

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044

Cryo2SSNet: 2-Stage Search Supervised Network for Cryo-electron Microscopy Single-Particle Reconstruction

Bingliang Zhang and Zitian Tang

IIIS, Tsinghua University, China
`{zhangbl19, tzt19}@mails.tsinghua.edu.cn`

Abstract. Cryogenic electron microscopy (cryo-EM) single-particle reconstruction is a powerful technique to elucidate the 3D structure of proteins and other biomolecules from a set of noisy and arbitrarily oriented 2D projection images. This paper presents *2-stage Search Supervised Network* (Cryo2SSNet), an auto-encoder architecture enhanced by a beam search module to jointly infer orientations and reconstruct a 3D density map from a random initialization in a 2-stage pipeline. In the framework, the encoder is to predict orientations for individual images while the decoder is to simulate the cryo-EM imaging process by directly maintaining the 3D density map of the target biomolecule. To stabilize predicted orientations, we propose *denoised heuristic beam search* (DHBS) to provide an explicit weak-supervision for encoder and refined orientations for reconstruction. We evaluate our method on several simulated protein datasets with different signal-to-noise ratios (SNR) and demonstrate the high precision for 3D density map reconstruction.

Keywords: Cryo-EM, Tomographic reconstruction, Auto-encoder

1 Introduction

Biomolecules perform most of the activities in living cells. As the function of a biomolecule is highly correlated to its structure, in order to investigate its activity, it is a key step to determine its structure. The most commonly used ways to do this include X-ray crystallography, nuclear magnetic resonance, and cryogenic electron microscopy (cryo-EM) [19]. With the long-term development, cryo-EM gains the ability to determine biomolecular structures at near-atomic resolution [12]. Though the structure resolution derived by cryo-EM cannot exceed that of X-ray crystallography, combining the cryo-EM structure with the chemical structure of the biomolecule, it is able to produce an informative atomic model or even directly locate every individual atom [19].

To determine the structure of a biomolecule using cryo-EM, about $10^{4\sim 7}$ copies of this biomolecule are ice-embedded to a thin vitreous slice. Cryo-EM imaging is then performed over these randomly oriented particles and results in a micrograph, containing the projections of each particle. These projections are

045 picked from the micrograph and represented as equal-size square images. The
046 projection of a particle can be approximated by an integral of its electron density
047 map, convoluted with the point-spread function (PSF) of the cryo-EM. The final
048 step is to reconstruct the density map based on these projection images, known
049 as tomographic reconstruction. In this procedure, each derived image is of a
050 single distinct particle, hence this process is called single-particle reconstruction
051 (SPR). This allows cryo-EM to capture the heterogeneous conformations of the
052 biomolecule, making cryo-EM advantageous over others. But in many works,
053 including this one, a unique conformation assumption is made to simplify the
054 problem.

055 The unknown particle poses and extremely strong noises make the tomo-
056 graphic reconstruction challenging. For each projection image, the orientation
057 of the captured particle is unknown, and there might be an unknown transla-
058 tional shift of particle center due to the inaccuracy of particle picking. Besides,
059 the major part of the signal detected by cryo-EM is from the surrounding ice
060 and other noisy sources, rather than the particles. This causes an extremely low
061 signal-to-noise ratio (SNR) of the captured images.

062 Many approaches have been proposed to solve the tomographic reconstruc-
063 tion problem, including classical ones used in current cryo-EM software pack-
064 ages, and those using deep learning emerging in recent years. In this work, a new
065 deep-learning-based method, *2-stage search supervised network* (Cryo2SSNet), is
066 developed towards tomographic reconstruction in cryo-EM. Specifically, it is an
067 auto-encoder-based model with a search module, whose encoder is a neural net-
068 work inferring the particle poses, and decoder is a cryo-EM simulator. To balance
069 pose estimation accuracy and efficiency in strongly noisy conditions, we propose
070 *denoised heuristic beam search* (DHBS) strategy as a search module. Besides,
071 our method adopts a 2-stage training pipeline, in the first of which the encoder
072 is weak-supervised by the search module and in the second of which the model
073 is trained in an end-to-end manner. Cryo2SSNet is tested on various simulated
074 protein datasets with low SNR and outperforms previous deep-learning-based
075 methods in reconstruction resolution.

076 2 Related Works

077 The cryo-EM SPR problem has been explored since decades ago. While most
078 current cryo-EM software packages are using classical methods like expectation-
079 maximization (EM) algorithm, some deep learning methods are proposed re-
080 cently to tackle this problem. These two categories are introduced here respec-
081 tively. Besides, projection-matching, a technique used in a class of methods,
082 which is improved in this work, is discussed in details.

083 2.1 Current Methods

084 A major factor making tomographic reconstruction hard is the unknown parti-
085 cle poses. Once the poses are provided, the density map can be reconstructed

by maximum-likelihood (ML) estimation accurately. Hence a straightforward and classical solution is to estimate the particle poses, among which projection-matching-based methods [14,15,1] are popular. Pose estimation and map refinement are performed iteratively. However, the extremely low SNR makes the estimation inaccurate, leading to bad reconstruction.

Instead of finding a certain pose, another class of methods [17,16] maintains a pose distribution for each particle and performs maximum-a-posteriori (MAP) optimization. EM algorithm is used to solve this problem, where E-step updates the pose distribution and M-step updates the density map. These two steps are performed iteratively until convergence. CryoSPARC [16] is one of the most widely used softwares. It proposes to solve the MAP via stochastic gradient descent (SGD) to derive a solution *ab initio*, and performs a subsequent refinement.

2.2 Deep Learning for Cryo-EM

Deep learning, having achieved state-of-the-art performances in a wide range of computer vision tasks, is already introduced into the cryo-EM SPR process, such as being used to denoise the micrographs [2] and to pick particles from the micrographs [23,26,21,22,3].

Works applying deep learning to tomographic reconstruction are proposed in recent years. [9] uses convolutional neural networks (CNNs) to refine pose estimation so that a more accurate density map can be reconstructed. CryoPoseNet [13] utilizes an auto-encoder framework to address the whole reconstruction problem, in which a CNN-based encoder learns to infer the pose of each particle, and a decoder maintains a density map and simulates the cryo-EM imaging process. A reconstruction loss calculated between the dataset images and the simulated reconstructed images so that the framework is trained in an end-to-end manner. CryoGAN [6] uses a paradigm of generative adversarial networks (GANs) [4] to sidestep pose estimation, but the price is that the reconstruction resolution is relatively low.

Some works aim to model the continuous heterogeneous conformation of a biomolecule [7,25]. [7] is a modification of CryoGAN, while [25] uses variational auto-encoder (VAE) [11] to infer a latent conformation representation for each particle.

2.3 Projection-Matching for Cryo-EM

Projection-matching is to find the particle pose via comparing the projections with dataset images. In basic projection-matching-based methods [14,15], the possible poses are discretized and grid search is performed. [5,17] use local search strategy to reduce computational cost, where grid search is only performed in a region near previous estimation, similar to our heuristic search strategy. Even though, grid-search-based projection-matching is still much time-consuming [18].

In [25], branch and bound (BNB) search is proposed. The search space is subdivided hierarchically and the branches with high matching loss are pruned.

Our beam search uses the same subdivision method but different pruning and stopping criterion.

3 Problem Statement

Tomographic reconstruction aims to reconstruct the electron density map $\phi \in \mathbb{R}^{n \times n \times n}$ from a given set of projection images $\{I_1, I_2, \dots, I_N\}$, where n is the box size of the map, equal to the side length of given images, and N is the dataset size. Since each $I_i \in \mathbb{R}^{n \times n}$ is derived from the cryo-EM imaging process, it can be approximately formulated as

$$I_i = C_{d_i} \otimes P_z \circ (R_i \circ \phi) + \epsilon_i, \quad (1)$$

where $\epsilon_i \in \mathbb{R}^{n \times n}$ is an additive noise, and C_{d_i} is convolutional coefficients called point-spread function (PSF), which is determined by the contrast transfer function (CTF) parameters d_i of the cryo-EM. R_i here is a rotation operator, rotating the map ϕ according to the underlying rotation matrix $r_i \in SO(3)$ of the i -th particle's orientation. P_z is a projection operator along z -axis, defined as

$$(P_z \circ \phi)_{x,y} = \int_{-\infty}^{+\infty} \phi_{x,y,z} dz. \quad (2)$$

Once the density map ϕ is represented as a 3-dimensional array, the integral in Eq. (2) can be approximated by summation.

Notice that the d_i is known parameter for all i , while the ϵ_i and R_i (and the underlying r_i) are unknown here. Besides, we omit the translational shift in particle pose. It can be added into Eq. (1) as a translation operator with unknown parameters if necessary.

4 Method

4.1 Overview

We consider $\Omega \subset \mathbb{R}^{n \times n}$ as the space of normalized 2D noisy projection images, $\tilde{\Omega} \subset \mathbb{R}^{n \times n}$ as the space of normalized noiseless ones, and $SO(3)$ as the special orthogonal group in 3-dimensional space. An auto-encoder architecture is adopted to simulate the imaging process of cryo-EM. Specifically, the encoder, $E_\theta : \Omega \rightarrow SO(3)$, is a convolutional network to predict the pose of an image, and decoder, $D_\phi : SO(3) \rightarrow \tilde{\Omega}$, is a differentiable cryo-EM physical simulator. For example, if the orientation operator of the i -th image is R_i together with CTF parameters d_i , the decoding process might be formulated as

$$D_\phi = N(C_{d_i} \otimes P_z \circ (R_i \circ \phi)), \quad (3)$$

where $N(\cdot)$ is normalization operator for 2D spacial pixels, transferring the mean and standard deviation of projection images to 0 and 1, respectively. Notice

that such normalization is effective to blind the numerical differences among various datasets, which further guarantees the generalizability of model hyper-parameters.

Given a set of noisy projection images $\{I_1, \dots, I_N\} \subset \Omega$, the ultimate goal is to reconstruct the 3D density map ϕ . The simplest way is to train a model in an end-to-end manner by minimizing a reconstruction loss shown in Eq. (4)

$$L(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N \ell_{rec}(I_i, (D_\phi \circ E_\theta)(I_i)). \quad (4)$$

However, the trick part here is that orientation is not an intrinsic property of an image. Instead, it depends on the orientation of the 3D density map. This makes encoder hard to learn before it implicitly learns the orientation information of the current density map by gradient information. So, direct end-to-end training would be even impossible to converge at all for complex biomolecules, which is shown in Section 5.2.

To alleviate this issue, we propose *denoised heuristic beam search* (DHBS) as a shortcut for encoder to locate density map orientation by directly learning searched best poses. When encoder adapts to the density map, it is then trained in an end-to-end manner.

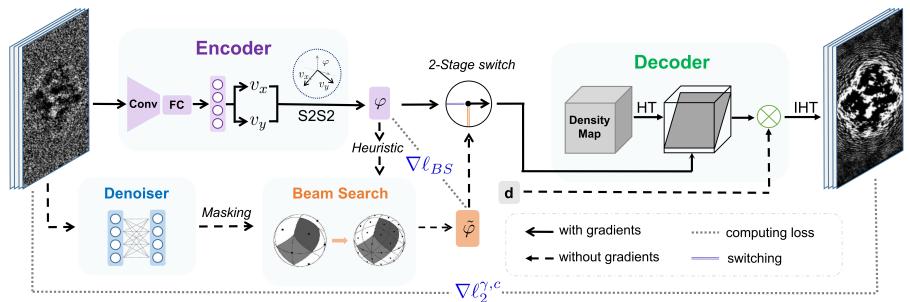


Fig. 1. Overall Cyro2SSNet architecture. Cyro2SSNet is composed of an auto-encoder framework and a search module. In the auto-encoder part, encoder is a CNN performing pose inference for each image, and decoder is a cryo-EM simulator, maintaining a current reconstructed density map and performing imaging process along given poses. The search module looks for a projection orientation matching a given image via DHBS. The training pipeline consists of 2 stages. In the first stage, the searched pose from DHBS supervises encoder and is fed to decoder. In the second stage, DHBS is turned off and the auto-encoder is trained in an end-to-end manner.

4.2 Auto-encoder Architecture

Our overall auto-encoder model consists of two parts, encoder E_θ and decoder D_ϕ . In practice, encoder E_θ takes in a batch of *normalized* noisy projection images and outputs predicted poses (a batch of rotation matrices), while decoder simulates the projection process given poses and outputs *normalized* noiseless projection images. One can refer to Fig. 1, and details are explained below.

Encoder is a convolutional neural network (CNN) followed by a pose prediction head. To output a rotation matrix, we use S2S2 as our parameterization method, which first maps the output feature vector to two 3-dimensional vectors and computes the uniquely determined rotation matrix.

Decoder is a differentiable cryo-EM simulator, modeling the density map as a voxel grid. The simulation is achieved efficiently with the help of the *Fourier slice theorem*, which claims that the Fourier transform of 2D projections is equivalent to the slices of Fourier transformed 3D map. However, complex number induced by Fourier transformation makes gradient backpropagation inconvenient. So we use Hartley transform instead, which holds the slice theorem while eliminates the complex number issue. As shown in Fig. 1, given a batch of poses, 3D Hartley transform is performed on the density map, and slices are extracted according to the poses. The final projection images are obtained by performing inverse Hartley transform on CTF-corrupted slices.

4.3 Denoised Heuristic Beam Search

Due to the highly non-convex property and no-reference problem of encoder, we present *denoised heuristic beam search* (DHBS) to alleviate the optimization difficulty. Instead of directly passing predicted poses from encoder for subsequent decoding, we stabilize and refine the predicted poses by a heuristic search. Specifically, if $r_i \in SO(3)$ is the predicted pose by encoder for the i -th image, DHBS finds \tilde{r}_i defined in Eq. (5), where $De(\cdot)$ denotes a pretrained denoiser and $D_\phi(\cdot)$ denotes decoder. Rather than a global exhaustive search, DHBS is a local heuristic beam search that dramatically reduces the search cost. There are two uses of the output \tilde{r}_i : 1) serves as a weak-supervision to r_i and 2) guides reconstruction. The training pipeline with DHBS is described in Algorithm 1 and each component is explained in detail then.

$$\tilde{r}_i \approx \operatorname{argmin}_{r \in Neighbour(r_i)} \ell_2[De(I_i), De(D_\phi(r))]. \quad (5)$$

Iterative discretization in $SO(3)$. To find \tilde{r}_i , the central problem is to model and discretize $Neighbour(r_i)$ in a reasonable manner. As introduced by [24], $SO(3)$ can be discretized iteratively by Hopf fibration. Let $G_k \subset SO(3)$ be a grid in resolution k , whose elements represent the centers of uniformly partitioned regions in $SO(3)$. Each region in resolution k can be further divided into 8 subregions uniformly. All the centers of these subregions compose G_{k+1} , the grid of next resolution level. Concretely, define the operator *SubDiv* as

$$SubDiv(e; k) = \{8 \text{ elements divide region of } e\} \subset G_{k+1}, \quad \forall e \in G_k, \quad (6)$$

then

$$G_{k+1} = \bigcup_{e \in G_k} SubDiv(e; k). \quad (7)$$

The base grid G_0 has 72 elements due to geometric property of $SO(3)$, so $|G_k| = 72 \times 8^k$. The naive brute force search is to find the pose globally by modeling $Neighbour(r_i) \approx G_k$ for a certain k , i.e.,

$$\tilde{r}_i^{BF} = \operatorname{argmin}_{r \in G_k} \ell_2[De(I_i), De(D_\phi(r))]. \quad (8)$$

However, the computational cost is exponential with respect to resolution k , hence a wiser search strategy is necessary for efficiency.

Heuristic Beam Search. As mentioned above, the search process should be done efficiently by trading off accuracy and efficiency, so we present *heuristic beam search* to evaluate poses iteratively and achieve *linear* time cost w.r.t. resolution level k . First we select c nearest neighbors of r_i in the base grid G_0 . Then these c poses are subdivided into $8c$ poses in the next level by $SubDiv(\cdot)$ operator. Compute the ℓ_2 costs for these poses and keep track of the top- c with minimal costs. This division-evaluation process is repeated for i times and the best pose among the last c poses is returned, termed as \tilde{r}_i^{BS} . The search width c and depth i can be tuned for an accuracy-efficiency tradeoff. The beam search degrades to greedy search when $c = 1$ and remains more currently sub-optimal poses when c is larger. Notice the heuristic trick that picking nearest neighbours of r_i instead of comparing all candidates in G_0 reduces nearly half of the time cost in practice.

Gradient Descent Enhancement. Though the searched pose usually results in a smaller cost compared with predicted pose by encoder, there is still a theoretical lower bound of distance between searched pose and the globally best one due to discrete nature of grid for a finite search depth. To tackle this problem, we first make the following assumption.

Assumption 1 (Local Convexity) *The cost landscape is locally convex near the globally best pose and the previously searched pose \tilde{r}_i^{BS} is close to globally best pose. So that it is within this convex region.*

There are some empirical evidence for the above assumption in Appendix B.3. Based on this, a further refinement after beam search is performed to update the pose locally and continuously by gradient descent. Specifically, we update the searched pose \tilde{r}_i^{BS} for n steps with learning rate η as formulated in Eq. (9). This gradient descent can be performed since the cost is differentiable with respect to \tilde{r}_i^{BS} .

$$\tilde{r}_i^{BS} \leftarrow \tilde{r}_i^{BS} - \eta \cdot \nabla_{\tilde{r}_i^{BS}} [\ell_2[De(I_i), De(D_\phi(\tilde{r}_i^{BS}))]]. \quad (9)$$

Denoiser. Due to the extremely low SNR of images, direct comparison between noisy and noiseless images may cause some issues (e.g. a sub-optimal search result). Hence, we denoise the images while computing the search cost. The denoiser we use is a simple neural network, described in Appendix A.2.

It is noticeable that perfect denoising quality is not required. The denoised images are only utilized to evaluate searched candidates, rather than to be reconstructed. Denoiser helps DHBS to focus on low-frequency features and find right poses in early stage and is turned off after a few epochs. This fits to a coarse-to-fine criterion, proven to be useful in many computer vision tasks.

Composing all the techniques above, an overall algorithm of DHBS is shown in Algorithm 1.

Algorithm 1 Denoised Heuristic Beam Search (DHBS) algorithm

Input: A noisy image I and predicted pose r , denoiser $De(\cdot)$, decoder $D_\phi(\cdot)$ and hyper-parameters (c, i, n, η)

Output: Refined pose \tilde{r}

- 1: $S_0 = \{e \mid e \in G_0 \text{ and } e \text{ is top } c \text{ nearest neighbour of } r\}$
- 2: **for** $k = 1, 2, \dots, i$ **do**
- 3: $\tilde{S}_k = \bigcup_{e \in S_{k-1}} SubDiv(e; k - 1)$ (sub-division)
- 4: $S_k = \text{low-}c_{e \in \tilde{S}_k} \ell_2[De(I), De(D_\phi(e))]$ (best c candidates)
- 5: **end for**
- 6: $\tilde{r} = \operatorname{argmin}_{e \in S_k} \ell_2[De(I), De(D_\phi(e))]$
- 7: **for** $j = 1, 2, \dots, n$ **do**
- 8: $\tilde{r} \leftarrow \tilde{r} - \eta \cdot \nabla_{\tilde{r}} \ell_2[De(I), De(D_\phi(\tilde{r}))]$
- 9: **end for**
- 10: Return: \tilde{r}

4.4 2-Stage Switch and Loss Objectives

First, a metric measuring distances in $SO(3)$ space is needed in our loss objectives. Proved in [8], there are a few equivalent metrics, among which *quaternion distance* is the one we use. For $R_1, R_2 \in SO(3)$, let $\varphi_1, \varphi_2 \in \mathbb{R}^4$ be quaternions of R_1, R_2 , respectively. The quaternion distance between R_1 and R_2 is

$$d_{\text{quat}}(R_1, R_2) = \min(\|\varphi_1 - \varphi_2\|_2, \|\varphi_1 + \varphi_2\|_2). \quad (10)$$

As mentioned before, DHBS helps encoder to implicitly learn the orientation and speeds up reconstruction by polishing predicted poses from encoder. When encoder learns to infer poses close to globally best ones, the optimization difficulty largely mitigates. Based on this fact, we present a 2-stage training pipeline to combine search-supervised stage and end-to-end training stage.

During the first stage, the gradient is not passed from decoder back to encoder. Instead, the searched poses from DHBS are used as weak-supervision for training encoder and passed to decoder. Specifically, if the predicted pose is r and searched pose is \tilde{r} , an auxiliary loss for weak-supervision is

$$\ell_{BS}(r, \tilde{r}) = d_{\text{quat}}^2(r, \tilde{r}). \quad (11)$$

As for the reconstruction loss, we notice that vanilla ℓ_2 loss is *infeasible* to define distance between space Ω and $\tilde{\Omega}$, since they are both normalized. Instead,

we proposed an unbiased estimation of vanilla ℓ_2 loss in Eq. (12), termed as $\ell_2^{\gamma,c} : \Omega \times \tilde{\Omega} \rightarrow \mathbb{R}^+$, where γ, c are obtained according to the SNR of images. Formally, if the SNR of space Ω is known, say s , then $\gamma = \sqrt{1 + 1/s}$ and $c = 1/s$. The x, y in the formula are two-dimensional indices. The derivation is shown in Appendix A.3. Once the SNR of a given dataset is unknown, we find a slightly overestimated SNR can make our model work as well, which makes it more flexible.

$$\ell_2^{\gamma,c}(I_i, \tilde{I}_i) = \frac{1}{n^2} \sum_{1 \leq x, y \leq n} (\gamma \cdot I_{i,x,y} - \tilde{I}_{i,x,y})^2 - c, \quad I_i \in \Omega, \quad \tilde{I}_i \in \tilde{\Omega}. \quad (12)$$

We notice that passing searched pose \tilde{r} to decoder is better than passing predicted pose r from encoder, since \tilde{r} is closer to the globally best pose hence better for reconstruction. By doing so, the gradient is cut between encoder and decoder, and $\ell_{BS}, \ell_2^{\gamma,c}$ optimize encoder and decoder respectively. They are combined by a tunable hyper-parameter, resulting in the loss function in the first stage, as

$$L_{S1}(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N \ell_2^{\gamma,c}(I_i, D_\phi(\tilde{r}_i)) + \lambda \cdot \ell_{BS}(r_i, \tilde{r}_i), \quad (13)$$

where: $r_i = E_\theta(I_i)$, $\tilde{r}_i = DHBS(r_i; D_\phi, I_i)$.

During the second stage, we assume encoder has already implicitly learned a rough orientation of 3D density map. Thus, encoder and decoder can be fine-tuned in an end-to-end manner with loss objective

$$L_{S2}(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N \ell_2^{\gamma,c}(I_i, D_\phi(E_\theta(I_i))) \quad (14)$$

We name our method *2-stage Search Supervised Network* (Cryo2SSNet). The overall algorithm pipeline of Cryo2SSNet is summarized in Algorithm 2.

5 Experimental Results

We now describe our experiments and results on several simulated protein datasets, including 4AKE, 5A1A, 7ACT and 7AD1. We derived the atomic models in PDB format from protein data bank¹ and generated their density maps using the open-source software EMAN2 [20], and built simulated datasets by CryoSPARC [16]. All simulated datasets use the default CTF parameters in CryoSPARC.

In each experiment, we reconstructed the density map with full dataset and calculated the FSC curve between the density map reconstruction and the ground truth. The resolution cutoff in FSC curve is 0.5.

¹ <https://www.rcsb.org/>

Algorithm 2 2-Stage Search Supervised Network (Cryo2SSNet)

Input: Noisy projection image dataset $\{I_1, I_2, \dots, I_N\}$, 2-stage epochs (E_1, E_2) .

Output: Reconstructed density map ϕ

```

405   1: Initialize  $E_\theta$  with random weight and  $D_\phi$  with scaled uniform random distribution
406      $U(0, 0.1)$  in region within an estimated protein radius and zero outside.
407   2: # Stage 1
408   3: for  $epoch = 1, 2, \dots, E_1$  do
409     4:   for  $\mathbf{I} \leftarrow$  sampled batch of images do
410       5:      $\mathbf{r} \leftarrow E_\theta(\mathbf{I})$ 
411       6:      $\tilde{\mathbf{r}} \leftarrow DHBS(\mathbf{r}; D_\phi, \mathbf{I})$ 
412       7:      $L = \ell_2^{c,c}(\mathbf{I}, D_\phi(\tilde{\mathbf{r}})) + \lambda \cdot \ell_{BS}(\mathbf{r}, \tilde{\mathbf{r}})$ 
413       8:     Update  $\theta, \phi$  with  $\nabla_{\theta, \phi} L$ 
414   9:   end for
415 10: end for
416 11: # Stage 2
417 12: for  $epoch = E_1 + 1, \dots, E_1 + E_2$  do
418    13:   for  $I \leftarrow$  sampled batch of images do
419      14:      $\mathbf{r} \leftarrow E_\theta(\mathbf{I})$ 
420      15:      $L = \ell_2^{c,c}(\mathbf{I}, D_\phi(\mathbf{r}))$ 
421      16:     Update  $\theta, \phi$  with  $\nabla_{\theta, \phi} L$ 
422    17:   end for
423 18: end for
424 19: Return:  $\phi$ 
425

```

427 5.1 Reconstruction

428 Experiment Settings For 4AKE, the protein used in CryoPoseNet [13], we
 429 generated a 1.6 Å density map of size $128 \times 128 \times 128$, with voxel size 0.8 Å,
 430 which corresponds to a Nyquist-Shannon limit of 1.6 Å for reconstruction. We
 431 created 1 noiseless dataset and 3 noisy datasets with SNRs equal to 10 dB, 0
 432 dB, -10 dB, respectively. Each dataset contains 10,000 projection images.

For 5A1A, the one used in CryoGAN [6], following their experimental settings, we generated a 2.5 Å density map of size $256 \times 256 \times 256$, with voxel size 1.25 Å, corresponding to a Nyquist-Shannon limit of 2.50 Å. 3 datasets with SNRs equal to -5.2 dB, -10 dB, -20 dB, respectively, were created. Each dataset contains 40,000 images. Though 5A1A is in a structure of D2 symmetry, unlike [6], we do not explicitly utilize this property in our reconstruction method.

One can refer to Appendix B for detailed settings and additional results for the remaining proteins.

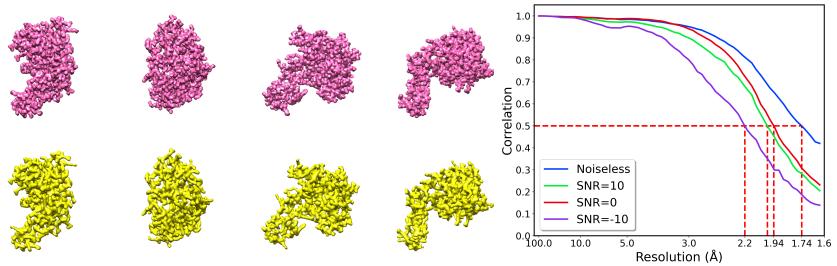
In each experiment, we ran Cryo2SSNet for 7500 steps in total with a batch size of 32, which corresponds to 6 epochs for 5A1A and 24 epochs for others. We divided the total number of epochs equally for two stages. As for DHBS, the search depth and search width were set to be 5 and 4, respectively. More detailed training parameters are listed in Appendix B.2.

448 Results Cryo2SSNet achieved a reconstruction resolution of 2.20 Å for 4AKE
 449 with SNR=-10 dB and 4.64 Å for 5A1A with SNR=-20 dB. Our reconstruc-

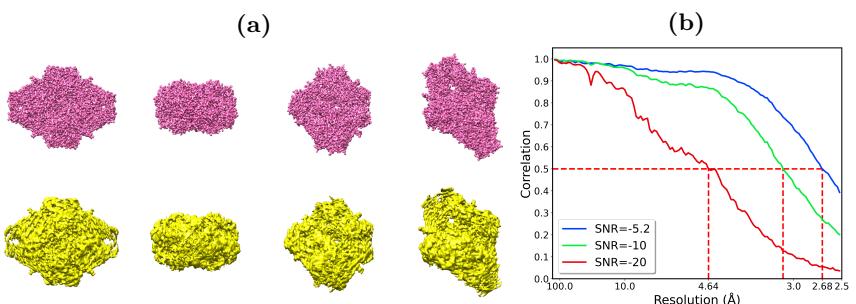
tion results are shown in Table 1. It demonstrates that our proposed method reconstructs the density maps of 4AKE and 5A1A with fine qualities close to the Nyquist-Shannon limit, and it can generalize to different noise levels and protein types. Visualization of our reconstruction results is shown in Fig. 2.

Table 1. The reconstruction resolutions of Cryo2SSNet. The results are compared with ones reported by previous proposed methods.

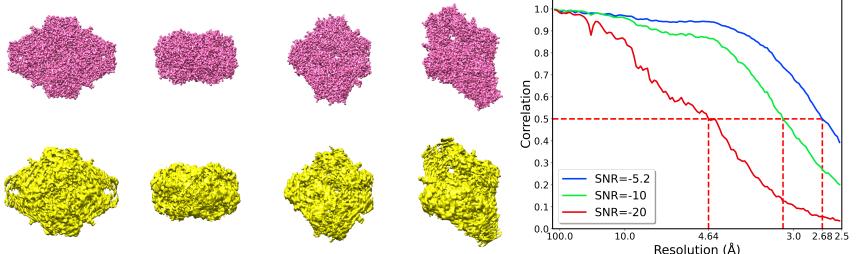
Protein	SNR (dB)	Resolution (\AA)	Previous Resolution (\AA)
4AKE	∞	1.74	/
	10	1.99	/
	0	1.94	2.57 (CryoPoseNet) ²
	-10	2.20	/
5A1A	-5.2	2.68	8.30 (CryoGAN)
	-10	3.14	/
	-20	4.64	15.30 (CryoGAN)



(a)



(b)



(c)

(d)

Fig. 2. Reconstruction results. (a) and (c) visualize our reconstructed density maps for 4AKE with SNR = -10 dB and 5A1A with SNR = -20 dB, respectively. The pink ones are ground truths while the yellow ones are reconstructions. (b) and (d) shows the FSC curves in various noise levels for 4AKE and 5A1A, respectively.

² The original paper of CryoPoseNet didn't mention the resolution of density map used to simulate dataset. So we compare it with the one of same SNR and use as reference only.

We notice that the reconstructed density map converges quickly during the first several epochs as illustrated in Fig. 4. It's indicated that the density map converges to a rough sketch first and is refined in details gradually. This also reveals that constructing a rough initial density map is much easier compared to structure refinement. The challenging part of tomographic reconstruction is, in order to derive a high-resolution reconstruction, an *accurate enough* pose for each particle is required to be inferred according to the extremely noisy images. After our first stage, the density map is almost converged. The second stage coordinates encoder and decoder by training them jointly at a fine pace. By utilizing DHBS to enhance supervision for encoder, we observe a high convergence speed as shown in Fig. 3, which is necessary for the second training stage.

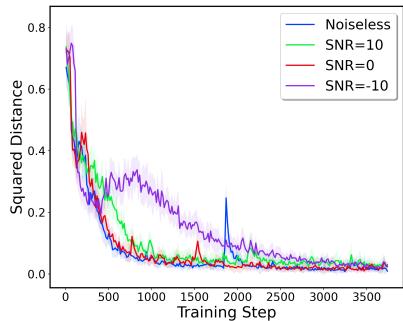


Fig. 3. Curves of squared distance between predicted and searched poses, tested on 4AKE.

The second stage coordinates encoder and decoder by training them jointly at a fine pace. By utilizing DHBS to enhance supervision for encoder, we observe a high convergence speed as shown in Fig. 3, which is necessary for the second training stage.

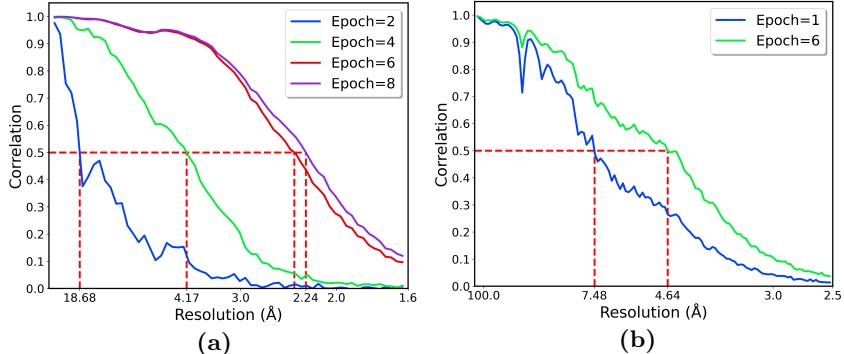


Fig. 4. Map resolution evolution. (a) FSC curves of 4AKE SNR=-10 dB of first several epochs. (b) Comparison of starting epochs and last epochs of 5A1A SNR=-20 dB.

5.2 Necessity of DHBS

It can be verified that DHBS is necessary for extremely noisy dataset, e.g., 5A1A with SNR=-20 dB. We run reconstruction by eliminating stage 1 but run twice of the total epochs (12 epochs) for time compensation. We notice, in this setting, the density map cannot converge at all. The comparison of reconstruction results is shown in Fig. 5. By eliminating DHBS, the result of end-to-end training is even worse than the first epoch of Cryo2SSNet,

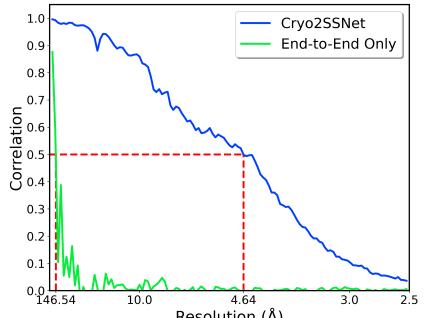


Fig. 5. Necessity of DHBS. End-to-end training does not work at all in 5A1A SNR=-20 dB.

referring to Fig. 4(b), which indicates that it fails to find a rough sketch of density map.

There are several possible reasons for this. 1) Strong noise in images confuses encoder to implicitly coordinate with decoder, 2) high variance of encoder predictions causes density map fluctuating, and 3) the whole model gets stuck in a bad local minimum. With the help of DHBS, these issues may be mitigated to a degree.

5.3 Effectiveness of Gradient Descent Enhancement

In this part we conduct experiments to verify the effectiveness of gradient descent (GD) enhancement. 5A1A with SNR=−10 dB is selected as testbed. We run our method with and without GD enhancement, keeping other settings the same as described in Appendix B. Also, to further reveal the potential of GD enhancement, we perform reconstruction using Adam [10] instead of vanilla GD in Cryo2SSNet, to update searched pose \tilde{r}^{BS} in each step. The results are shown in Fig. 6, where the curve “w/ GD” is the result of Cryo2SSNet. It is shown that gradient descent enhancement is necessary for Cryo2SSNet to achieve high resolution without enlarging search depth. As indicated by results, GD enhancement with Adam outperforms the vanilla one. One can substitute GD in Cryo2SSNet with Adam for better results, while GD is used throughout this paper for simplicity.

5.4 Effectiveness of Denoiser

As we mention in Section 4.3, the denoiser allows a coarse-to-fine process when the images are extremely noisy. The denoising results are visualized in Fig. 7(a). As we can tell, denoiser doesn’t remove all the noises, but does emphasize and expose a rough structure of the particle in a certain view, which helps a lot in looking for correct poses in the early stage. In practice, instead of denoising the entire image, we first take a circular mask of a rough particle radius to filter out the noises outside. This radius can be easily estimated by taking numerical average over the entire dataset. Applying this mask helps denoiser to focus on the region where the particle locates.

Additionally, we conduct quantitative experiments using protein 7AD1 with SNR=−10 dB, whose structure is the most complicated among the 4 proteins. Details of this dataset are described in Appendix B.1. To highlight the effect of denoiser, we use search width of 3 and search depth of 3, instead of 4 and 5 in other experiments, and keep other parameters unchanged. In this experiment, we control the number of epochs applying denoiser. The results are shown in

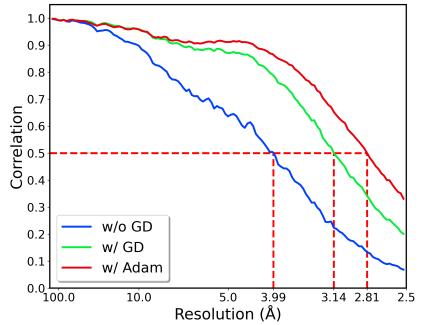


Fig. 6. Effectiveness test of GD. Removing GD enhancement destroys the performance.

Fig. 7(b). Since there are 12 epochs in the first stage, “Stop Epoch=0” (in short as $D0$ later) means not to apply denoiser and “Stop Epoch=12” means to apply denoiser through the entire first stage. As shown, $D0$ fails completely and $D12$ degrades the performance compared to $D4$ and $D8$. This result coincides with our expectation: denoiser helps to search right poses, but, as a compromise, loses high-frequency information. Therefore, applying denoiser in early stage and turning it off later mimics a coarse-to-fine paradigm in computer vision.

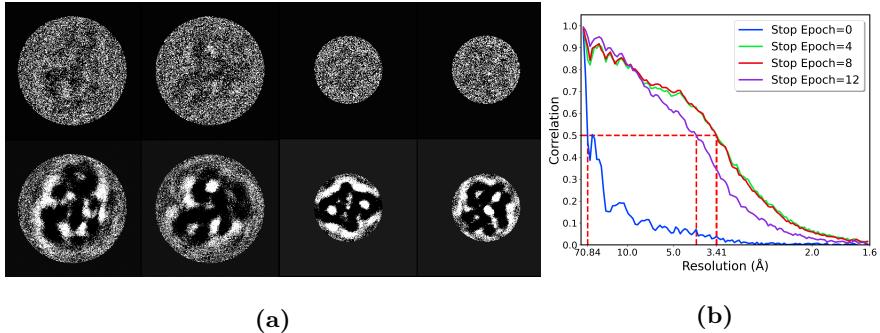


Fig. 7. Qualitative and quantitative analysis of denoising effectiveness. (a) Denoising results. The first row is the raw images of 7AD1 (first two; SNR = -10 dB) and 5A1A (last two; SNR = -20 dB), while the second row is the denoising results. Denoiser focuses on a particle-size centered region. (b) FSC curves of 7AD1 reconstruction with SNR = -10 dB. Applying denoiser in early epochs and turning it off outperforms other settings.

6 Conclusion

In this paper, we present Cryo2SSNet, a search supervised auto-encoder tackling cryo-EM single-particle reconstruction problem. Specifically, a DHBS method is proposed to help pose inference. The search module is used to supervise the encoder in early stage and turned off later, resulting in a 2-stage training pipeline.

Cryo2SSNet is shown to be able to achieve higher reconstruction resolution than previous deep-learning-based methods on a variety of simulated protein datasets, within comparable time cost. Our experiments indicate that DHBS helps encoder learn to infer poses more accurately, which is necessary to reconstruct high-resolution density maps. Also, we claim that Cryo2SSNet can generalize to different datasets, and can be further tuned for a tradeoff between accuracy and efficiency.

However, there are still limitations of our work. The second stage of Cryo2SSNet can make only minor refinement of the density map, which doesn't meet our expectations. Besides, omission of translational shift in particle pose makes it hard to adapt to real datasets. We leave solving them as future works.

630 References

- 632 1. Baker, T.S., Cheng, R.H.: A model-based approach for determining orientations of
633 biological macromolecules imaged by cryo-electron microscopy. *Journal of structural
634 biology* **116**(1), 120–130 (1996)
- 635 2. Bepler, T., Kelley, K., Noble, A.J., Berger, B.: Topaz-denoise: general deep denoising
636 models for cryoem and cryoet. *Nature communications* **11**(1), 1–12 (2020)
- 637 3. Bepler, T., Morin, A., Rapp, M., Brasch, J., Shapiro, L., Noble, A.J., Berger,
638 B.: Positive-unlabeled convolutional neural networks for particle picking in cryo-
639 electron micrographs. *Nature methods* **16**(11), 1153–1160 (2019)
- 640 4. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S.,
641 Courville, A., Bengio, Y.: Generative adversarial networks. *Communications of the
642 ACM* **63**(11), 139–144 (2020)
- 643 5. Grigorieff, N.: Frealign: high-resolution refinement of single particle structures.
644 *Journal of structural biology* **157**(1), 117–125 (2007)
- 645 6. Gupta, H., McCann, M.T., Donati, L., Unser, M.: Cryogan: A new re-
646 construction paradigm for single-particle cryo-em via deep adversarial learning.
647 *IEEE Transactions on Computational Imaging* **7**, 759–774 (2021).
<https://doi.org/10.1109/TCI.2021.3096491>
- 648 7. Gupta, H., Phan, T.H., Yoo, J., Unser, M.: Multi-cryogan: Reconstruction of con-
649 tinuous conformations in cryo-em using generative adversarial networks. In: European
650 Conference on Computer Vision. pp. 429–444. Springer (2020)
- 651 8. Huynh, D.Q.: Metrics for 3d rotations: Comparison and analysis. *Journal of Math-
652 ematical Imaging and Vision* **35**(2), 155–164 (2009)
- 653 9. Jiménez-Moreno, A., Strelák, D., Filipović, J., Carazo, J., Sorzano, C.: Deepalign,
654 a 3d alignment method based on regionalized deep learning for cryo-em. *Journal
655 of Structural Biology* **213**(2), 107712 (2021)
- 656 10. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint
657 arXiv:1412.6980 (2014)
- 658 11. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint
659 arXiv:1312.6114 (2013)
- 660 12. Kühlbrandt, W.: The resolution revolution. *Science* **343**(6178), 1443–1444 (2014)
- 661 13. Nashed, Y.S.G., Poitevin, F., Gupta, H., Woppard, G., Kagan, M., Yoon, C.H.,
662 Ratner, D.: Cryoposenet: End-to-end simultaneous learning of single-particle orien-
663 tation and 3d map reconstruction from cryo-electron microscopy data. In: Pro-
664 ceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)
665 Workshops. pp. 4066–4076 (October 2021)
- 666 14. Penczek, P., Radermacher, M., Frank, J.: Three-dimensional reconstruction of sin-
667 gle particles embedded in ice. *Ultramicroscopy* **40**(1), 33–53 (1992)
- 668 15. Penczek, P.A., Grassucci, R.A., Frank, J.: The ribosome at improved resolution:
669 new techniques for merging and orientation refinement in 3d cryo-electron mi-
670 croscopy of biological particles. *Ultramicroscopy* **53**(3), 251–270 (1994)
- 671 16. Punjani, A., Rubinstein, J.L., Fleet, D.J., Brubaker, M.A.: cryosparc: algorithms
672 for rapid unsupervised cryo-em structure determination. *Nature methods* **14**(3),
673 290–296 (2017)
- 674 17. Scheres, S.H.: Relion: implementation of a bayesian approach to cryo-em structure
675 determination. *Journal of structural biology* **180**(3), 519–530 (2012)
- 676 18. Sigworth, F.J.: Principles of cryo-em single-particle image processing. *Microscopy*
677 **65**(1), 57–67 (2016)

- 675 19. Singer, A., Sigworth, F.J.: Computational methods for single-particle electron cryo-
676 omicroscopy. *Annual Review of Biomedical Data Science* **3**, 163–190 (2020)
677 20. Tang, G., Peng, L., Baldwin, P.R., Mann, D.S., Jiang, W., Rees, I., Ludtke, S.J.:
678 Eman2: an extensible image processing suite for electron microscopy. *Journal of
679 structural biology* **157**(1), 38–46 (2007)
680 21. Tegunov, D., Cramer, P.: Real-time cryo-em data pre-processing with warp.
681 BioRxiv p. 338558 (2018)
682 22. Wagner, T., Merino, F., Stabrin, M., Moriya, T., Antoni, C., Apelbaum, A., Hagel,
683 P., Sitsel, O., Raisch, T., Prumboam, D., et al.: Spire-cryolo is a fast and accurate
684 fully automated particle picker for cryo-em. *Communications biology* **2**(1), 1–13
685 (2019)
686 23. Wang, F., Gong, H., Liu, G., Li, M., Yan, C., Xia, T., Li, X., Zeng, J.: Deeppicker:
687 A deep learning approach for fully automated particle picking in cryo-em. *Journal
688 of Structural Biology* **195**(3), 325–336 (2016)
689 24. Yershova, A., Jain, S., Lavalle, S.M., Mitchell, J.C.: Generating uniform incremen-
690 tal grids on so (3) using the hopf fibration. *The International journal of robotics
691 research* **29**(7), 801–812 (2010)
692 25. Zhong, E.D., Bepler, T., Davis, J.H., Berger, B.: Reconstructing continuous distri-
693 butions of 3d protein structure from cryo-em images. In: *International Conference
694 on Learning Representations* (2019)
695 26. Zhu, Y., Ouyang, Q., Mao, Y.: A deep convolutional neural network approach to
single-particle recognition in cryo-electron microscopy. *BMC bioinformatics* **18**(1),
1–10 (2017)

720 A Appendix-Method

721 A.1 Encoder Architecture

722 Here we describe the architecture used in encoder. The CNN is composed of
 723 3 convolutional blocks followed by 2 fully connected layers (MLP) and a pose
 724 prediction head. Each convolutional block consists of 2 Conv-ReLU modules
 725 and 1 max-pooling layer. The numbers of channels for each block are 32, 64, 128,
 726 respectively. Each block down-samples feature map by a factor of 2. The hidden
 727 dimensions for 2 subsequent fully connected layers are both 512.

728 The pose is then predicted by S2S2 parameterization. Specifically, two 3D
 729 dimensional vectors v_x and v_y are predicted from obtained feature vector through
 730 two fully-connected layers. In S2S2, v_y is projected to orthogonal direction of
 731 v_x . Then v_x and projected v_y are normalized, termed as \hat{v}_x and \hat{v}_y . Then \hat{v}_z
 732 is defined as the cross product of \hat{v}_x and \hat{v}_y . Finally $[\hat{v}_x \hat{v}_y \hat{v}_z]$ is the predicted
 733 rotation matrix.

736 A.2 Denoiser

737 The denoiser used in this work is a simple MLP, pretrained by the normalized
 738 noisy projection images unsupervisedly. It is a 5-layer auto-encoder model. The
 739 first layer down-samples the image's dimension (i.e. the total number of pixels)
 740 by a factor of f , a tunable hyper-parameter. The 3 consequent hidden layers keep
 741 the dimension and the last one recovers the original image size. The denoiser is
 742 optimized to minimize the ℓ_2 reconstruction loss.

743 In the MLP, down-sampling serves as an information bottleneck filtering
 744 amount of noise out. CNN is not used here since MLP is capable enough to
 745 achieve reasonably good denoising result compared with CNN. Hence we only
 746 use an MLP denoiser for simplicity and efficiency.

747 The denoiser is trained for 7500 steps with batch size of 64, using Adam
 748 optimizer with initial learning rate of $1e-4$. We use unbiased ℓ_2 loss proposed
 749 in Eq. (12) as the optimization objective. Training time for denoiser is negligible
 750 compared to reconstruction time (less than 5%). The downsample factors f for
 751 each protein are shown in Table 3.

755 A.3 Unbiased Estimation of ℓ_2 Loss

756 Here we give the derivation of unbiased estimation of ℓ_2 loss. We represent
 757 noiseless projection image (corrupted by CTF already) as a random variable
 758 defined in $\mathbb{R}^{n \times n}$, where n is the edge length of image. Besides, we model the
 759 noise as a zero-mean random variable. We assume SNR of the images is known.
 760 Specifically, let $d, e \in \mathbb{R}^{n \times n}$ be random variables for projection image and noise,
 761 following a joint distribution $p(d, e)$. A noisy image in dataset is $x = d + e$. Let
 762 $f(\cdot)$ denote our auto-encoder model. $\mathbb{E}(\cdot)$ stands for the expectation of a certain
 763 random variable.

Definition 1 (ℓ_2 Loss). For a noiseless image d and noise e , the expected squared error between $f(x)$ and d is ℓ_2 loss shown in Eq. (15).

$$\ell_2(d, f(x)) = \mathbb{E}_{d,e \sim p(d,e)} \left[\frac{1}{n^2} \sum_{1 \leq i,j \leq n} (d - f(x))_{i,j}^2 \right] \quad (15)$$

The expectation is taken over the joint probability distribution $p(d, e)$. We do not explicitly write this later. One problem here is that we would like to perform reconstruction over a normalized output for generalization. Here we give the exact definition for spatial normalization.

Definition 2 (Spacial Normalization). Given an instance $y \in \mathbb{R}^{n \times n}$, we define $M(y), S(y)$ as the mean and standard deviation over n^2 positions, shown in Eq. (16). And normalization is simply $N(y) = (y - M(y))/S(y)$.

$$M(y) = \frac{1}{n^2} \left(\sum_{1 \leq i,j \leq n} y_{i,j} \right), \quad S(y) = \sqrt{\frac{1}{n^2} \left(\sum_{1 \leq i,j \leq n} (y - M(y))_{i,j}^2 \right)} \quad (16)$$

According to the definition above, we know $N(x) \in \Omega$ and $N(f(x)) \in \tilde{\Omega}$. So what we want is to find out an estimation of $\ell_2(N(d), N(f(x)))$, since this is inaccessible because we do not obtain d from the dataset. So we have to estimate the above term with the help of the entire dataset. Unfortunately, this is unachievable for general joint probability distribution $p(d, e)$. The relation between d and e should be considered also.

Specifically, considering noise e and image d , we assume e is *independent* of d . This assumption is natural. Even though this may not be correct in real dataset, it still serves as a good approximation. On this assumption, we claim the following property.

Assumption 2 (Independence) The image and noise are independent, i.e., $p(d, e) = p(d) \cdot p(e)$. With this property, for any given point-wise function $\psi(\cdot)$, we have

$$\mathbb{E}_{d,e}[\psi(d) \cdot e] = \mathbb{E}_{d,e}[\psi(d)] \cdot \mathbb{E}_{d,e}[e] = 0. \quad (17)$$

Here point-wise function means applying the same function to each position repeatedly.

Then we give the definition of our proposed unbiased ℓ_2 loss and the main theorem.

Definition 3 (Unbiased L2 Loss). Given $SNR=s$, a noisy image x and a model $f(\cdot)$, let $\gamma = \sqrt{1 + 1/s}$ and $c = 1/s$, then we define $\ell_2^{\gamma,c}$ as

$$\ell_2^{\gamma,c}(N(x), N(f(x))) = \mathbb{E}_{d,e} \left[\frac{1}{n^2} \sum_{1 \leq i,j \leq n} (\gamma \cdot N(x) - N(f(x)))_{i,j}^2 - c \right]. \quad (18)$$

810 **Theorem 1 (Equivalence).** *Given the definition above, we have*

$$813 \quad \ell_2^{\gamma,c}(N(x), N(f(x))) = \ell_2(N(d), N(f(x))). \quad (19)$$

815 *Proof (Equivalence).* First according to the assumptions, we know $\mathbb{E}[e_{i,j}] = 0$
 816 for $\forall 1 \leq i, j \leq n$, and

$$818 \quad \text{SNR} = s = \mathbb{E}_{d,e} \left[\frac{S^2(d)}{S^2(e)} \right]. \quad (20)$$

820 Then we have

$$\begin{aligned} 822 \quad & \ell_2^{\gamma,c}(N(x), N(f(x))) \\ 823 \quad &= \mathbb{E}_{d,e} \left[\frac{1}{n^2} \sum_{1 \leq i, j \leq n} (\gamma \cdot N(x) - N(f(x)))_{i,j}^2 - c \right] \\ 824 \quad &= \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \mathbb{E}_{d,e} \left[\left(N(f(x)) - \gamma \frac{x - M(x)}{S(x)} \right)_{i,j}^2 \right] - c \\ 825 \quad &= \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \mathbb{E}_{d,e} \left[\left(N(f(x)) - \gamma \sqrt{\frac{s}{1+s}} \cdot \frac{d + e - M(d)}{S(d)} \right)_{i,j}^2 \right] - c \\ 826 \quad &= \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \mathbb{E}_{d,e} \left[\left(N(f(x)) - N(d) - \frac{e}{S(d)} \right)_{i,j}^2 \right] - c \\ 827 \quad &= \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \mathbb{E}_{d,e} [(N(f(x)) - N(d))_{i,j}^2] + \mathbb{E}_{d,e} \left[\frac{S^2(e)}{S^2(d)} \right] - c + \\ 828 \quad &\quad \text{red term} \\ 829 \quad &= \frac{2}{n^2} \sum_{1 \leq i, j \leq n} \mathbb{E}_{d,e} \left[\left(\frac{N(f(x)) - N(d)}{S(d)} \cdot e \right)_{i,j} \right] \\ 830 \quad &= \ell_2(N(d), N(f(x))) + 0 + \\ 831 \quad &\quad \text{blue term} \\ 832 \quad &= \ell_2(N(d), N(f(x))). \end{aligned} \quad (21)$$

847 The red term is the one we want and the blue term equals zero due to the defini-
 848 tion of SNR. The cyan term equals zero due to our independence assumption.
 849 Then the proof is done.

851 One thing to notice is that Eq. (12) is an empirical estimation of unbiased
 852 ℓ_2 loss over the dataset. The meaning is clear with the help of the context. The
 853 significance of this theorem is that we are able to optimize an unbiased loss
 854 objective even without access to noiseless images.

B More Experimental Results

B.1 Remaining Protein Settings

We have mentioned the setting of simulated dataset of 4AKE and 5A1A in Section 5.1. Here we describe the setting of other protein simulated datasets.

As for 7ACT, we generated a 1.0 Å density map of size $128 \times 128 \times 128$ with voxel size 0.5 Å, which corresponds to a Nyquist-Shannon limit of 1.0 Å. Only noiseless and noisy dataset of SNR=−10 dB were considered, where each one contains 10,000 projection images.

As for 7AD1, density map resolution was 1.6 Å and size was $256 \times 256 \times 256$, with voxel size 0.8 Å. We still only considered noiseless and noisy dataset of SNR=−10 dB with 10,000 images for each. A summary of 4 protein datasets is in Table 2.

Table 2. This gives the summary of simulated protein dataset. The Nyquist-Shannon limit is the twice of voxel size.

Protein	Map Resolution (Å)	Voxel Size (Å)	Set Size
4AKE	1.6	0.8	10,000
5A1A	2.5	1.25	40,000
7AD1	1.6	0.8	10,000
7ACT	1.0	0.5	10,000

As for CTF parameters, we use the default values in CryoSPARC for all of our simulated datasets. Specifically, the accelerating voltage is 300 kV, the amplitude contrast is 0.07, and the spherical aberration is 2.7 mm. The mean of two defocus is uniformly drawn from 15000 Å to 20000 Å, and the difference between two defocus is uniformly drawn from 100 Å to 500 Å. The defocus angle is uniformly drawn from 0.2 rad to 1.4 rad.

B.2 Training Details

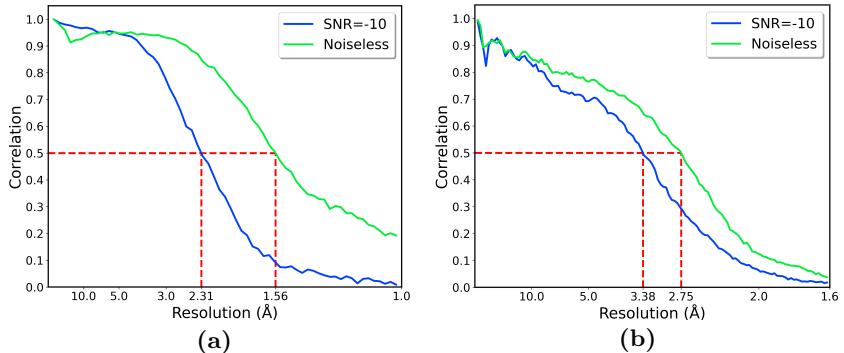
In this section, we describe the training details of our experiments. We use Pytorch for automatic gradient computation. Encoder was initialized by default in Pytorch and the density map in decoder was initialized by a uniform distribution between 0 and 0.1 within a circular region. The outside region was initialized to 0. Protein-sized mask radiiuses are listed in Table 3, which stands for the proportion to half-edge length.

The total number of training steps was fixed as 7500, and the batch size is 32. The number of epochs was 6 for 5A1A and 24 for others. The number of epochs in stage 1 and stage 2 were partitioned equally. As a coarse-to-fine process, we replaced denoiser with an identity function after the first 1/3 of stage 1. The search depth and search width were set to be 5 and 4, respectively. Adam was used as the optimizer with initial learning rate of $1e-4$ for encoder and $1e-2$ for decoder. We only decayed the learning rate of decoder by a factor 10 when

900 entering stage 2. The balance factor between the two loss terms in stage 1 was
 901 set to 1.

902 As for gradient descent enhancement, we used vanilla GD algorithm for 20
 903 steps with learning rate 0.05.

904 All the experiments were running on a single NVIDIA Tesla K80 GPU, and
 905 the training time is shown in Table 3.



906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917
Fig. 8. The reconstruction results of 7ACT (a) and 7AD1 (b) on noiseless and noisy
 918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971
 972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025
 1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079
 1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133
 1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187
 1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241
 1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295
 1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349
 1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403
 1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457
 1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619
 1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673
 1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727
 1728
 1729
 1730
 1731
 1732
 1733
 1734
 1735
 1736
 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781
 1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835
 1836
 1837
 1838
 1839
 1840
 1841
 1842
 1843
 1844
 1845
 1846
 1847
 1848
 1849
 1850
 1851
 1852
 1853
 1854
 1855
 1856
 1857
 1858
 1859
 1860
 1861
 1862
 1863
 1864
 1865
 1866
 1867
 1868
 1869
 1870
 1871
 1872
 1873
 1874
 1875
 1876
 1877
 1878
 1879
 1880
 1881
 1882
 1883
 1884
 1885
 1886
 1887
 1888
 1889
 1890
 1891
 1892
 1893
 1894
 1895
 1896
 1897
 1898
 1899
 1900
 1901
 1902
 1903
 1904
 1905
 1906
 1907
 1908
 1909
 1910
 1911
 1912
 1913
 1914
 1915
 1916
 1917
 1918
 1919
 1920
 1921
 1922
 1923
 1924
 1925
 1926
 1927
 1928
 1929
 1930
 1931
 1932
 1933
 1934
 1935
 1936
 1937
 1938
 1939
 1940
 1941
 1942
 1943
 1944
 1945
 1946
 1947
 1948
 1949
 1950
 1951
 1952
 1953
 1954
 1955
 1956
 1957
 1958
 1959
 1960
 1961
 1962
 1963
 1964
 1965
 1966
 1967
 1968
 1969
 1970
 1971
 1972
 1973
 1974
 1975
 1976
 1977
 1978
 1979
 1980
 1981
 1982
 1983
 1984
 1985
 1986
 1987
 1988
 1989
 1990
 1991
 1992
 1993
 1994
 1995
 1996
 1997
 1998
 1999
 2000
 2001
 2002
 2003
 2004
 2005
 2006
 2007
 2008
 2009
 2010
 2011
 2012
 2013
 2014
 2015
 2016
 2017
 2018
 2019
 2020
 2021
 2022
 2023
 2024
 2025
 2026
 2027
 2028
 2029
 2030
 2031
 2032
 2033
 2034
 2035
 2036
 2037
 2038
 2039
 2040
 2041
 2042
 2043
 2044
 2045
 2046
 2047
 2048
 2049
 2050
 2051
 2052
 2053
 2054
 2055
 2056
 2057
 2058
 2059
 2060
 2061
 2062
 2063
 2064
 2065
 2066
 2067
 2068
 2069
 2070
 2071
 2072
 2073
 2074
 2075
 2076
 2077
 2078
 2079
 2080
 2081
 2082
 2083
 2084
 2085
 2086
 2087
 2088
 2089
 2090
 2091
 2092
 2093
 2094
 2095
 2096
 2097
 2098
 2099
 2100
 2101
 2102
 2103
 2104
 2105
 2106
 2107
 2108
 2109
 2110
 2111
 2112
 2113
 2114
 2115
 2116
 2117
 2118
 2119
 2120
 2121
 2122
 2123
 2124
 2125
 2126
 2127
 2128
 2129
 2130
 2131
 2132
 2133
 2134
 2135
 2136
 2137
 2138
 2139
 2140
 2141
 2142
 2143
 2144
 2145
 2146
 2147
 2148
 2149
 2150
 2151
 2152
 2153
 2154
 2155
 2156
 2157
 2158
 2159
 2160
 2161
 2162
 2163
 2164
 2165
 2166
 2167
 2168
 2169
 2170
 2171
 2172
 2173
 2174
 2175
 2176
 2177
 2178
 2179
 2180
 2181
 2182
 2183
 2184
 2185
 2186
 2187
 2188
 2189
 2190
 2191
 2192
 2193
 2194
 2195
 2196
 2197
 2198
 2199
 2200
 2201
 2202
 2203
 2204
 2205
 2206
 2207
 2208
 2209
 2210
 2211
 2212
 2213
 2214
 2215
 2216
 2217
 2218
 2219
 2220
 2221
 2222
 2223
 2224
 2225
 2226
 2227
 2228
 2229
 2230
 2231
 2232
 2233
 2234
 2235
 2236
 2237
 2238
 2239
 2240
 2241
 2242
 2243
 2244
 2245
 2246
 2247
 2248
 2249
 2250
 2251
 2252
 2253
 2254
 2255
 2256
 2257
 2258
 2259
 2260
 2261
 2262
 2263
 2264
 2265
 2266
 2267
 2268
 2269
 2270
 2271
 2272
 2273
 2274
 2275
 2276
 2277
 2278
 2279
 2280
 2281
 2282
 2283
 2284
 2285
 2286
 2287
 2288
 2289
 2290
 2291
 2292
 2293
 2294
 2295
 2296
 2297
 2298
 2299
 2300
 2301
 2302
 2303
 2304
 2305
 2306
 2307
 2308
 2309
 2310
 2311
 2312
 2313
 2314
 2315
 2316
 2317
 2318
 2319
 2320
 2321
 2322
 2323
 2324
 2325
 2326
 2327
 2328
 2329
 2330
 2331
 2332
 2333
 2334
 2335
 2336
 2337
 2338
 2339
 2340
 2341
 2342
 2343
 2344
 2345
 2346
 2347
 2348
 2349
 2350
 2351
 2352
 2353
 2354
 2355
 2356
 2357
 2358
 2359
 2360
 2361
 2362
 2363
 2364
 2365
 2366
 2367
 2368
 2369
 2370
 2371
 2372
 2373
 2374
 2375
 2376
 2377
 2378
 2379
 2380
 2381
 2382
 2383
 2384
 2385
 2386
 2387
 2388
 2389
 2390
 2391
 2392
 2393
 2394
 2395
 2396
 2397
 2398
 2399
 2400
 240

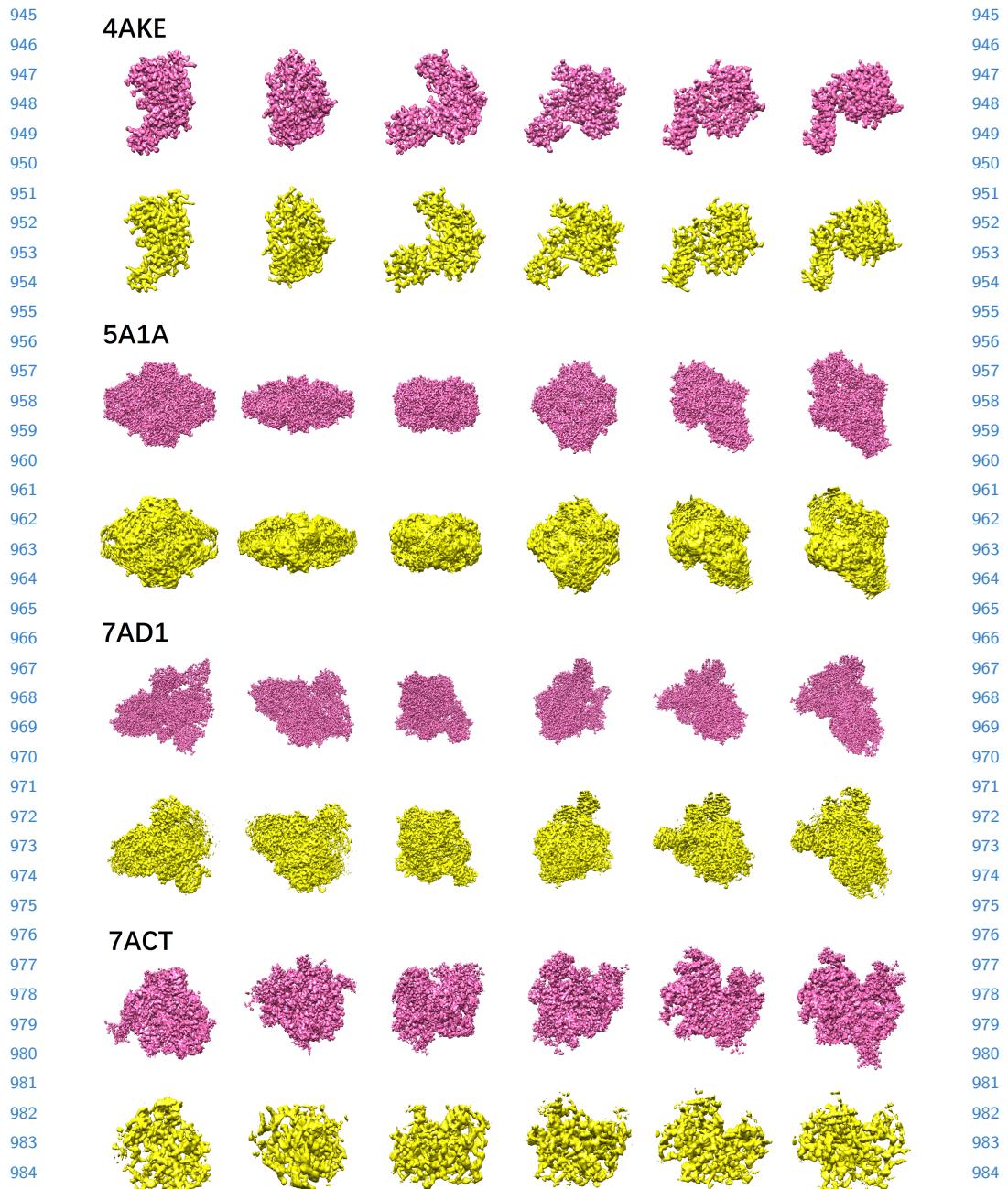


Fig. 9. Visualization of reconstructed density maps. SNR for 5A1A is -20 dB, while those for others are -10 dB. The pink ones are ground truths while the yellow ones are our reconstructions.

Table 4. The reconstruction resolutions of Cryo2SSNet. The results are compared with ones reported by previous proposed methods.

Protein	SNR (dB)	Resolution (\AA)	Previous Resolution(\AA)
4AKE	∞	1.74	/
	10	1.99	/
	0	1.94	2.57 (CryoPoseNet)
	-10	2.20	/
5A1A	-5.2	2.68	8.30 (CryoGAN)
	-10	3.14	/
	-20	4.64	15.30 (CryoGAN)
7AD1	∞	2.75	/
	-10	3.38	/
7ACT	∞	1.56	/
	-10	2.31	/

conduct reconstruction with the same hyper-parameters described in Section B.2 except creating couple of new datasets with density map resolutions 3.0 \AA and 5.0 \AA . As shown in Fig. 10, Cryo2SSNet works well among all these settings. It achieves 2.90 \AA and 3.90 \AA resolution for datasets of SNR=−10 dB with density map resolution 3.0 \AA and 5.0 \AA , respectively, which degrades an acceptable fraction under gradually worse datasets.

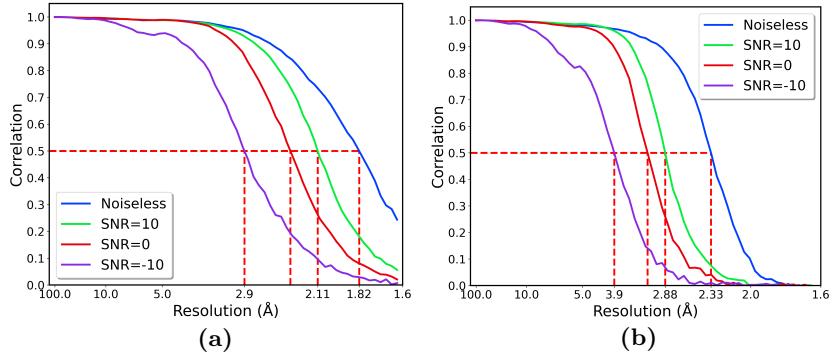


Fig. 10. Reconstruction of 4AKE datasets from density map resolution of 3.0 \AA (a) and 5.0 \AA (b). Cryo2SSNet performs consistently well and generalizes among different SNR.

Visualization of Loss Landscape To illustrate the reason why it is feasible to perform gradient descent enhancement after beam search rather than from scratch, we visualize the loss landscape when fixing the density map and changing the projection orientation. We pick one image from the 5A1A dataset with SNR=−20 dB and fetch the ground truth pose and density map. While the

quaternion representation of the pose is on a \mathbb{R}^4 sphere, two dimensions are perturbed in certain ranges (and normalized) to form a grid. At each grid point, the unbiased ℓ_2 loss between the projection and the original noisy image is calculated, by which a loss landscape is plotted. We change the perturbation range from 2.5 down to 2.5×10^{-3} , and the results are shown in Fig. 11.

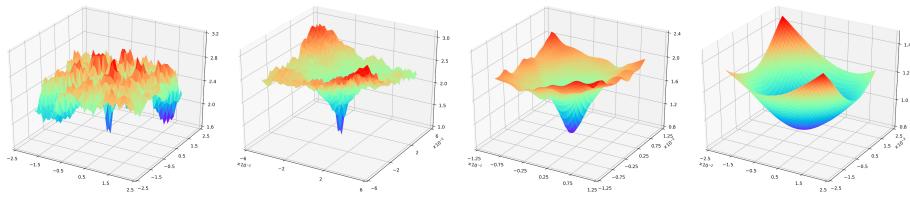


Fig. 11. Visualization of loss landscape. The perturbation ranges are from 2.5 (left most) to 2.5×10^{-3} (right most).

As we can see, the ground truth poses almost locates at the global minimum, while the landscape outside this valley is in a mess. Even the losses at points near the ground truth can be pretty large. When zooming in, we can see the landscape is becoming smooth and even convex, which is a friendly condition for gradient descent. Therefore, DHBS performs gradient descent enhancement after searching into a narrow region, otherwise, it cannot reach a good enough destination.