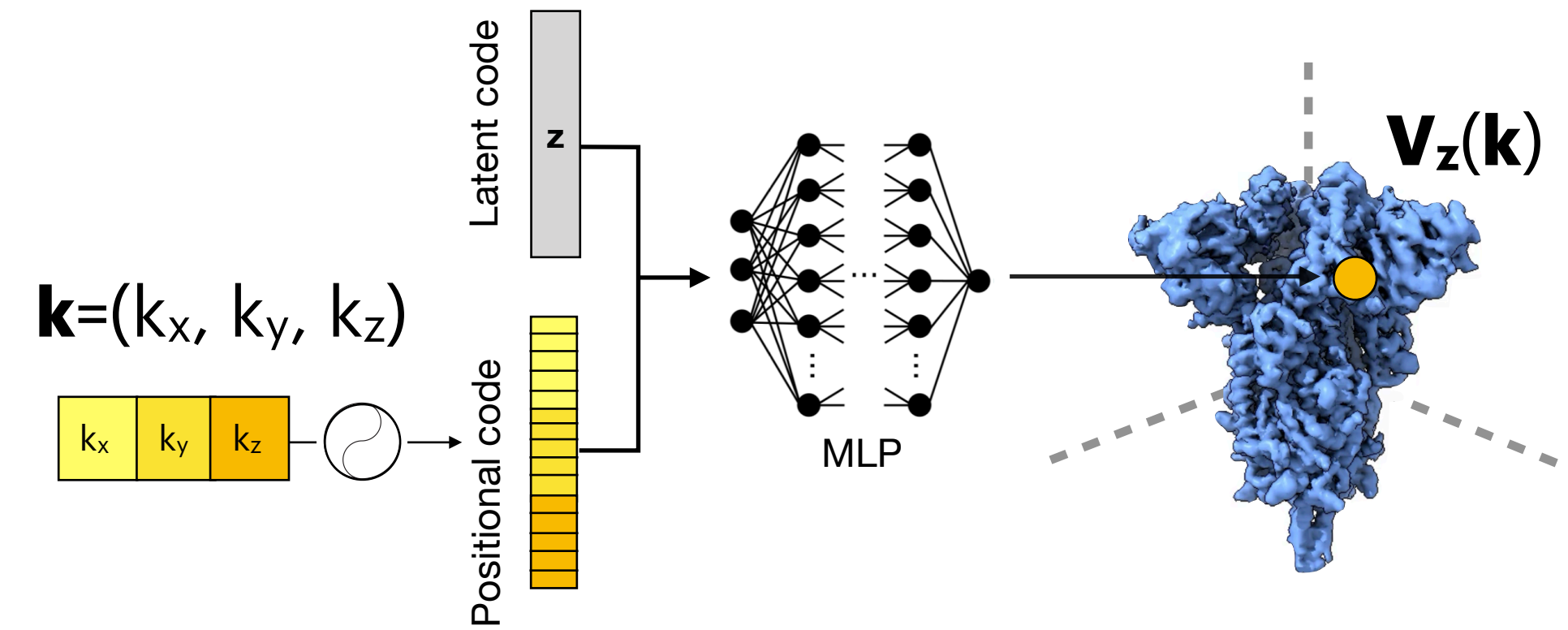


8.5 Å, 1,255 particles
0.3% of the dataset

Available in cryoDRGN 1.0

Summary

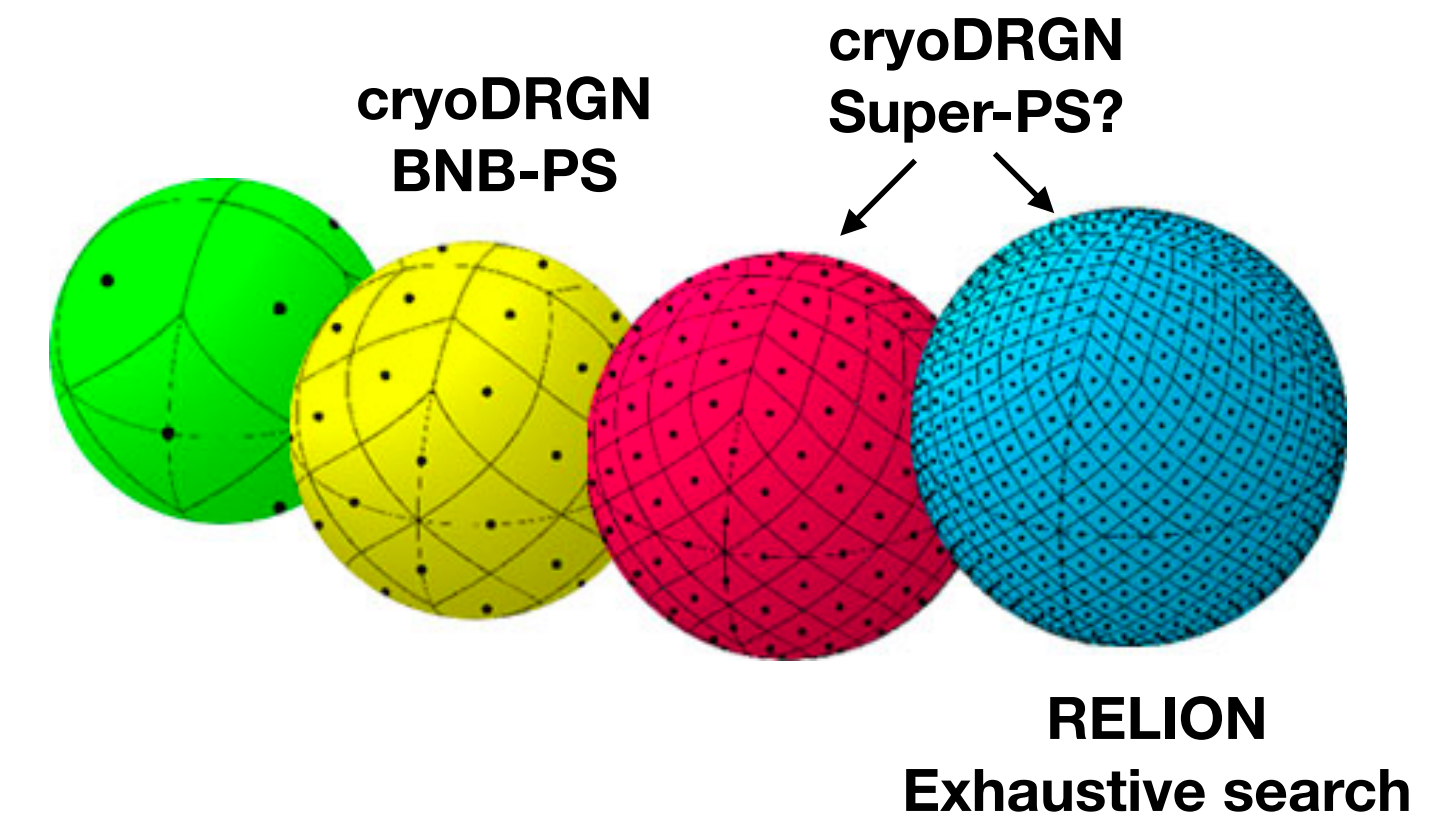
- Unsupervised reconstruction of a **continuous** distribution of protein structures from cryo-EM images
- A new neural representation for modeling high-resolution density maps
- The deep generative model provides a **general, flexible framework** for modeling heterogeneity
 - Discovery of new structures
 - Visualization of continuous dynamics
- **Novel structures** and molecular motions from cryo-EM data
- **Future outlook:** A nascent area of ML for protein structure determination



Catalytic trajectory of PchE. Wang et al 2022

A nascent area in ML algorithms for cryo-EM reconstruction

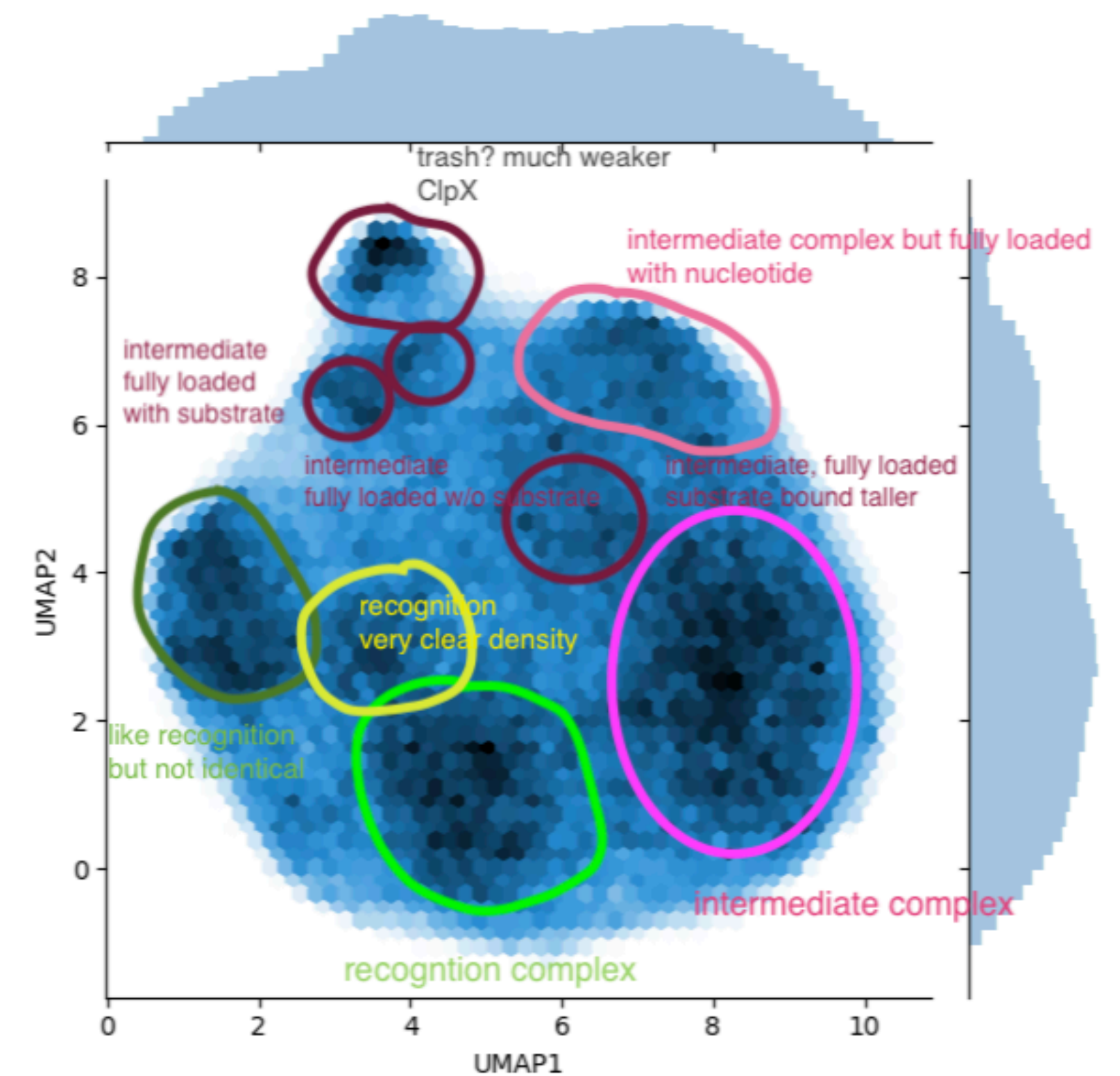
- *Ab initio* cryoDRGN on real datasets
 - Engineering, optimization, and identifiability challenges



A nascent area in ML algorithms for cryo-EM reconstruction

- *Ab initio* cryoDRGN on real datasets
 - Engineering, optimization, and identifiability challenges
- Characterizing distributions of protein structure
 - Methods for exploratory data analysis, benchmarks, and atomic modeling

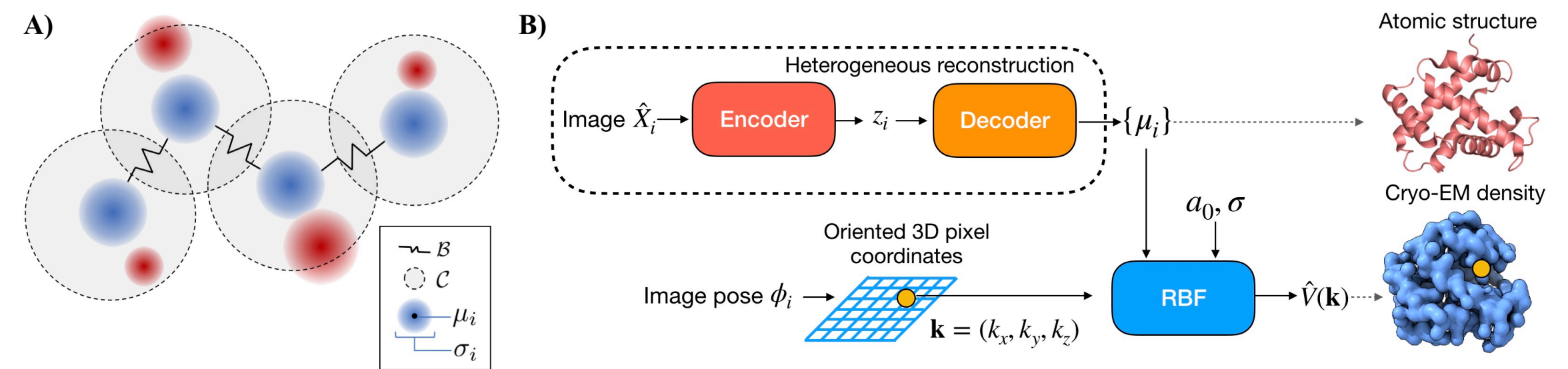
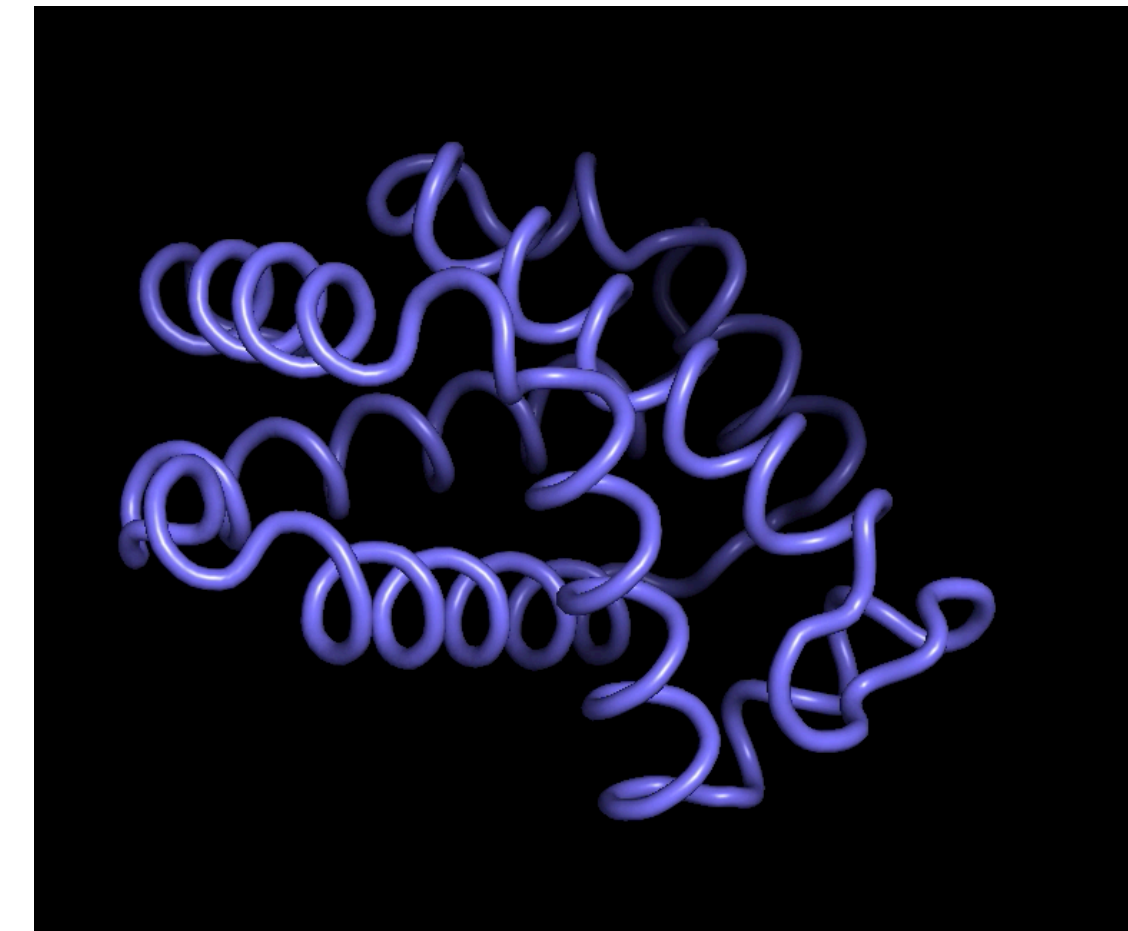
How can we gain insight from high-dimensional biological datasets?



A nascent area in ML algorithms for cryo-EM reconstruction

- *Ab initio* cryoDRGN on real datasets
 - Engineering, optimization, and identifiability challenges
- Characterizing distributions of protein structure
 - Methods for exploratory data analysis, benchmarks, and atomic modeling
- New representations and generative modeling paradigms
 - Better inductive biases for protein motion/dynamics; Exploiting information from structure/sequence databases

Radial basis function representation:

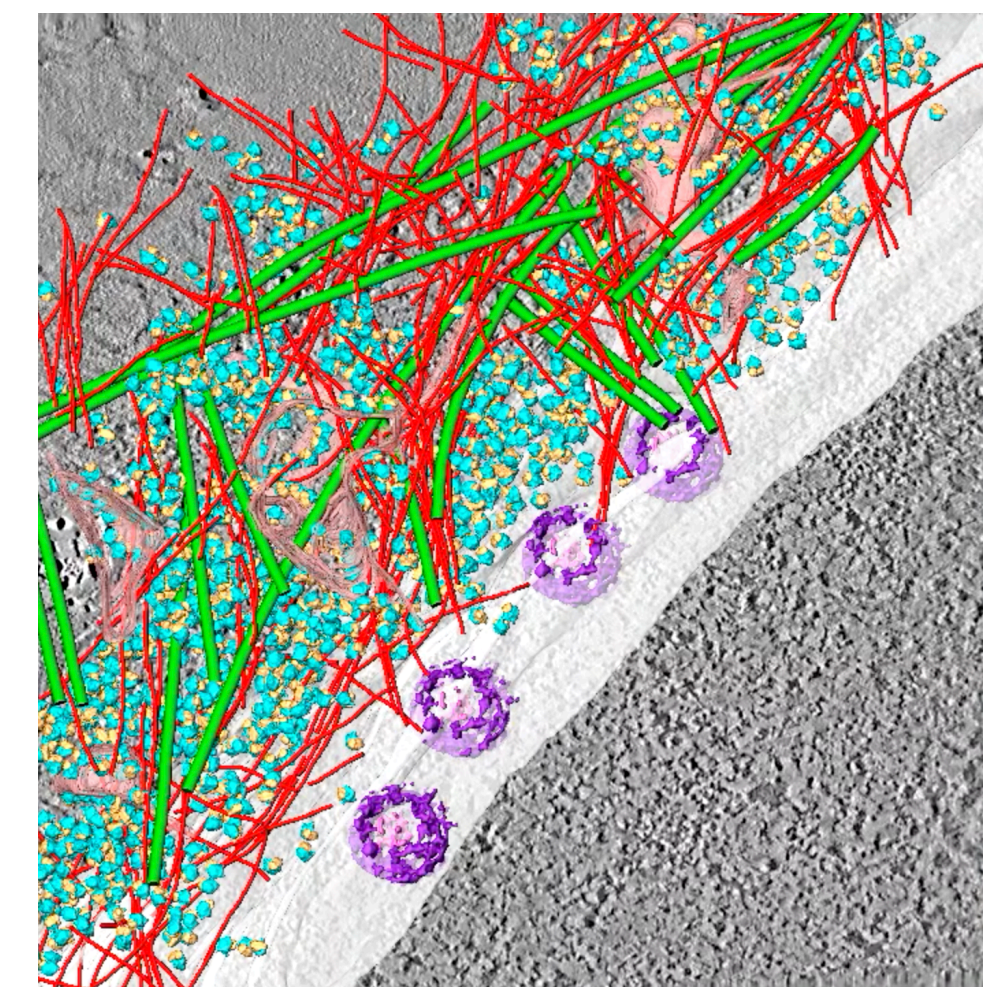


A nascent area in ML algorithms for cryo-EM reconstruction

- *Ab initio* cryoDRGN on real datasets
 - Engineering, optimization, and identifiability challenges
- Characterizing distributions of protein structure
 - Methods for exploratory data analysis, benchmarks, and atomic modeling
- New representations and generative modeling paradigms
 - Better inductive biases for protein motion/dynamics, Exploiting information from structure/sequence databases
- *In situ* cryoDRGN
 - Towards in situ structural biology with cryo-electron tomography (cryo-ET)



<https://pdb101.rcsb.org/sci-art/goodsell-gallery/escherichia-coli-bacterium>

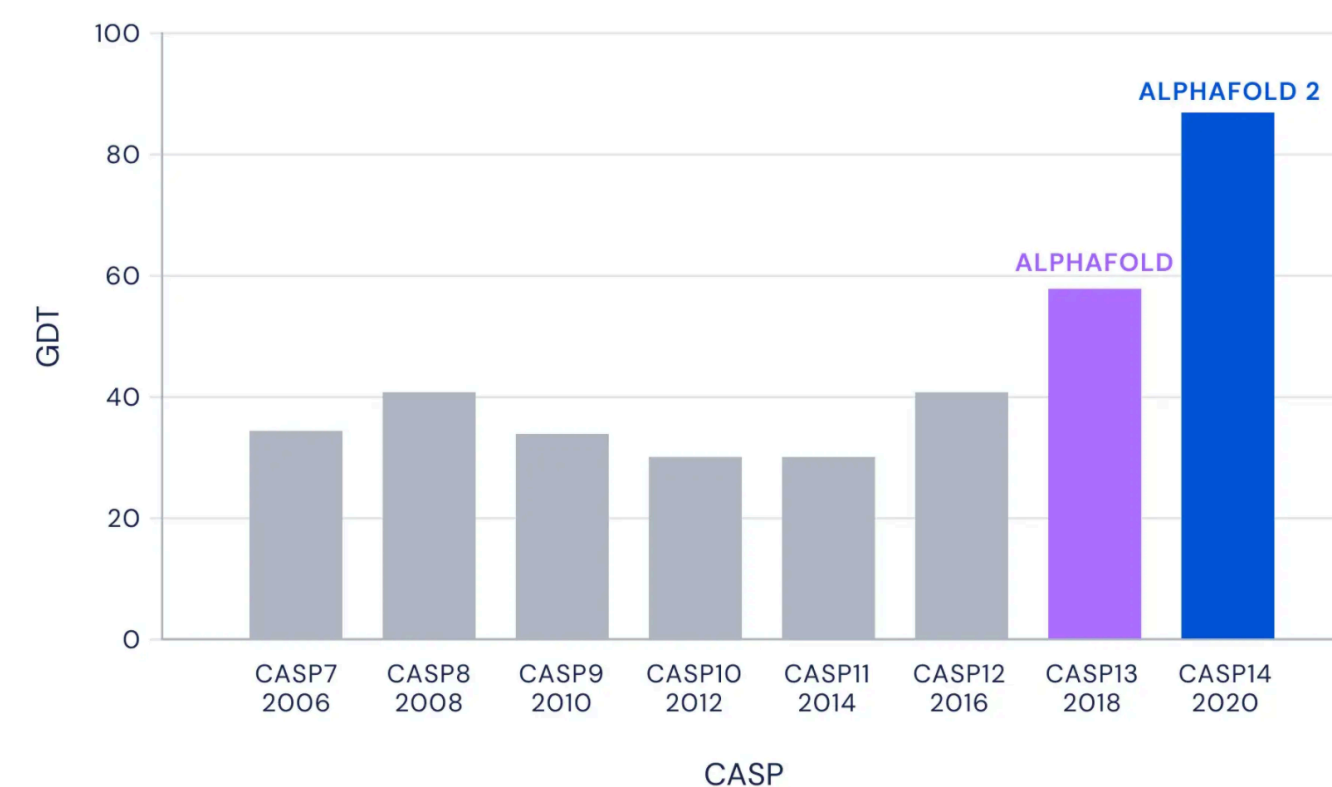


Visualizing the molecular sociology at the HeLa cell nuclear periphery
Mahamid et al, Science 2016

A nascent area in ML algorithms for cryo-EM reconstruction

- *Ab initio* cryoDRGN on real datasets
 - Engineering, optimization, and identifiability challenges
- Characterizing distributions of protein structure
 - Methods for exploratory data analysis, benchmarks, and atomic modeling
- New representations and generative modeling paradigms
 - Better inductive biases for protein motion/dynamics, Exploiting information from structure/sequence databases
- *In situ* cryoDRGN
 - Towards in situ structural biology with cryo-electron tomography (cryo-ET)
- **Outlook in the post-AlphaFold2 era?**

Median Free-Modelling Accuracy



Thank you for listening!

E.Z. Lab

- [Ramya Rangan](#)
- [Axel Levy](#)
- Rish Raghu
- Ryan Feathers
- Vineet Bansal

Stanford

- Fred Poitevin
- Gordon Wetzstein

NVIDIA

- Tim Dockhorn
- Karsten Kreis

Flatiron Institute

- Pilar Cossio
- Sonya Hanson

ET/Thermo Fischer

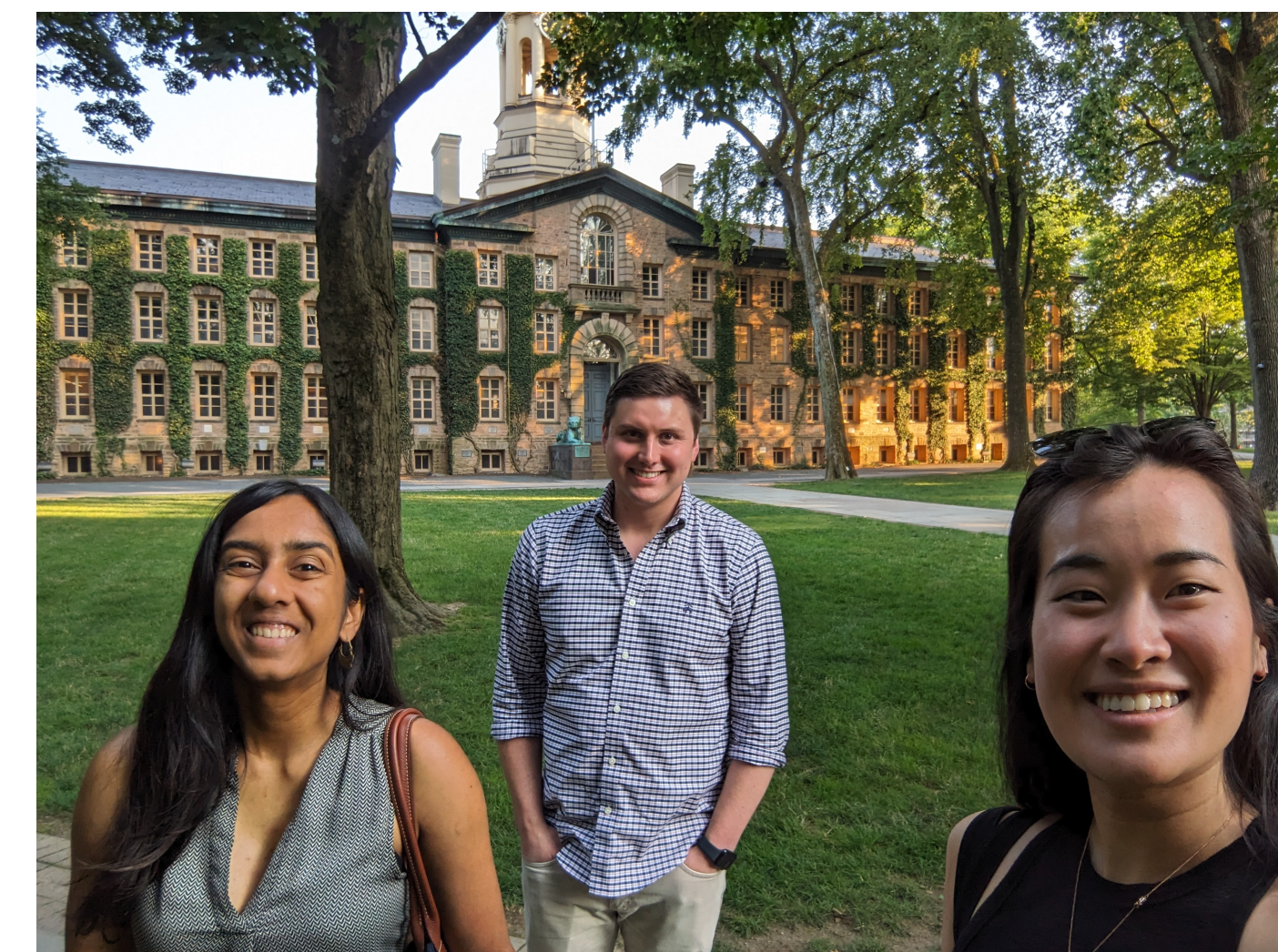
- Sagar Svenekar (MPI Biochemistry)
- Jake Johnston (Columbia/NYSBC)
- Adam Lerer (DeepMind)
- Martin Obr (TF)
- Ron Kelley (TF)
- Abhay Kotecha (TF)

CryoDRGN1

- Tristan Bepler
- Bonnie Berger (MIT Mathematics)
- Joey Davis (MIT Biology)
- Ashwin Narayan
- Barrett Powell
- Laurel Kinman



EZ Lab at NeurIPS Machine Learning for Structural Biology Workshop



PRINCETON
UNIVERSITY