

Title: CryoDRGN: Reconstruction of heterogeneous structures from cryo-electron micrographs using neural networks

Authors: Ellen D. Zhong^{1,2}, Tristan Bepler^{1,2}, Bonnie Berger^{2,3*}, Joseph H. Davis^{1,4*}

Author information: ¹Computational and Systems Biology, ²Computer Science and Artificial Intelligence Laboratory, ³Department of Mathematics, ⁴Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139.

*Correspondence: bab@mit.edu, jhdavis@mit.edu

Running Title: Determination of highly heterogeneous molecular structures.

Keywords: cryo-electron microscopy, macromolecular complexes, structural biology, machine learning.

27,212 characters (including spaces), 3,891 words, 6 figures, 7 supplemental figures

Abstract

Cryo-EM single-particle analysis has proven powerful in determining the structures of rigid macromolecules. However, many protein complexes are flexible and can change conformation and composition as a result of functionally-associated dynamics. Such dynamics are poorly captured by current analysis methods. Here, we present cryoDRGN, an algorithm that for the first time leverages the representation power of deep neural networks to efficiently reconstruct highly heterogeneous complexes and continuous trajectories of protein motion. We apply this tool to two synthetic and three publicly available cryo-EM datasets, and we show that cryoDRGN provides an interpretable representation of structural heterogeneity that can be used to identify discrete states as well as continuous conformational changes. This ability enables cryoDRGN to discover previously overlooked structural states and to visualize molecules in motion.

1 **Main**

2 Single particle cryo-electron microscopy (cryo-EM) is a rapidly maturing method for high-
3 resolution structure determination of large macromolecular complexes^{1,2}. Major advances in
4 hardware³⁻⁵ and software⁴⁻⁹ have streamlined the collection and analysis of cryo-EM datasets, such
5 that structures of rigid macromolecules can routinely be solved at near atomic resolution^{10,11}.
6 However, a major computational bottleneck remains when conformational or compositional
7 heterogeneity is present in the sample.

8 The crux of cryo-EM structure determination is the computational task of reconstruction,
9 where algorithms must learn the 3D density or densities from the recorded dataset of 2D particle
10 images¹². While the standard formulation of reconstruction assumes that each 2D image is
11 generated from a single, static structure, in reality, each image contains a unique snapshot of the
12 molecule of interest. While this heterogeneity complicates reconstruction, it presents an
13 opportunity for single particle cryo-EM to reveal the conformational landscape of dynamic
14 macromolecular complexes.

15 Existing tools for heterogeneous reconstruction often make strong assumptions on the type
16 of heterogeneity in the dataset. Most commonly, heterogeneity is modeled as though it originates
17 from a small number of independent, *discrete* states¹³⁻¹⁶, consistent with molecules undergoing
18 cooperative conformational changes. However, because the number of underlying structural states
19 are unknown, such discrete classification approaches are error-prone and often result in the
20 omission of potentially relevant structures. Moreover, this approach fails to model molecules that
21 undergo continuous conformational changes. In these conformationally heterogeneous systems,
22 user-defined masks have been employed to resolve isolated rigid-body motions¹⁷, but these
23 approaches require assumptions on the types and location of molecular motions. Additionally, new

24 theoretical methods have been proposed to model global continuous heterogeneity¹⁸⁻²⁰, however
25 no such tools have been made available.

26 Here, we present cryoDRGN (Deep Reconstructing Generative Networks), a cryo-EM
27 reconstruction method that uses a deep neural network to directly approximate the molecule's
28 continuous 3D density function (**Fig. 1a**). We designed this tool based on the reasoning that deep
29 neural networks, which are known for their ability to model continuous, nonlinear functions, might
30 effectively capture dynamical structures. We show that this neural network representation of
31 structure, which we call a *deep coordinate network*, can efficiently learn heterogeneous ensembles
32 of high-resolution structures from single particle cryo-EM datasets²¹.

33 To learn this representation, cryoDRGN introduces an image-encoder/volume-decoder
34 framework to learn a latent representation of heterogeneity from single particle cryo-electron
35 micrographs. Once trained, users can visualize the dataset in the low-dimensional latent space,
36 which we find reflects the structural heterogeneity of the imaged molecule. This structural
37 heterogeneity can then be interrogated by generating 3D density maps at an arbitrary number of
38 desired positions within the latent space, which can be used to visualize continuous structural
39 trajectories.

40 CryoDRGN is a powerful and general approach for analyzing heterogeneity in imaging
41 datasets and can be used to reconstruct both compositionally and conformationally heterogeneous
42 structures. We demonstrate its efficacy by reconstructing and analyzing structures of the
43 eukaryotic ribosome, the assembling bacterial ribosome, and the pre-catalytic spliceosome
44 complex. In these machines, we discover new conformational states and observe dynamic
45 molecular motions. CryoDRGN is distributed as an open-source tool²² that can be easily integrated
46 in existing pipelines and is freely available at cryodrgn.csail.mit.edu.

47 **Results**

48 CryoDRGN architecture and training

49 CryoDRGN performs heterogeneous reconstruction by learning a neural network
50 representation of 3D structure from single particle cryo-EM micrographs. In contrast to traditional
51 reconstruction algorithms, which represent the 3D density map on a discretized voxel array,
52 cryoDRGN uses a neural network to predict density as a function of 3D Cartesian coordinates. We
53 call this architecture²³⁻²⁵ a *deep coordinate network*. To model heterogeneity, the deep coordinate
54 network can be extended to predict density as a function of both 3D coordinates and continuous
55 latent variables, z , which define a n -dimensional manifold of heterogeneous structures (**Fig. 1a**).
56 Coordinates are featurized with a positional encoding function before they are input to the deep
57 coordinate network (Methods). This choice of model assumes that structures can be embedded
58 within a continuous low-dimensional space, *i.e.* the latent space, where the dimensionality of the
59 latent space is defined by the user.

60 To train this neural network representation of 3D structure from a single particle cryo-EM
61 dataset, we develop an encoder–decoder architecture based on the Variational Autoencoder
62 (VAE)^{26,27} (**Fig. 1b**). The structure is represented in the Fourier domain in order to relate 2D
63 images as central slices out of a density map according to the Fourier slice theorem²⁸. For a given
64 image, X , the encoder neural network predicts a distribution of possible latent variable values,
65 $q(z|X)$. The image's oriented 3D pixel coordinates are computed from the image's pose
66 assignment provided by a previously determined consensus reconstruction. Then, given a sample
67 from the encoder distribution $z \sim q(z|X)$ and the 3D coordinates of the slice, the image is
68 reconstructed pixel-by-pixel through the deep coordinate network. The networks are trained jointly
69 using an objective function that seeks to optimize a variational upper bound on the data likelihood
70 as in standard VAEs²⁵. This objective function consists of the image reconstruction error and a

71 regularization term on the latent space. The parameters of the neural networks are iteratively
72 updated by gradient descent on this objective function.

73 After training, the encoder network is used to map images into the low-dimensional latent
74 space, where we define $\hat{z} = \operatorname{argmax}_z q(z|X)$ as each image's "latent encoding" (**Fig. 1c**). The
75 full distribution of latent encodings can then be visualized to study the latent space data manifold.
76 To explore the ensemble of structures, the deep coordinate network can directly reconstruct a 3D
77 density map given a desired value of the latent variable z and the 3D coordinates of a voxel array.

78 Deep coordinate networks can learn static structures from homogeneous datasets

79 To test the efficacy of neural networks in representing 3D structure, we first trained a deep
80 coordinate network with no latent variable input to learn the homogenous structure of the
81 *Plasmodium falciparum* 80S (*Pf*80S) ribosome from the EMPIAR 10028 dataset²⁹. The network
82 was trained on full resolution images ($D=360$, Nyquist limit of 2.7 Å), where image poses were
83 obtained from a consensus reconstruction in cryoSPARC³⁰. We found that the deep coordinate
84 network produced a structure qualitatively matching the traditional reconstruction (**Fig. 2a**) at
85 resolutions up to ~4.0 Å at an FSC=0.5 threshold (**Fig. 2b**).

86 As neural networks have a fixed capacity for representation that is constrained by their
87 architecture, we compared architectures of different sizes to evaluate the tradeoff between
88 representation power and training speed. We found that larger architectures converged to lower
89 values of the objective function (**Fig. 2c**) and correlated with the traditionally reconstructed map
90 at higher resolution (**Fig. 2b**). These improvements in the resulting structure came at the cost of
91 extended training times, suggesting that the architecture and the image size should be tuned to suit
92 the desired balance of training speed and achievable resolution (**Supplementary Fig. 1**).

93 CryoDRGN models both discrete and continuous structural heterogeneity

94 We next used simulated single particle cryo-EM datasets to test if the complete cryoDRGN
95 framework could perform *heterogeneous* reconstruction. This simulation-based approach allowed
96 us to quantitatively evaluate the method's performance by comparing the reconstructed density
97 maps to the ground truth density maps. To simulate continuous motions, we constructed an atomic
98 model of a hypothetical protein complex and iteratively rotated one bond's dihedral angle,
99 resulting in a series of 50 distinct but closely-related atomic models. We then generated density
100 maps along this reaction coordinate to serve as the ground truth density maps (**Fig. 3a**). Cryo-EM
101 micrographs were generated by projecting the ground truth maps with random poses, followed by
102 application of the contrast transfer function (CTF) and the addition of noise (see Methods). To
103 simulate a compositionally heterogeneous dataset, this procedure was repeated by mixing images
104 generated from the bacterial 30S, 50S, and 70S ribosomal density maps (**Fig. 3d**). The cryoDRGN
105 networks were then provided these simulated images and their corresponding poses, and were
106 trained with 1-dimensional (1D) latent variable models.

107 When trained on the dataset with continuous heterogeneity, we found that cryoDRGN
108 accurately modeled the full continuum of structures as assessed by two criteria. First, the latent
109 encoding of each image produced by the encoder network correlated well with the dihedral angle
110 of the underlying model (Spearman $r = -0.996$), which we characterize as the ground truth
111 reaction coordinate (**Fig. 3b**). Second, when provided a series of latent variable values, the deep
112 coordinate network produced structures that correlated with the ground-truth maps (**Fig. 3c**). We
113 note that the deep coordinate network can generate an arbitrary number of conformations along
114 the trajectory, and found that for 100 images equally spaced along the reaction coordinate, the
115 generated structure at each image's predicted latent encoding correlated well with its ground truth
116 map (**Supplementary Fig. 2**).

117 When cryoDRGN was trained on the compositionally heterogeneous dataset, we observed
118 that the encoder network mapped particles to three distinct clusters in latent space (**Fig. 3e**). These
119 clusters aligned with the ground truth class assignments from the 30S, 50S, and 70S ribosome
120 (classification accuracy of 99.9%), and the appropriate ribosomal structures were generated by the
121 deep coordinate network when provided with latent variable values at the corresponding cluster
122 centers (**Fig. 3f, Supplementary Fig. 2**).

123 CryoDRGN uncovers residual heterogeneity in a high-resolution cryo-EM reconstruction

124 We next evaluated cryoDRGN's ability to learn heterogeneous structures from real cryo-
125 EM data, which contains structured noise and imaging artifacts that are difficult to simulate. When
126 analyzing a homogeneous reconstruction of the *Pf80S* ribosome, Wong *et al.* observed flexibility
127 in the small subunit head region and missing density for peripheral rRNA expansion segment
128 elements that prevented completion of an atomic model in these regions²⁹. To explore if this
129 unresolved density resulted from residual heterogeneity, we trained a 10-dimensional (10D) latent
130 variable model with cryoDRGN on their deposited dataset (EMPIAR-10028), using poses from a
131 consensus reconstruction in cryoSPARC³⁰. We then visualized the dataset's latent encodings using
132 principal component analysis (PCA) (**Fig. 4a**) and observed a subset of particles separated along
133 PC1. A density map generated by the deep coordinate network from this region of latent space
134 revealed a distinct conformation of the 40S subunit, which was rotated relative to the 60S subunit
135 (**Fig. 4b,c**). Concomitant with the inter-subunit rotation, we observed the disappearance of the
136 inter-subunit bridge formed by the C-terminal helix of eL8, which is consistent with Sun *et al.*'s
137 characterization of *Pf80S* dynamics³¹. We further explored structural heterogeneity by performing
138 *k*-means clustering of the latent encodings with *k*=20 clusters and subsequently generating
139 structures at the cluster centers. We observed diverse structures including those bearing a rotated

140 *Pf*40S head group, those missing the head group, and those with clearly resolved rRNA helices
141 that were absent from the homogeneous reconstruction (**Fig. 4c**).

142 CryoDRGN automatically partitions assembly states of the bacterial ribosome

143 Next, we assessed cryoDRGN's ability to analyze and reconstruct density maps from a
144 dataset known to contain substantial compositional and conformational heterogeneity. For this
145 assessment, we investigated a highly heterogeneous mixture of assembly intermediates of the *E.*
146 *coli* large ribosomal subunit (LSU). This dataset (EMPIAR 10076) had previously been analyzed
147 through multiple expert-guided rounds of hierarchical 3D classification resulting in 13 discrete
148 structures that were grouped into 4 major classes³². These particles were obtained by crudely
149 fractionating a lysate with the explicit goal of imaging and later analyzing the full ensemble of
150 cellular assembly intermediates. As such, a substantial fraction of the published particle stack
151 corresponds to non-ribosomal impurities that were discarded during 3D classification in the
152 original analysis (26,575 out of 131,899 images). Despite this heterogeneity, a homogeneous
153 reconstruction of the full dataset produced a consensus structure of the mature LSU (GSFSC
154 resolution of 3.2 Å), suggesting that even in the presence of these impurities, the heterogeneous
155 ribosomal particles could be aligned to the rigid core of the LSU, which enabled analysis using
156 cryoDRGN.

157 To assess the degree of heterogeneity in the data, we first trained a 1D latent variable model
158 on down-sampled images (D=128, Nyquist limit of 6.6 Å) using image poses from a consensus
159 reconstruction. After model training, the encoder network was used to predict the latent encoding
160 for each particle, and the resulting histogram of the full dataset's encodings revealed five distinct
161 peaks. Four of the peaks corresponded to each of the four major classes of the LSU, and the fifth
162 peak near $z = -2$ captured particles that were unassigned by Davis *et al.* (**Fig. 5a**)³². This clear
163 separation in latent space suggested that cryoDRGN can identify sample impurities without

164 supervision. When using the subset of particles from this region (assigned $z \leq -1$), neither 2D
165 class averages nor a traditional 3D reconstruction produced structures consistent with assembling
166 ribosomes (**Supplementary Fig. 3**). As we do not wish to model these impurities, we filtered the
167 dataset by the latent variable, keeping 101,604 images with $z > -1$ for further analysis.

168 To explore the heterogeneity within these major assembly states, we trained a 10D latent
169 variable model on the remaining high-resolution images ($D=256$, Nyquist limit of 3.3 Å). We
170 visualized the resulting 10D latent encodings using UMAP³³, and observed particle super-clusters
171 corresponding to major classes of LSU assembly (**Fig. 5b**) and sub-clusters that aligned with Davis
172 *et al.*'s minor class assignments (**Supplementary Fig. 4**)³². We found that when provided latent
173 codes from these clusters, the decoder network generated structures matching the major (**Fig. 5c**)
174 and minor (**Supplementary Fig. 5**) assembly states of the LSU. With the 10D latent variable
175 model, we also noted a clearly separated cluster of particles assigned to class A, and structures
176 sampled from this region of latent space reconstructed the 70S ribosome, an impurity in the dataset
177 (**Supplementary Fig. 6**). Finally, we identified a small cluster of ~1,200 particles in latent space
178 adjacent to the class C cluster whose particles were classified into class E by Davis *et al.*
179 (**Supplementary Fig. 4**). The density map reconstructed by the deep coordinate network from this
180 region revealed a previously unreported assembly intermediate that we newly call class C4. Like
181 the other class C structures, class C4 lacked the central protuberance, but bore clearly resolved
182 density for helix 68, which was only present in classes E4 and E5 from Davis *et al.*³². Traditional
183 voxel-based back-projection of the particle images constituting this cluster reproduced a similar,
184 albeit lower-resolution structure, confirming the existence of this structural state in the original
185 dataset (**Supplementary Fig. 6**).

186 CryoDRGN reveals dynamic continuous motions in the pre-catalytic spliceosome

187 Finally, we evaluated the performance of cryoDRGN in analyzing micrographs of the pre-
188 catalytic spliceosome (EMPIAR 10180)³⁴. Plaschka *et al.* employed extensive expert-guided
189 focused classifications to reconstruct a composite map for this complex and suggested that the
190 complex sampled a continuum of conformations³⁴. To understand how cryoDRGN would encode
191 such continuous structural heterogeneity in latent space, we first trained a 10D latent variable
192 model on the downsampled images (D=128, Nyquist limit of 8.5 Å) using image poses derived
193 from a consensus reconstruction. Multiple clusters were observed in the latent space encodings
194 (**Fig. 6a**). After sampling structures from the latent space, we observed expected spliceosome
195 conformations from the largest cluster, poorly resolved structures from the leftmost cluster,
196 structures lacking density for the SF3b domain from a third cluster, and additional density
197 consistent with particle aggregation from the uppermost cluster (**Fig. 6b**). To focus our analysis
198 on bone-fide pre-catalytic spliceosome particles, we leveraged the latent space representation to
199 eliminate any particles that mapped to the undesired clusters from two replicate runs, which
200 resulted in a final particle stack of 150,098 images.

201 With the filtered particle stack, we trained a 10D model on higher resolution images
202 (D=256, Nyquist limit of 3.4 Å), and visualized the dataset's latent encodings in 2D using PCA
203 and UMAP (**Fig. 6a,c**). The visualized data manifold was unfeatured, consistent with a molecule
204 undergoing non-cooperative conformational changes. By generating structures along the first
205 principal component of the latent space encodings, we reconstructed a trajectory of the SF3b and
206 helicase domains in motion (**Fig. 6d**), which smoothly transitioned from an elongated state to one
207 compressed against the body of the spliceosome. A similar traversal along the second PC produced
208 a continuous trajectory of the SF3b and helicase domains moving in opposition (**Supplementary**
209 **Fig. 7**). This anticorrelated motion of the SF3b and helicase domains in PC2, together with their
210 correlated motion in PC1, suggested that the two domains move independently in the imaged

211 ensemble. Finally, although trajectories along latent space PCs provide a summary of the extent of
212 variability in the structure, cryoDRGN can also generate structures at arbitrary points from the
213 latent space. By traversing along the k -nearest neighbor graph of the latent encodings and
214 generating structures at the visited nodes, cryoDRGN generated a plausible trajectory of the
215 conformations adopted by the pre-catalytic spliceosome (**Supplemental Movie 1**), highlighting
216 the potential of single particle cryo-EM to uncover the conformational dynamics of molecular
217 machines.

218 **Discussion**

219 This work introduces cryoDRGN, a new method using neural networks to reconstruct 3D
220 density maps from heterogeneous single particle cryo-EM datasets. The power of this approach
221 lies in its ability to represent heterogeneous structures without simplifying assumptions on the type
222 of heterogeneity. In principle, cryoDRGN is able to represent any distribution of structures that
223 can be approximated by a deep neural network, a broad class of function approximators³⁵. This
224 flexibility contrasts with existing methods that impose strong assumptions on the types of
225 structural heterogeneity present in the sample. For example, traditional 3D classification assumes
226 a mixture of discrete structural classes, whereas multibody refinement assumes conformational
227 changes are composed strictly of rigid-body motions. Although these approaches have proven
228 useful, they are inherently unable to model true structural heterogeneity and thus often introduce
229 bias into reconstructions. In contrast, we empirically show that deep coordinate networks can
230 model both discrete compositional heterogeneity and continuous conformational changes without
231 the aforementioned assumptions. For example, by using this less biased approach, we discovered
232 heterogeneous states of the *Pf80S* ribosome that were originally averaged out of the homogeneous
233 reconstruction. When analyzing the assembling *E. Coli* LSU dataset, cryoDRGN learned an
234 ensemble of LSU assembly states without *a priori* specification of the number of states as is

235 required for 3D classification. Finally, when analyzing the pre-catalytic spliceosome, we found
236 that the continuous conformational changes cryoDRGN reconstructed lack the rigid-body
237 boundary artifacts introduced from multibody refinement's mask-based approach¹⁷.

238 Interpretation of the latent space

239 A key feature of cryoDRGN is its ability to provide a low-dimensional representation of
240 the dataset's heterogeneity, which is given by each particle's latent encoding. Subject to
241 optimization, cryoDRGN organizes the latent space such that structurally related particles are in
242 close proximity. Thus, visualization of the distribution of latent encodings can be informative in
243 understanding the structural heterogeneity within the imaged ensemble. In both simulated and real
244 datasets we find that continuous motions are embedded along a continuum in latent space (**Fig.**
245 **3b, 6c**) and that compositionally distinct states manifest as clusters (**Fig. 3e, 5b**). This observation
246 suggests an interpretation of the latent encodings as an approximate conformational landscape,
247 with regions of high-particle occupancy corresponding to low-energy states, and regions of lower-
248 particle occupancy denoting higher energy states. We note however that structures reconstructed
249 from unoccupied regions will not in general correspond to true physical intermediates, as
250 cryoDRGN optimizes the likelihood of the observed data and these intermediates are not observed.
251 Finally, in real datasets, there may exist images that do not originate from the standard single
252 particle image formation model, for example, false positives encountered during particle picking⁹.
253 We demonstrated the utility of the latent space encodings in identifying such impurities, ice
254 artifacts, and other such out-of-distribution particles that may be filtered out in subsequent analyses
255 (**Fig. 5a, 6a**).

256 Visualizing structural trajectories

257 In addition to encoding particles in an unsupervised manner, cryoDRGN can reconstruct
258 3D density maps from user-defined positions in latent space. Because cryoDRGN learns a

259 generative model for structure, an unlimited number of structures can be generated and analyzed,
260 thus enabling visualization of structural trajectories. By leveraging the latent encodings of the
261 particle images, users can directly traverse the data manifold and only sample structures from
262 regions of latent space with significant particle occupancy. Indeed, we applied a well-established
263 graph-traversal algorithm³⁶ to visualize a data-supported path of the *Pf*80S ribosome, bL17-
264 independent assembly of the bacterial ribosome, and the pre-catalytic spliceosome (**Supplemental**
265 **Movies 1,2,3,4**).

266 Practical considerations in choosing training hyperparameters

267 Although this method emphasizes an unsupervised approach to analyzing structural
268 heterogeneity, cryoDRGN does require that the user define the dimensionality of the latent space
269 and the architecture of both the encoder and decoder networks. We find that in practice, a 1D latent
270 space is effective at distinguishing bona-fide particles from contaminants and imaging artifacts
271 (**Fig. 5a**), and we recommend users initially employ such a model to filter their dataset.
272 Additionally, we find that in our tested datasets, a 10D latent space provides sufficient
273 representation capacity to effectively model structural heterogeneity, and that this 10D space can
274 be readily visualized with PCA or UMAP. Notably, we recommend the use of such as 10D latent
275 space instead of lower dimensional space as we have found that 10D spaces result in much more
276 rapid overall training, which is consistent with similar observations of related overparameterized
277 neural network architectures^{37,38}. Finally, users must specify the number of nodes and layers in the
278 neural networks. Here, we find an inverse relationship between neural network size and the
279 achievable resolution of a given structure (**Supplemental Fig. 1**). Training larger networks on
280 larger images is significantly slower, and we recommend that users perform an initial assessment
281 using down-sampled images and relatively small networks before proceeding to high-resolution
282 reconstructions.

283 Discovering new states using cryoDRGN

284 CryoDRGN can be used to identify novel clusters of structurally-related particles, which
285 can then be visualized by sampling a 3D structure from that region of latent space. Indeed, in
286 analyzing the bL17-depleted LSU assembly dataset, we noted a completely new structural class,
287 which like the C-classes, lacked the central protuberance, but like the most mature E classes,
288 clearly bore a functionally critical inter-subunit helix (h68). This state was completely missed in
289 traditional hierarchical classification³², and provides structural evidence that this vital intersubunit
290 helix can dock in a native conformation in the absence of the central protuberance
291 (**Supplementary Fig. 6**). Notably, we could validate the existence of this class by performing
292 traditional back-projection using ~1,000 particles from this cluster (**Supplementary Fig. 6**).

293 In future work, we envision using cryoDRGN to reveal the number of discrete classes, their
294 constituent particles, and to produce initial 3D models that could be used as inputs for a traditional
295 3D reconstruction. Given the mature state of such tools^{39,40}, this unbiased classification approach
296 followed by traditional homogeneous reconstruction, particle polishing, and higher order image
297 aberration correction, has the potential to produce very high-resolution structures of the full
298 spectrum of discrete structural states without the need for expert-guided classification.

299 Fully unsupervised 3D reconstruction

300 As implemented, cryoDRGN uses pose estimates resulting from a traditional consensus 3D
301 reconstruction. In analyzing three publicly available datasets, we found that such consensus pose
302 estimates were sufficiently accurate to generate meaningful latent space encodings and to produce
303 interpretable density maps of distinct structures. It is clear, however, that this approach will fail if
304 the degree of structural heterogeneity in the dataset results in inaccurate pose estimates. For
305 example, a mixture of structurally unrelated complexes will align poorly to a consensus structure,
306 and thus produce poor pose estimates. Notably, our framework is differentiable with respect to

307 pose variables, which, in principle, should allow for on-the-fly pose-refinement or *de novo* pose
308 estimation²⁵, and future work will explore the efficacy of incorporating such features.

Acknowledgments

The authors thank Ben Demeo, Ashwin Narayan, Adam Lerer, Roy Lederman, Sam Rodriques, Bob Sauer, Phil Sharp, and Kotaro Kelley for helpful discussions and feedback. This work was funded by the National Science Foundation Graduate Research Fellowship Program, NIH grant R01-GM081871 to BB, NIH grant R00-AG050749 to JD, and the MIT J-Clinic for Machine Learning and Health to JD and BB.

Author Contributions

EZ, TB, BB, and JD conceived of the work. EZ, TB, and BB developed the representation learning method. EZ and JD tailored the method to cryo-EM data, and designed and analyzed the described experiments. EZ implemented the software and performed the experiments. EZ, JD, and BB wrote the manuscript.

Competing Interests Statement

The authors declare no competing financial interests.

References

1. Nogales, E. The development of cryo-EM into a mainstream structural biology technique. *Nat. Methods* **13**, 24–27 (2016).
2. Cheng, Y. Single-particle cryo-EM—How did it get here and where will it go. *Science* **361**, 876–880 (2018).
3. Bammes, B. E., Rochat, R. H., Jakana, J., Chen, D.-H. & Chiu, W. Direct electron detection yields cryo-EM reconstructions at resolutions beyond 3/4 Nyquist frequency. *J. Struct. Biol.* **177**, 589–601 (2012).

4. Suloway, C. *et al.* Automated molecular microscopy: The new Legimon system. *J. Struct. Biol.* **151**, 41–60 (2005).
5. Li, X. *et al.* Electron counting and beam-induced motion correction enable near-atomic-resolution single-particle cryo-EM. *Nat. Methods* **10**, 584–590 (2013).
6. Zhang, K. Gctf: Real-time CTF determination and correction. *J. Struct. Biol.* **193**, 1–12 (2016).
7. Brubaker, M. A., Punjani, A. & Fleet, D. J. Building Proteins in a Day: Efficient 3D Molecular Reconstruction. 3099–3108 (2015).
8. Scheres, S. H. W. A Bayesian view on cryo-EM structure determination. *J. Mol. Biol.* **415**, 406–418 (2012).
9. Bepler, T. *et al.* Positive-unlabeled convolutional neural networks for particle picking in cryo-electron micrographs. *Nat Methods* **16**, 1153–1160 (2019).
10. Ahmed, T., Yin, Z. & Bhushan, S. Cryo-EM structure of the large subunit of the spinach chloroplast ribosome. *Sci Rep* **6**, 1–13 (2016).
11. Wrapp, D. *et al.* Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science* **367**, 1260–1263 (2020).
12. Sigworth, F. J. Principles of cryo-EM single-particle image processing. *Microscopy (Oxf)* **65**, 57–67 (2016).
13. Scheres, S. H. W. *et al.* Maximum-likelihood Multi-reference Refinement for Electron Microscopy Images. *J. Mol. Biol.* **348**, 139–149 (2005).
14. Lyumkis, D., Brilot, A. F., Theobald, D. L. & Grigorieff, N. Likelihood-based classification of cryo-EM images using FREALIGN. *J. Struct. Biol.* **183**, 377–388 (2013).

15. Scheres, S. H. W. *et al.* Disentangling conformational states of macromolecules in 3D-EM through likelihood optimization. *Nat. Methods* **4**, 27–29 (2007).
16. Haselbach, D. *et al.* Structure and Conformational Dynamics of the Human Spliceosomal Bact Complex. *Cell* **172**, 454–464.e11 (2018).
17. Nakane, T., Kimanius, D., Lindahl, E. & Scheres, S. H. Characterisation of molecular motions in cryo-EM single-particle data by multi-body refinement in RELION. *Elife* **7**, e36861 (2018).
18. Frank, J. & Ourmazd, A. Continuous changes in structure mapped by manifold embedding of single-particle data in cryo-EM. *Methods* **100**, 61–67 (2016).
19. Moscovich, A., Halevi, A., Andén, J. & Singer, A. Cryo-EM reconstruction of continuous heterogeneity by Laplacian spectral volumes. *arXiv.org eess.IV*, (2019).
20. Lederman, R. R. & Singer, A. Continuously heterogeneous hyper-objects in cryo-EM and 3-D movies of many temporal dimensions. *arXiv.org cs.CV*, (2017).
21. Iudin, A., Korir, P. K., Salavert-Torres, J., Kleywegt, G. J. & Patwardhan, A. EMPIAR: a public archive for raw electron microscopy image data. *Nat. Methods* **13**, 387–388 (2016).
22. GPLv3 GNU General Public License. Free Software Foundation (2007).
23. Bricman, P. A. & Ionescu, R. T. CocoNet: A deep neural network for mapping pixel coordinates to color values. *arXiv.org cs.CV*, (2018).
24. Bepler, T., Zhong, E., Kelley, K., Brignole, E. & Berger, B. Explicitly disentangling image content from translation and rotation with spatial-VAE. *Advances in Neural Information Processing Systems* **32** 15435–15445 (2019).

25. Zhong, E. D., Bepler, T., Davis, J. H. & Berger, B. Reconstructing continuous distributions of 3D protein structure from cryo-EM images. In *International Conference on Learning Representations* (2020).
26. Kingma, D. P. & Welling, M. Auto-Encoding Variational Bayes. *arXiv.org stat.ML*, (2013).
27. Rezende, D. J., Mohamed, S. & Wierstra, D. Stochastic Backpropagation and Approximate Inference in Deep Generative Models. *arxiv.org* (2014).
28. Bracewell, R. N. Strip Integration in Radio Astronomy. *Aust. J. Phys.* **9**, 198–217 (1956).
29. Wong, W. *et al.* Cryo-EM structure of the Plasmodium falciparum 80S ribosome bound to the anti-protozoan drug emetine. *Elife* **3**, e01963 (2014).
30. Punjani, A., Rubinstein, J. L., Fleet, D. J. & Brubaker, M. A. cryoSPARC: algorithms for rapid unsupervised cryo-EM structure determination. *Nat. Methods* **14**, 290–296 (2017).
31. Sun, M. *et al.* Dynamical features of the Plasmodium falciparum ribosome during translation. *Nucleic Acids Research* **2**;43(21):10515-24 (2015).
32. Davis, J. H. *et al.* Modular Assembly of the Bacterial Large Ribosomal Subunit. *Cell* **167**, 1610–1622.e15 (2016).
33. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arxiv.org* (2018).
34. Plaschka, C., Lin, P.-C. & Nagai, K. Structure of a pre-catalytic spliceosome. *Nature* **546**, 617–621 (2017).
35. Hornik, K. Approximation capabilities of multilayer feedforward networks. *Neural Networks* **4**, 251–257 (1991).

36. Cormen, Thomas H.; Leiserson, Charles E.; Rivest, Ronald L.; Stein, Clifford (2001). "Section 24.3: Dijkstra's algorithm". *Introduction to Algorithms* (Second ed.). MIT Press and McGraw–Hill. pp. 595–601. ISBN 0-262-03293-7.
37. Zhang, C., Bengio, S., Hardt, M., Recht, B. & Vinyals, O. Understanding deep learning requires rethinking generalization. *arxiv.org* (2016).
38. Buhai, R.-D., Risteski, A., Halpern, Y. & Sontag, D. Benefits of Overparameterization in Single-Layer Latent Variable Generative Models. (2019).
39. Zivanov, J. *et al.* RELION-3: new tools for automated high-resolution cryo-EM structure determination. *bioRxiv* 421123 (2018). doi:10.1101/421123
40. Punjani, A., Zhang, H. & Fleet, D. J. Non-uniform refinement: Adaptive regularization improves single particle cryo-EM reconstruction. *bioRxiv* **179**, 2019.12.15.877092 (2019).

309 **Methods**

310 **The cryoDRGN method**

311 Deep coordinate networks to represent 3D structure

312 The cryoDRGN method performs heterogeneous cryo-EM reconstruction by learning a
313 neural network representation of 3D structure. In particular, we use a neural network to
314 approximate the function $V: \mathbb{R}^{3+n} \rightarrow \mathbb{R}$, which models structures as generated from an n -
315 dimensional continuous latent space. We call this architecture¹⁻³ a *deep coordinate network* as we
316 explicitly model the volume as a function of Cartesian coordinates.

317 Without loss of generality, we model volumes on the domain $[-0.5, 0.5]^3$. Instead of
318 directly supplying the 3D Cartesian coordinates, \mathbf{k} , to the deep coordinate network, coordinates
319 are featurized with a fixed positional encoding function consisting of sinusoids whose wavelengths
320 follow a geometric progression from 1 up to the Nyquist limit:

$$pe^{(2i)}(k_j) = \sin\left(k_j D \pi \left(\frac{2}{D}\right)^{\frac{i}{\left(\frac{D}{2}-1\right)}}\right), i = 0, \dots, \frac{D}{2} - 1; k_j \in \mathbf{k}$$

$$pe^{(2i+1)}(k_j) = \cos\left(k_j D \pi \left(\frac{2}{D}\right)^{\frac{i}{\left(\frac{D}{2}-1\right)}}\right), i = 0, \dots, \frac{D}{2} - 1; k_j \in \mathbf{k}$$

where D is set to the image size¹ used in training. Empirically, we found that excluding the highest frequencies of the positional encoding led to better performance when training on noisy data, and we provide an option to modify the positional encoding function by increasing all wavelengths by a factor of 2π .

Training system

This parametric representation of 3D structure is learned via an image-encoder/volume-decoder architecture based on the variational autoencoder (VAE)^{4,5}. We follow the standard image formation model in single particle cryo-EM³ where observed images are generated from projections of a volume at a random unknown orientation, $R \in SO(3)$. We use an additive Gaussian white noise model. Volume heterogeneity is generated from a continuous latent space, modeled by the latent variable \mathbf{z} , where the dimensionality of \mathbf{z} is a hyperparameter of the model.

Given an image X , the variational encoder, $q_{\xi}(\mathbf{z}|X)$, produces a mean and variance, $\mu_{\mathbf{z}|X}$ and $\Sigma_{\mathbf{z}|X}$, statistics that parameterize a Gaussian distribution with diagonal covariance, as the variational approximation to the true posterior $p(\mathbf{z}|X)$. The prior on the latent variable is a standard normal distribution $\mathcal{N}(0, \mathbf{I})$. The deep coordinate network architecture is used as the probabilistic decoder, $p_{\theta}(V|\mathbf{k}, \mathbf{z})$, and models structures in frequency space. Given Cartesian coordinate $\mathbf{k} \in \mathbb{R}^3$ and latent variable \mathbf{z} , the probabilistic decoder predicts a Gaussian distribution over $V(\mathbf{k}, \mathbf{z})$.

¹ Number of pixels along one dimension of the image, i.e. a $D \times D$ image

340 The encoder and decoder are parameterized with fully connected neural networks with parameters
341 ξ and θ , respectively.

342 Since 2D projection images can be related to volumes as 2D central slices in Fourier space⁶,
343 oriented 3D coordinates for a given image can be obtained by rotating a $D \times D$ lattice spanning
344 $[-0.5, 0.5]^2$ originally on the x-y plane by R , the orientation of the volume during imaging. Then,
345 given a sample out of $q_\xi(\mathbf{z}|X)$ and the oriented coordinates, an image can be reconstructed pixel-
346 by-pixel through the decoder. The reconstructed image is then translated by the image's in-plane
347 shift, and the CTF is applied before it is compared to the input image. The negative log likelihood
348 of a given image under our model is computed as the mean square error between the reconstructed
349 image and the input image. Following the standard VAE framework, the optimization objective is
350 the variational lower bound of the model evidence:

$$351 \quad \mathcal{L}(X; \xi, \theta) = E_{q_\xi(\mathbf{z}|X)}[\log p(X|\mathbf{z})] - KL(q_\xi(\mathbf{z}|X)||p(\mathbf{z}))$$

352 where the first term is the reconstruction error estimated with one Monte Carlo sample and the
353 second term is a regularization term on the latent representation. By training on many 2D slices
354 with sufficiently diverse orientations, the 3D volume can be learned through feedback from the 2D
355 views. For further details, we refer the reader to a preliminary version of the method described in
356 the proceedings of the International Conference for Learning Representations³. The results
357 presented here employ the training regime described in Zhong *et al.* using previously determined
358 poses from a consensus reconstruction³.

359 **Datasets**

360 *Simulated homogeneous dataset generation*

361 The 50S subunit of the *E. coli* ribosome was extracted from PDB 4YBB in PyMOL⁷. A
362 density map was generated from the atomic model using the molmap command in Chimera⁸ at a
363 grid spacing of 1.5 Å/pix and a resolution of 4.5 Å. The resulting volume was padded to a box size

364 of $D=256$, where D is the width in pixels along one dimension. Simulated micrographs were
365 generated with custom Python scripts as follows: 50k projection images were generated by rotating
366 the density map with a random rotation sampled uniformly from $SO(3)$, projecting along the z -
367 axis, and shifting the image with an in-plane translation sampled uniformly from $[-20,20]^2$ pixels.
368 Projection images were multiplied with the CTF in Fourier space, where the CTF was computed
369 from defocus values randomly sampled from those given in EMPIAR-10028, no astigmatism, an
370 accelerating voltage of 300 kV, a spherical aberration of 2mm, and an amplitude contrast ratio of
371 0.1. An envelope function with a B-factor of 100 \AA^2 was applied. Noise was added with a signal
372 to noise ratio (SNR) of 0.1 where the noise-free signal images were defined as the entire $D \times D$
373 image. To generate the dataset with $D=128$, the $D=256$ noiseless projection images of the 50S
374 were downsampled by Fourier clipping, followed by addition of CTF and noise as above.

375 Discrete3 heterogeneous dataset generation

376 To generate the “Discrete3” dataset, 10k, 15k, and 25k simulated micrographs of the 30S,
377 50S, and 70S ribosome, respectively, were combined. 15k micrographs from the homogeneous
378 50S dataset were used, and micrographs of the 30S and 70S ribosome were generated using the
379 same procedure starting from the atomic model extracted from PDB 4YBB, and extracting either
380 the 30S or 70S subunits. Images were downsampled to $D=128$, corresponding to a Nyquist limit
381 of 6 \AA .

382 Linear1D heterogeneous dataset generation

383 To generate the “Linear1D” dataset, 50 density maps were generated along a reaction
384 coordinate defined by rotation of a dihedral angle in an atomic model of a hypothetical protein
385 complex. Each model was generated at 0.03 radian increments of the bond rotation, leading to a
386 total range of 1.5 radians. Density maps were generated in Chimera at a grid spacing of 6 $\text{\AA}/\text{pix}$
387 and resolution of 12 \AA , and padded to a box size of $D=128$. 1000 projection images were generated

388 with random orientations and in-plane translations from $[-10,10]^2$ pixels for each map leading to a
389 final particle stack of 50k images. CTF and noise at an SNR=0.1 were added using the same
390 procedure described above.

391 Real cryo-EM datasets

392 Processed shiny particles and the star file containing CTF parameters were downloaded
393 from the Electron Microscopy Public Image Archive (EMPIAR) ⁹ for datasets EMPIAR-10028,
394 EMPIAR-10076, and EMPIAR-10180. Particle images were resized to either D=96, 128, or 256
395 by clipping in Fourier space with a custom Python script. These various images sizes resulted in
396 the following Nyquist limits:

Dataset name	EMPIAR ID	Image size, D (pixels)	Nyquist limit (Å)	Figure
80S ribosome	10028	96	10.1	N/A
80S ribosome	10028	256	3.8	4
Assembling LSU ribosome	10076	128	6.6	5
Assembling LSU ribosome	10076	256	3.3	5
Pre-catalytic spliceosome	10180	128	8.5	6
Pre-catalytic spliceosome	10180	256	4.3	6

397

398 **Traditional homogeneous reconstruction**

399 3D reconstruction of the 80S ribosome (EMPIAR-10028) was performed in cryoSPARC
400 v2.4¹⁰ using the ab-initio reconstruction job followed by the homogeneous refinement job with
401 default parameters. The final reconstruction reported a GSFSC_{0.143}¹¹ resolution of 3.1 Å with a
402 tight mask and 4.1 Å unmasked. The density map was sharpened using the published B-factor of
403 -80.1 Å² for visualization.

404 Homogeneous 3D reconstruction of the L17-depleted ribosome assembly intermediates
405 (EMPIAR 10076) was performed as above, leading to a final structure with a GSFSC_{0.143}¹¹
406 resolution of 3.2 Å with a tight mask and 4.8 Å unmasked.

407 **Deep coordinate network training on homogeneous structures**

408 For each dataset and for each architecture, a separate deep coordinate network with no
409 latent variable was trained for 25 epochs, where an epoch is defined as one pass through the dataset.
410 The tested architectures were fully connected networks with ReLU activations, where the network
411 size was either 3 layers of dimension 128 (128 nodes/layer x 3 layers), 3 layers of dimension 256
412 (256x3), 3 layers of dimension 1024 (1024x3), or 10 layers of dimension 1024 (1024x10). Image
413 poses were set to either the ground truth poses for the simulated datasets, or poses obtained from
414 a traditional homogeneous reconstruction in cryoSPARC. Networks were trained on minibatches
415 of 8 images using the Adam¹² optimizer with a learning rate of 0.0001. Once training completed,
416 the deep coordinate network was evaluated on the 3D coordinates of a $D \times D \times D$ voxel array
417 spanning $[-0.5, 0.5]^3$, where D is the image size in pixels along one dimension. The density map
418 was sharpened using the published B-factor of -80.1 \AA^2 for visualization¹³.

419 **Map-to-map FSC**

420 Fourier shell correlation curves were computed between the ground truth density maps and
421 the neural network reconstructed density maps using a custom Python script. For the homogeneous
422 reconstruction of EMPIAR 10028, the map-to-map FSC was computed between the neural
423 network structure and the traditional homogeneous reconstruction in cryoSPARC after applying a
424 real space mask and with phase randomization at frequencies above 3.1 \AA , the $\text{GSFSC}_{0.143}$ of the
425 cryoSPARC reconstruction. The real space mask was defined by first thresholding the volume at
426 half of the 99.99th percentile density value. The mask was then dilated by 15 \AA from the original
427 boundary, and a soft cosine edge was used to taper the mask to 0 at 25 \AA from the original
428 boundary.

429 **CryoDRGN heterogeneous reconstruction**

430 CryoDRGN encoder-decoder networks were trained from their randomly initialized values
431 for each single particle cryo-EM dataset. Unless otherwise specified, all networks were trained on
432 minibatches of 8 images using the Adam optimizer with a learning rate of 0.0001. After training,
433 the dataset was evaluated through the encoder, and the *maximum a posteriori* value of $q(\mathbf{z}|X)$ was
434 defined as the latent encoding for each image. Visualization of the latent encodings with PCA and
435 UMAP and analysis with k-means clustering was performed with scikit-learn¹⁴. Density maps were
436 generated by evaluating the decoder on a desired value of the latent variable z and the 3D
437 coordinates of a $D \times D \times D$ voxel array spanning $[-0.5, 0.5]^3$.

438 Heterogeneous reconstruction of simulated datasets

439 For each simulated heterogeneous dataset, a 1D latent variable model was trained for 100
440 epochs. The encoder architecture was 256x3 (nodes/layer x layers) and the decoder architecture
441 was 512x5. The image poses used for training were the ground truth image poses. Structures shown
442 in Figure 3b were generated at the 5th, 23rd, 41st, 59th, 77th, and 95th percentile values of the
443 latent encodings, and sharpened by a B-factor of -100 \AA^2 . Structures shown in Figure 3e were
444 generated at the k -means cluster centers after performing k -means clustering with $k=3$ on the latent
445 encodings, and sharpened by a B-factor of -100 \AA^2 .

446 Heterogeneous reconstruction of the 80S ribosome (EMPIAR-10028)

447 *Pilot experiments:* A 10D latent variable model was trained on downsampled images
448 ($D=96$, $4.91 \text{ \AA}/\text{pix}$) from EMPIAR-10028 for 50 epochs. The encoder and decoder architectures
449 were 128×10 , and the mini-batch size was 5. Image poses were obtained from a traditional
450 homogeneous reconstruction in cryoSPARC.

451 *Particle filtering:* After training, k -means clustering with $k=20$ was performed on the
452 predicted latent encodings for the dataset. One cluster contained 860 particles that were outliers
453 when viewing the projected encodings along the first and second principle component. This

454 observation was reproducible, and the particles belonging to the outlier cluster from either of two
455 replicates (960 particles in total) were removed from the dataset.

456 *High resolution training:* After particle filtering, a 10D latent variable model was trained
457 on the remaining 104,280 images (D=256, 1.84 Å/pix) for 150 epochs. The encoder and decoder
458 architectures were 1024x3.

459 *Analysis:* After training, k -means clustering with $k=20$ was performed on the predicted
460 latent encodings for the dataset, and volumes were generated at the cluster centers using the
461 decoder network. Representative structures were manually selected for visualization in Figure 4.

462 *Heterogeneous reconstruction of the L17-depleted ribosome assembly intermediates (EMPIAR-*
463 *10076)*

464 *Pilot experiments:* A 10D latent variable model was trained on downsampled images
465 (D=128, 3.3 Å/pix) from EMPIAR 10076 for 50 epochs. The encoder and decoder architectures
466 were 256x3. Image poses were obtained from a traditional homogeneous reconstruction in
467 cryoSPARC.

468 *Particle filtering:* Particles with $\mathbf{z} \leq -1$ were removed from subsequent analysis.

469 *High resolution training:* A 10D latent variable model was trained on the remaining
470 101,604 images (D=256, 1.7 Å/pix) for 50 epochs. The encoder and decoder architectures were
471 1024x3.

472 *Analysis:* After training, the dataset's latent encodings were viewed in 2D with UMAP¹⁵.
473 Density maps corresponding to the major and minor assembly states were generated at the mean
474 latent encoding for each class, *i.e.* $\hat{\mathbf{z}}_M = \frac{1}{|M|} \sum_{i \in M} \mathbf{z}_i$, where M is the set of particles assigned to a
475 given class in the published 3D classification. Instead of evaluating the volume decoder at $\hat{\mathbf{z}}_M$, we
476 find the latent encoding of the dataset closest in Euclidean distance to $\hat{\mathbf{z}}_M$ as the “on data”
477 representative encoding.

478 *New assembly state:* Particles corresponding to the new assembly state (C4) were manually
479 selected from the UMAP embeddings with an interactive lasso tool in a custom visualization script.
480 The mean latent encoding of the resulting 1,211 selected particles was used to generate the
481 structure representative for this new assembly state.

482 *Voxel-based back-projection:* The particles associated with class C4 and their
483 corresponding poses were used to reconstruct a structure via traditional voxel-based back-
484 projection using a custom Python script. In this simplified implementation, images were first phase
485 flipped to correct for the CTF. Then each image was centered by its in-plane translation and aligned
486 in 3D space based on its 3D rotation. The density for each voxel was computed using a linear
487 interpolation kernel. The structure was then low-pass filtered to 8 Å for visual clarity.

488 *Heterogeneous reconstruction of the pre-catalytic spliceosome (EMPIAR-10180)*

489 *Pilot experiments:* A 10D latent variable model was trained on downsampled images
490 (D=128, 4.25 Å/pix) from EMPIAR 10180 for 50 epochs. The encoder and decoder architectures
491 were 256x3. Poses were obtained from the consensus reconstruction values given in the
492 consensus_data.star deposited to EMPIAR 10180.

493 *Particle filtering:* The UMAP embeddings showed multiple clusters where the largest
494 cluster corresponded to fully formed pre-catalytic spliceosomes. Particles corresponding to other
495 clusters were removed from subsequent analyses by first performing *k*-means clustering with *k*=20
496 on the latent encodings, and removing *k*-means clusters whose structure did not resemble the fully
497 formed pre-catalytic spliceosome (11 out of 20 *k*-means clusters in one replicate, and 10 out of 20
498 in a second replicate).

499 *High resolution training:* A 10D latent variable model was trained on the remaining
500 150,098 images (D=256, 2.1 Å/pix) for 50 epochs. The encoder and decoder architectures were
501 1024x3.

502 *Analysis:* After training, the dataset’s latent encoding was viewed in 2D with UMAP (Fig.
503 6a) and PCA (Fig. 6c). Density maps in Figure 6d were generated at the latent encoding values
504 that traverse PC1 at five equally spaced points between the 5th and 95th percentile of PC1 values.
505 Density maps in Extended Fig. 7 were generated at the latent encoding values that traverse PC2 at
506 five equally spaced points between the 5th and 95th percentile of PC2 values.

507 **Latent space graph traversal for generating trajectories**

508 Trajectories were generated by first creating a nearest-neighbors graph from the latent
509 encodings of the images, where a neighbor was defined if the Euclidean distance was below a
510 threshold computed from the statistics of all pairwise distances. We choose a value such that the
511 average number of neighbors across all nodes is 5. Edges were then pruned such that a given node
512 does not have more than 10 neighbors. Then, Dijkstra’s algorithm was used to find the shortest
513 path along the graph connecting a series of anchor points, and density maps were generated at the
514 \mathbf{z} value of the visited nodes. Anchor points were set to be the “on-data” cluster centers after
515 performing k -means clustering of the latent encodings with $k=20$. Instead of using the mean value
516 of each k -means cluster, we define the latent encoding closest in Euclidean distance to the k -means
517 cluster center as the “on-data” cluster center.

518 To generate Supplemental Movie 1 of the 80S ribosome, 113 density maps were generated
519 by following the protocol above, and we visualized a representative sequence of 60 density maps
520 that contained the 40S rotated state. To generate Supplemental Movies 2 and 3 of the assembling
521 bacterial LSU, anchor points were manually chosen from an interactive tool provided in
522 cryoDRGN to create a path along the C-class assembly pathway and the D-class assembly
523 pathway. To generate Supplemental Movie 4 of the pre-catalytic spliceosome, 132 density maps
524 were generated following the above protocol.

525 **2D class averages**

526 2D classification was performed in cryoSPARC¹⁰ using all default options except for the
527 number of 2D classes, which was set to 20.

528 **Data availability**

529 Trained cryoDRGN models for all experiments, simulated datasets, and indices of filtered
530 particles of EMPIAR-10028, EMPIAR-10076, and EMPIAR-10180 are available upon request.

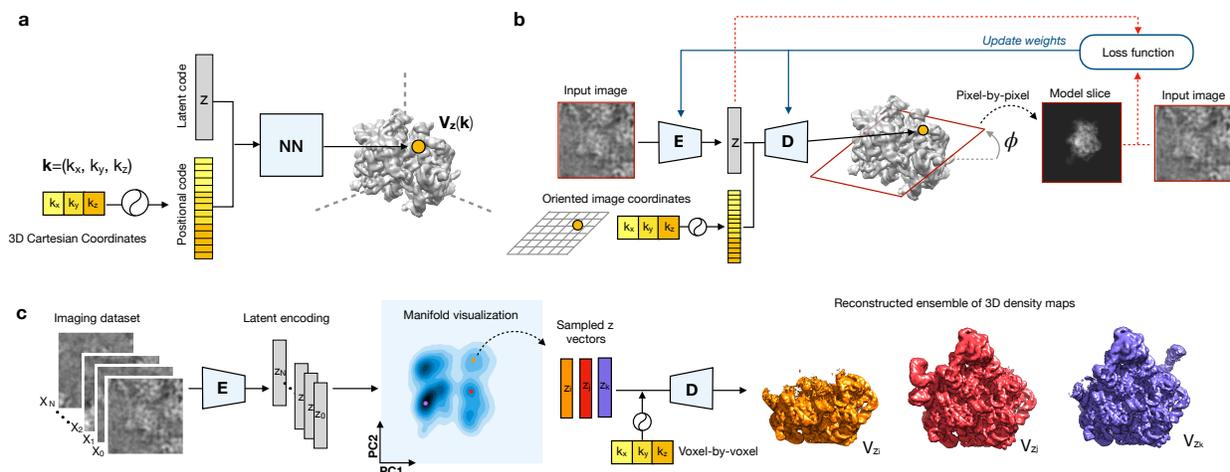
531 **Software availability**

532 All software and analysis scripts are implemented in custom Python code using PyTorch¹⁶
533 and are available at cryodrgn.csail.mit.edu.

References

1. Bricman, P. A. & Ionescu, R. T. CocomNet: A deep neural network for mapping pixel coordinates to color values. *arXiv.org cs.CV*, (2018).
2. Bepler, T., Zhong, E., Kelley, K., Brignole, E. & Berger, B. Explicitly disentangling image content from translation and rotation with spatial-VAE. 15435–15445 (2019).
3. Zhong, E. D., Bepler, T., Davis, J. H. & Berger, B. Reconstructing continuous distributions of 3D protein structure from cryo-EM images. *arXiv.org q-bio.QM*, (2019).
4. Kingma, D. P. & Welling, M. Auto-Encoding Variational Bayes. *arXiv.org stat.ML*, (2013).
5. Rezende, D. J., Mohamed, S. & Wierstra, D. Stochastic Backpropagation and Approximate Inference in Deep Generative Models. *arxiv.org* (2014).
6. Bracewell, R. N. Strip Integration in Radio Astronomy. *Aust. J. Phys.* **9**, 198–217 (1956).
7. The PyMOL Molecular Graphics System, Version 2.3 Schrödinger, LLC.

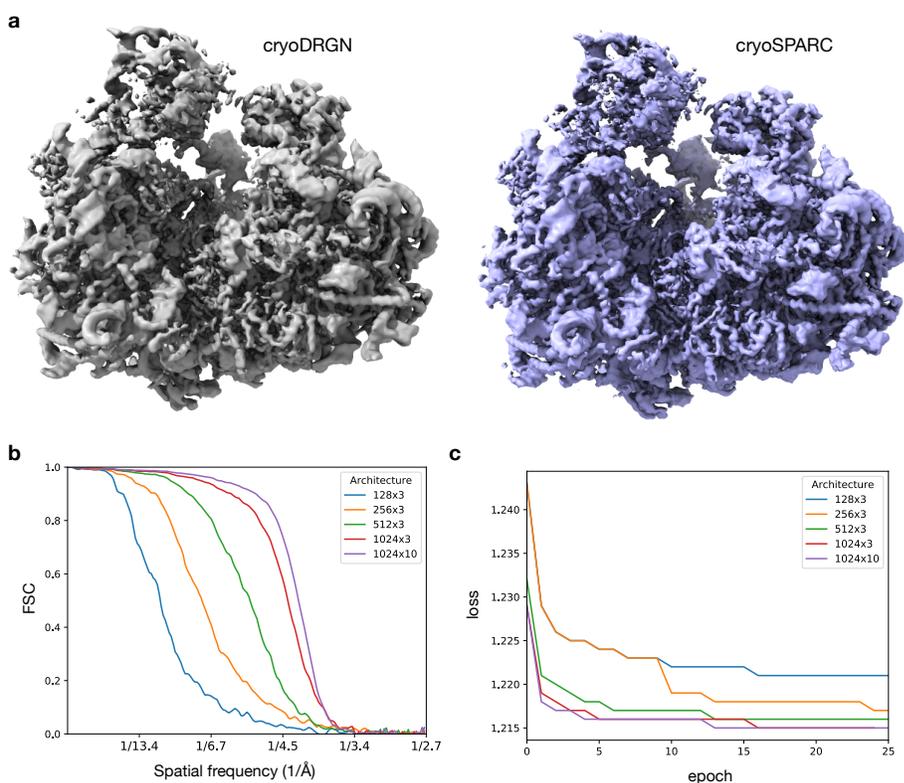
8. Pettersen, E. F. *et al.* UCSF Chimera—A visualization system for exploratory research and analysis. *Journal of Computational Chemistry* **25**, 1605–1612 (2004).
9. Iudin, A., Korir, P. K., Salavert-Torres, J., Kleywegt, G. J. & Patwardhan, A. EMPIAR: a public archive for raw electron microscopy image data. *Nat. Methods* **13**, 387–388 (2016).
10. Punjani, A., Rubinstein, J. L., Fleet, D. J. & Brubaker, M. A. cryoSPARC: algorithms for rapid unsupervised cryo-EM structure determination. *Nat. Methods* **14**, 290–296 (2017).
11. Rosenthal, P. B. & Henderson, R. Optimal determination of particle orientation, absolute hand, and contrast loss in single-particle electron cryomicroscopy. *J. Mol. Biol.* **333**, 721–745 (2003).
12. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arxiv.org* (2014).
13. Wong, W. *et al.* Cryo-EM structure of the Plasmodium falciparum 80S ribosome bound to the anti-protozoan drug emetine. *Elife* **3**, e01963 (2014).
14. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011).
15. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arxiv.org* (2018).
16. Paszke, A. *et al.* PyTorch: An Imperative Style, High-Performance Deep Learning Library. 8026–8037 (2019).



534

535 **Figure 1. The cryoDRGN method for heterogeneous single particle cryo-EM reconstruction.**

536 **a)** A *deep coordinate network* approximates a molecule's density as a function of featurized 3D
537 Cartesian coordinates and continuous latent variables, z , which define a continuous manifold of
538 heterogeneous structures. **b)** The overall cryoDRGN training framework consists of two neural
539 networks structured in an encoder/decoder architecture. Data is represented in the Fourier domain
540 in order to relate 2D images as slices out of the 3D density map. During training, an input image
541 is encoded in latent space by the encoder network (E). A 2D lattice is rotated by the image's
542 previously determined pose, ϕ , to represent the 3D coordinates of the image slice. Given the
543 coordinates and a sample of the predicted latent variable z , the image is reconstructed pixel-by-
544 pixel through the decoder (D), *i.e.* the deep coordinate network. The loss function is a variational
545 upper bound on the data likelihood and consists of the image reconstruction error and latent loss
546 (red arrows), which is used to update neural network weights by stochastic gradient descent (blue
547 arrows). **c)** After training, the encoder can be used to visualize the dataset's distribution in latent
548 space (manifold visualization), and the decoder can be used to directly reconstruct structures at
549 arbitrary points from the latent code. Example micrographs and reconstructed density maps from
550 EMPIAR 10076³².



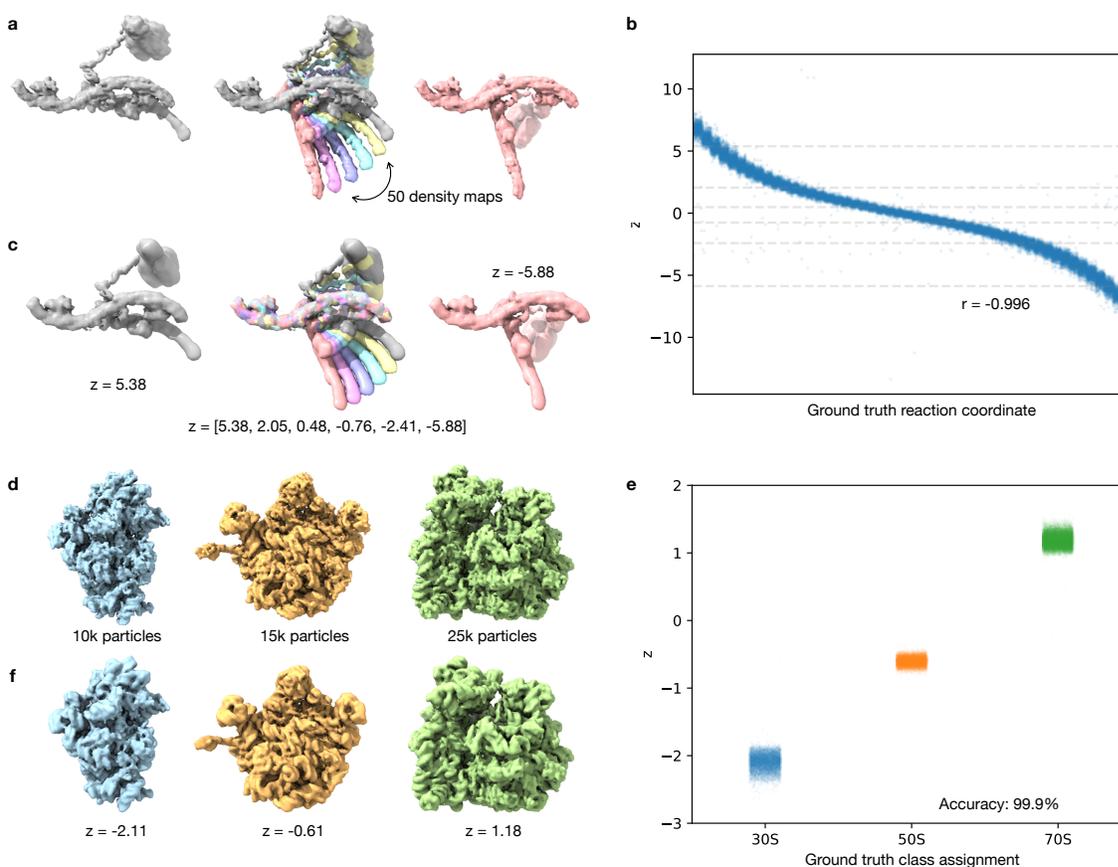
551

552 **Figure 2. Deep coordinate network representation of static structure. a)** Reconstructed density
553 map produced by a deep coordinate network with 10 hidden layers of dimension 1024 trained on
554 particle images from EMPIAR 10028²⁹ (D=360, Nyquist limit of 2.7 Å) and a traditional
555 homogeneous reconstruction in cryoSPARC²⁹. **b)** Fourier shell correlation (FSC) curves between
556 the density map produced by deep coordinate networks of varying architectures (nodes x hidden
557 layers) and the traditional homogeneous reconstruction in (a) after 25 epochs of training. **c)**
558 Average loss over the dataset during training deep coordinate networks of varying architectures
559 on EMPIAR 10028²⁹.

560

561

562



563

564 **Figure 3. CryoDRGN heterogeneous reconstruction of simulated datasets with continuous**

565 **and discrete heterogeneity. a)** Ground truth density maps sampled along a reaction coordinate

566 that describes the transition from the leftmost to rightmost structure used to simulate a dataset with

567 continuous heterogeneity. **b)** Predicted latent encoding for each image of the dataset from (a) after

568 training a cryoDRGN 1D latent variable model versus the ground truth reaction coordinate

569 describing the motion (Spearman $r = -0.996$). **c)** Reconstructed structures at specified values of

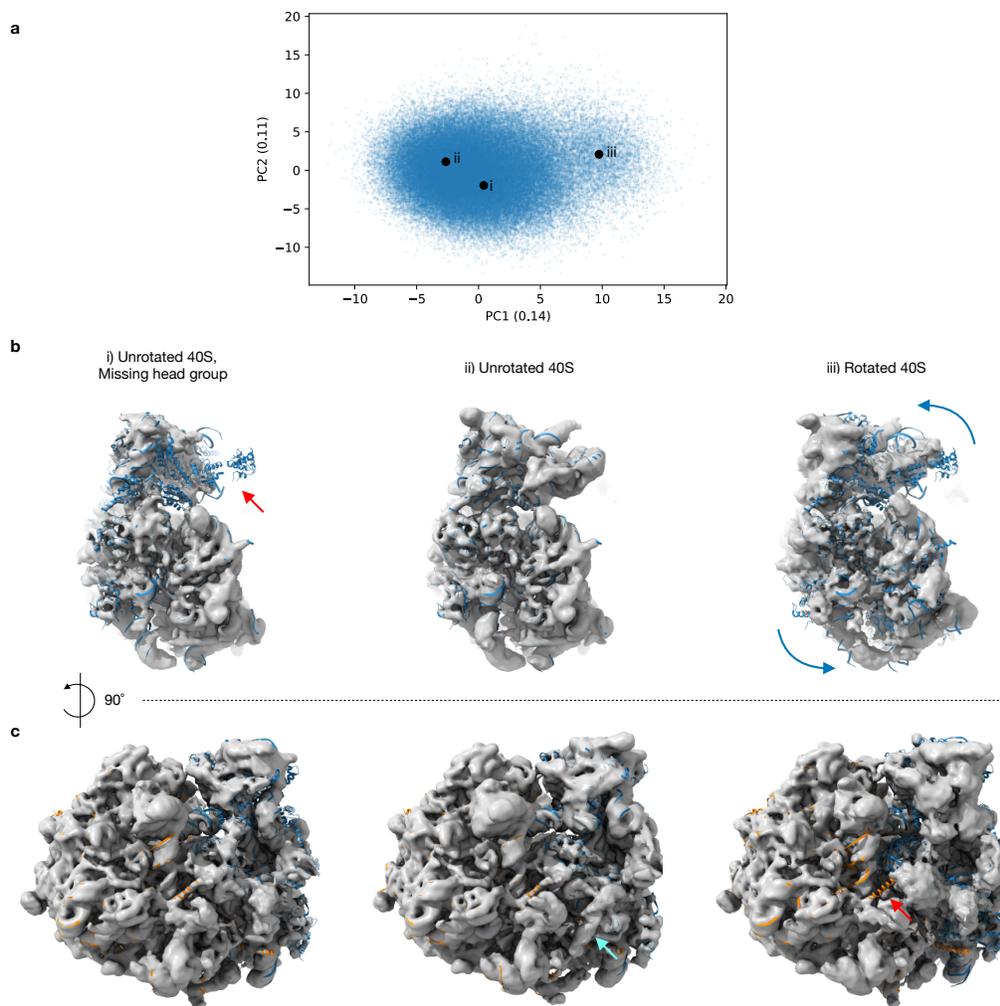
570 the latent variable, shown as dotted lines in (b). **d)** Ground truth density maps of the bacterial 30S,

571 50S, and 70S ribosome used to simulate a dataset with discrete heterogeneity. **e)** Predicted latent

572 encoding for each particle image of (d) variable after training a cryoDRGN 1D latent variable

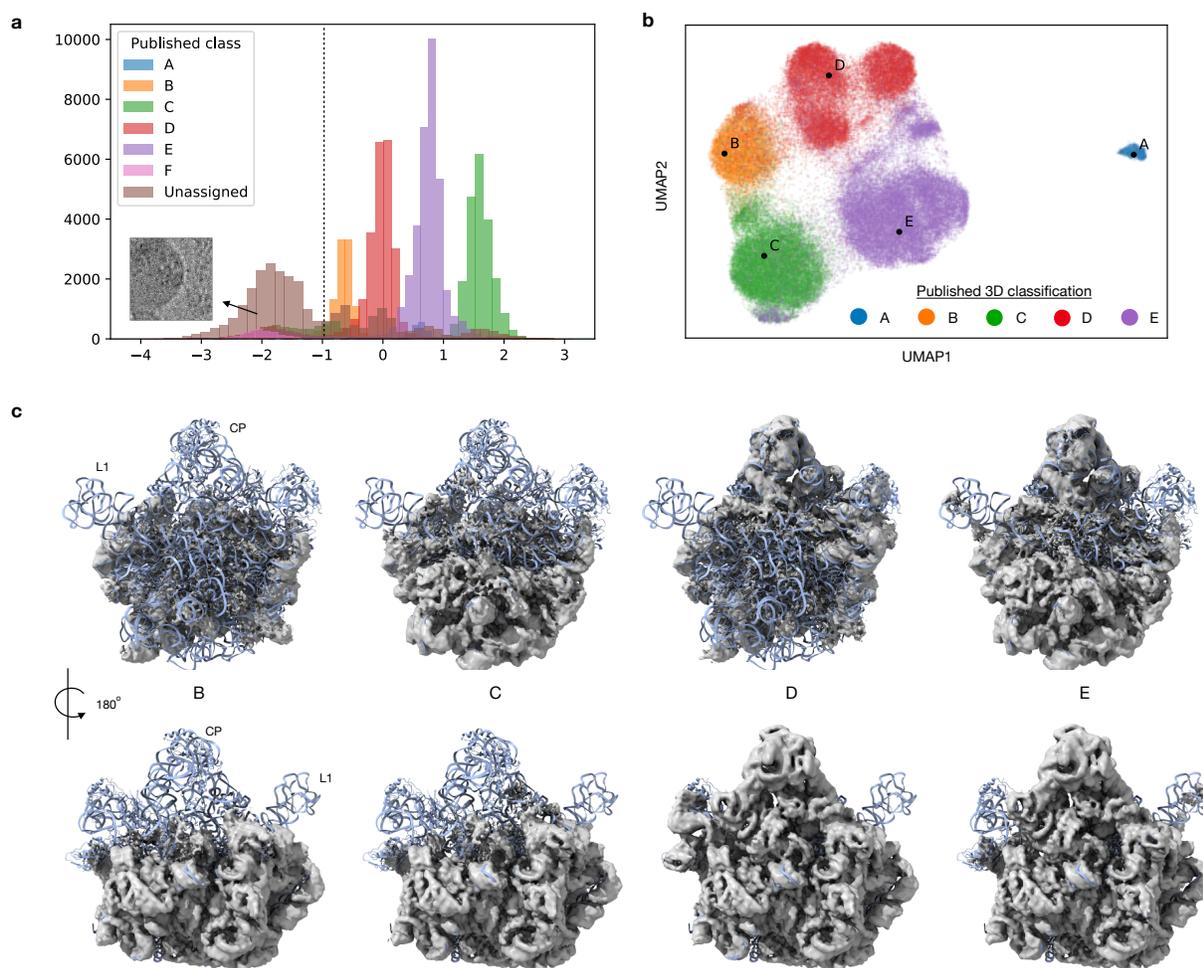
573 model vs. its ground truth class assignment (classification accuracy of 99.9%). **f)** Reconstructed

574 structures at specified values of the latent variable from (e).

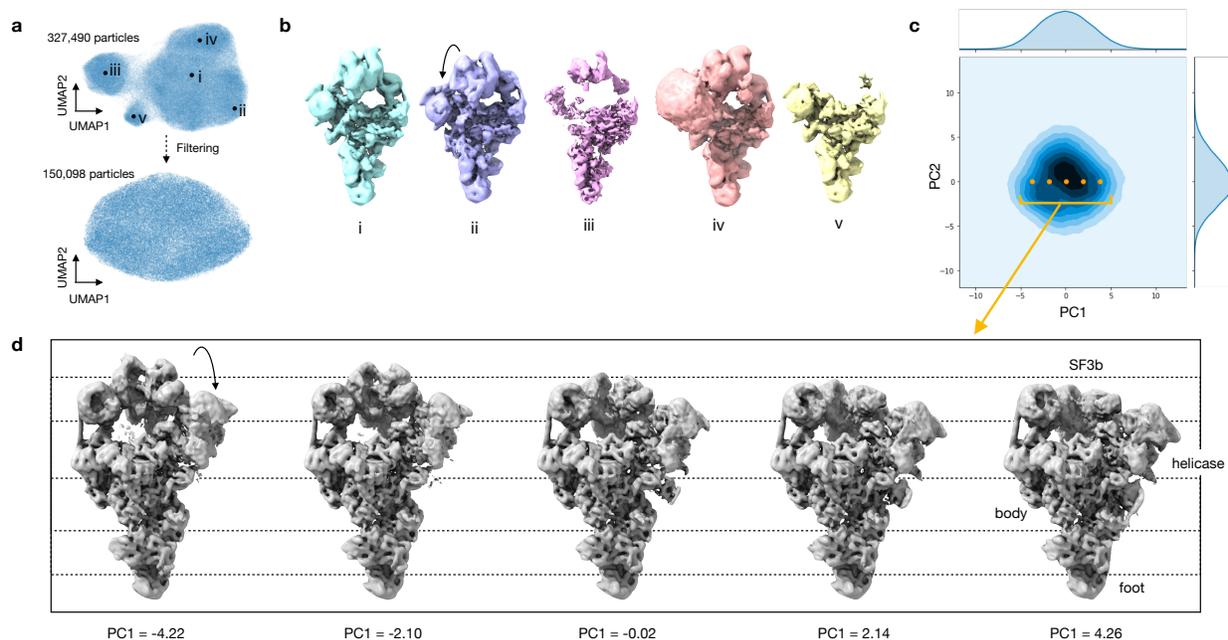


575

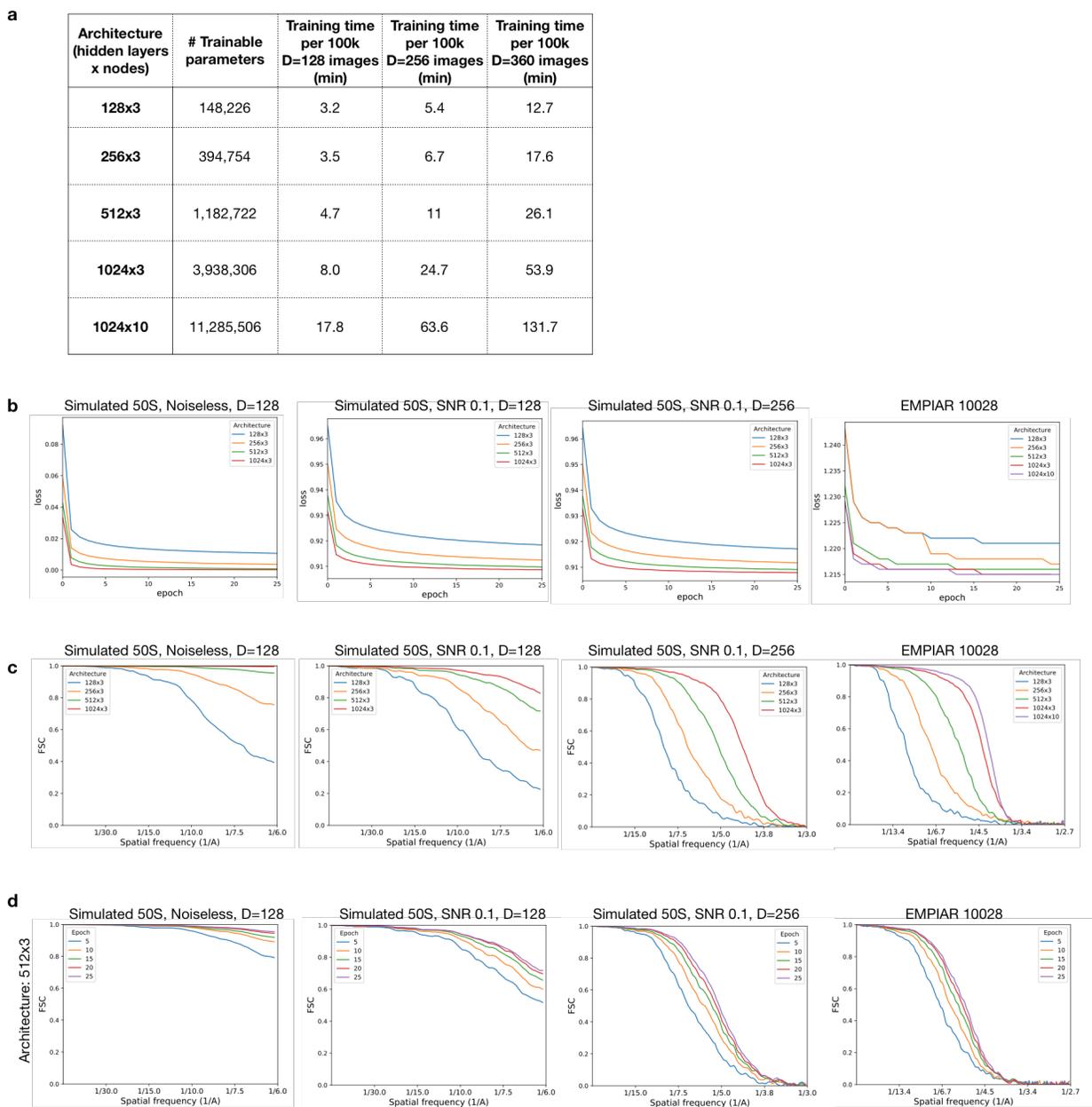
576 **Figure 4. CryoDRGN heterogeneous reconstruction of the *Pf*80S ribosome. a)** PCA projection
577 of latent space encodings after training a 10D latent variable model on particle images from
578 EMPAIR-10028²⁹. **b)** Three representative density maps that were reconstructed at the points
579 depicted in (a) are shown with a docked atomic model (PDB 3J79, 3J7A) of the 40S (blue). The
580 red arrow highlights the missing 40S head group, and the blue arrow depicts the rotation of the
581 40S relative to the 60S. **c)** Additional views of the structures shown in (b), with atomic model of
582 the 60S colored in orange. The cyan arrow notes the presence of an additional RNA helix not
583 present in the homogeneous reconstruction, and the red arrow notes the disappearance of the C-
584 terminus of eL8 in the rotated state.



585
 586 **Figure 5. CryoDRGN heterogeneous reconstruction of the assembling large ribosomal**
 587 **subunit from *E. coli*.** **a)** Histograms of latent encodings of particle images from EMPIAR 10076³²
 588 after training a cryoDRGN 1D latent variable model. Overlaid histograms are shown for particles
 589 from each published major class assignments from *Davis et al*³¹. A cutoff of $z = -1$ was used to
 590 filter impurities from the dataset for subsequent analyses. Example image of an ice artifact
 591 predicted at $z = -2$. **b)** UMAP visualization of latent encodings after training a cryoDRGN 10D
 592 latent variable model, colored by the published major class assignments³². **c)** CryoDRGN
 593 reconstructed density maps of the major assembly states of the LSU generated from points B-E
 594 shown in (b) along with a docked atomic model (PDB 4YBB).



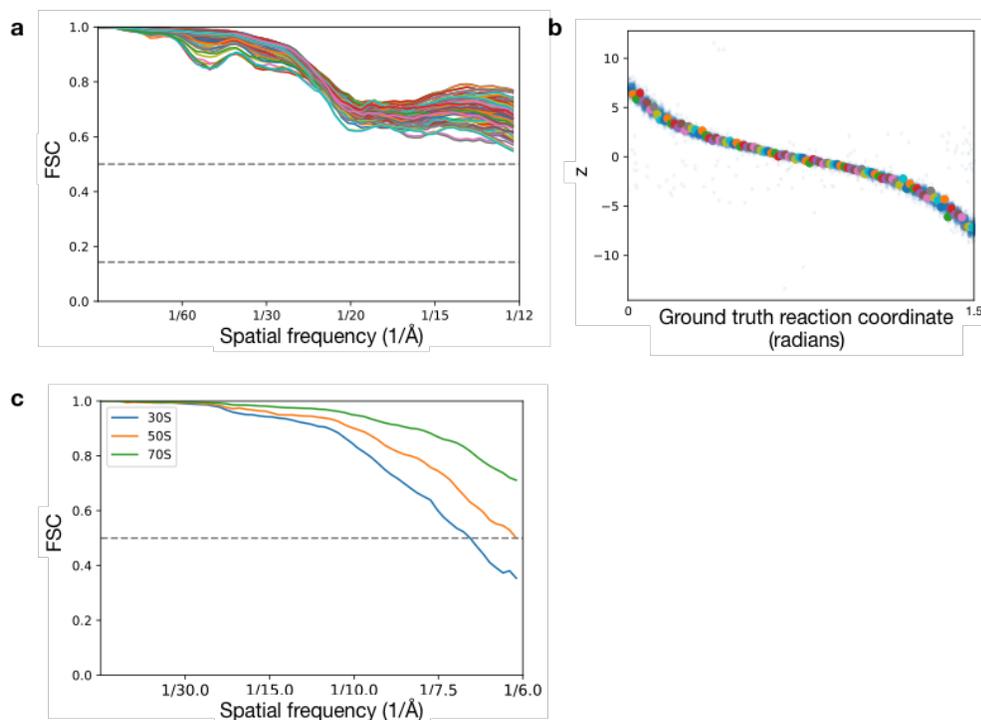
595
596 **Figure 6. CryoDRGN heterogeneous reconstruction of the pre-catalytic spliceosome. a)**
597 UMAP visualization of latent encodings after training a 10D latent variable model with cryoDRGN
598 on EMPIAR 10180³³, before (top) and after (bottom) particle filtering. **b)** Representative structures
599 generated at points shown in (a) which depict the expected structures (i,ii), broken particles (iii),
600 particles with apparent aggregation (iv), and the complex lacking the SF3b domain (v). **c)** PCA
601 projection of latent space encodings after training a 10D latent variable model on the filtered
602 images. **d)** Structures generated by traversing along PC1 of the latent space encodings at points
603 shown in (c).



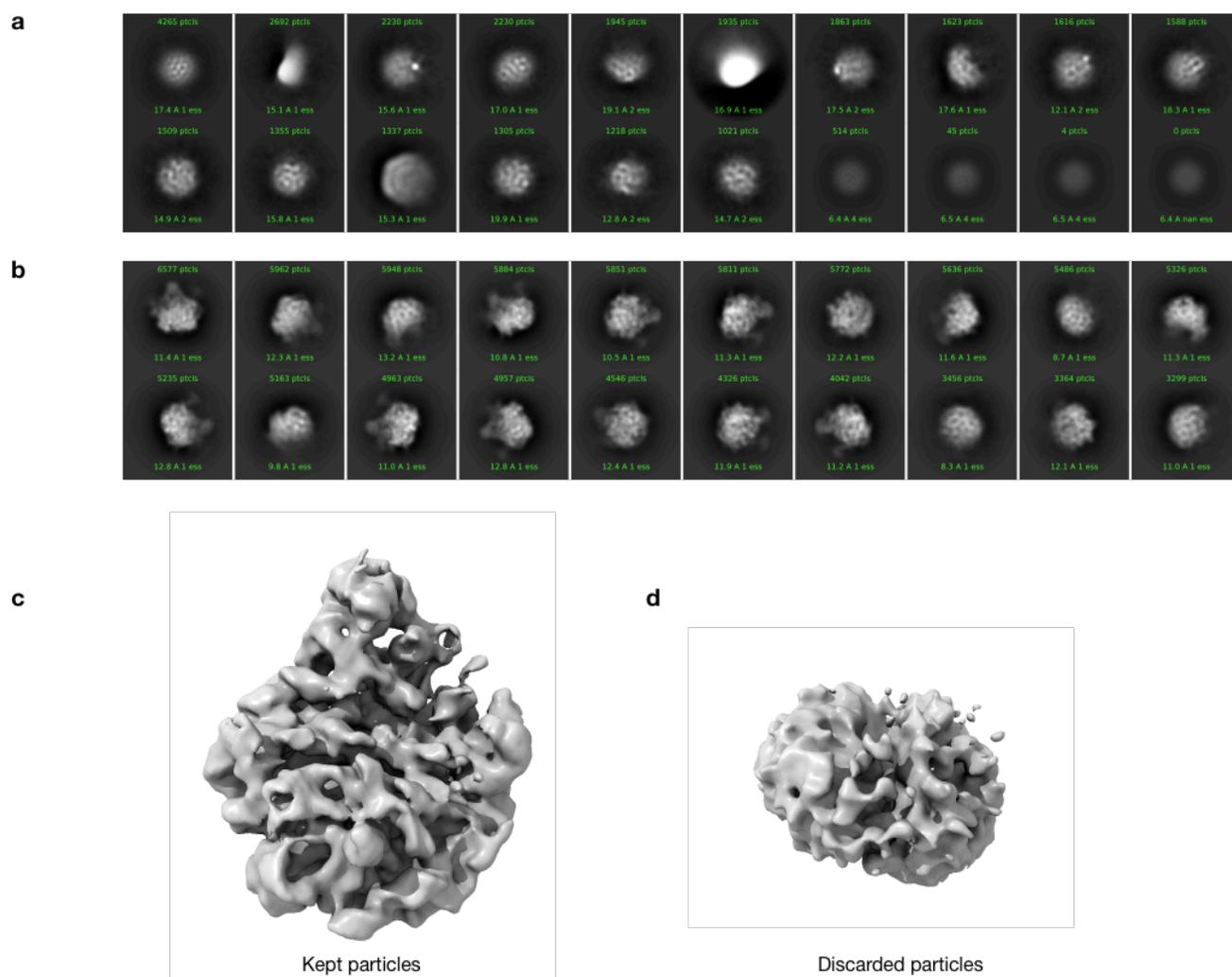
Supplementary Figure 1. Neural network training statistics for homogeneous reconstruction.

a) Training time in minutes per 100k images for different architectures and image sizes on a single Nvidia V100 GPU. **b)** Loss curve for training deep coordinate networks of varying architectures on four different datasets: 50k simulated noiseless projection images of the 50S ribosome (D=128, Nyquist limit of 6 Å), 50k simulated micrographs of the 50S ribosome (D=128), 50k simulated micrographs of the 50S ribosome (D=256, Nyquist limit of 3 Å), and 104,249 micrographs of the

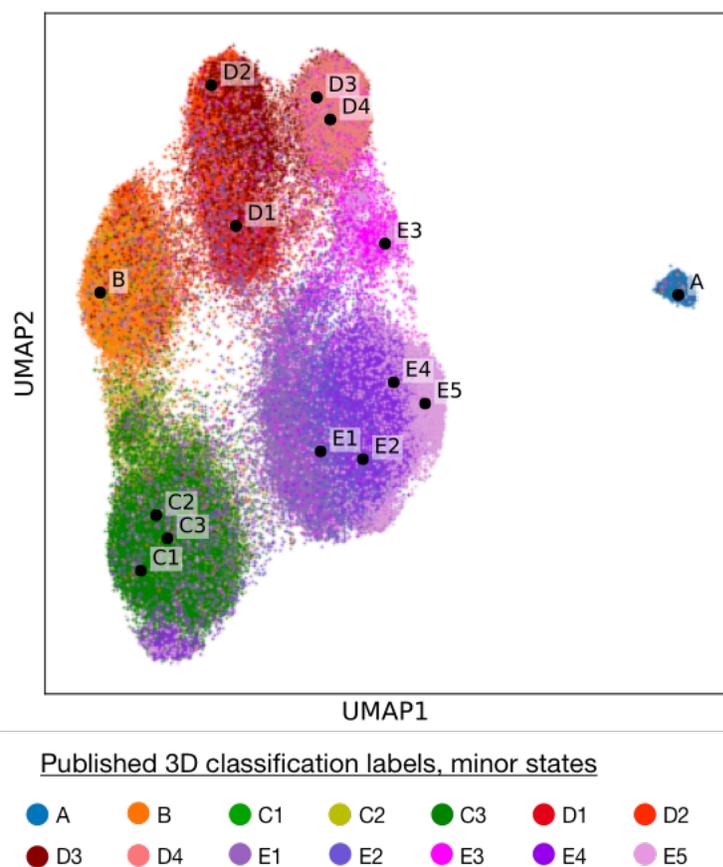
80S ribosome from EMPIAR 10028 ($D=360$, Nyquist limit of 2.68 \AA). **c)** FSC curve between the ground truth density map and the learned density map after 25 epochs of training deep coordinate networks of varying architectures. **d)** FSC curve between the ground truth density map and the learned density map at different epochs of training a deep coordinate network with 3 hidden layers of dimension 512. We use the traditionally reconstructed map as the ground truth structure for EMPIAR 10028.



Supplementary Figure 2. FSC curves between ground truth maps and density maps from cryoDRGN trained on heterogeneous simulated datasets. a) 100 FSC curves between generated and ground truth density maps. The density maps are generated at the value of the latent variable predicted for a given image, and compared against the ground truth density map that generated the image. Images are uniformly sampled along the reaction coordinate. **b)** The predicted latent encoding for the 100 images along the ground truth reaction coordinate for the density maps in (a). **c)** FSC curve between the generated density maps shown in Figure 3f and the ground truth 30S, 50S, and 70S density map.

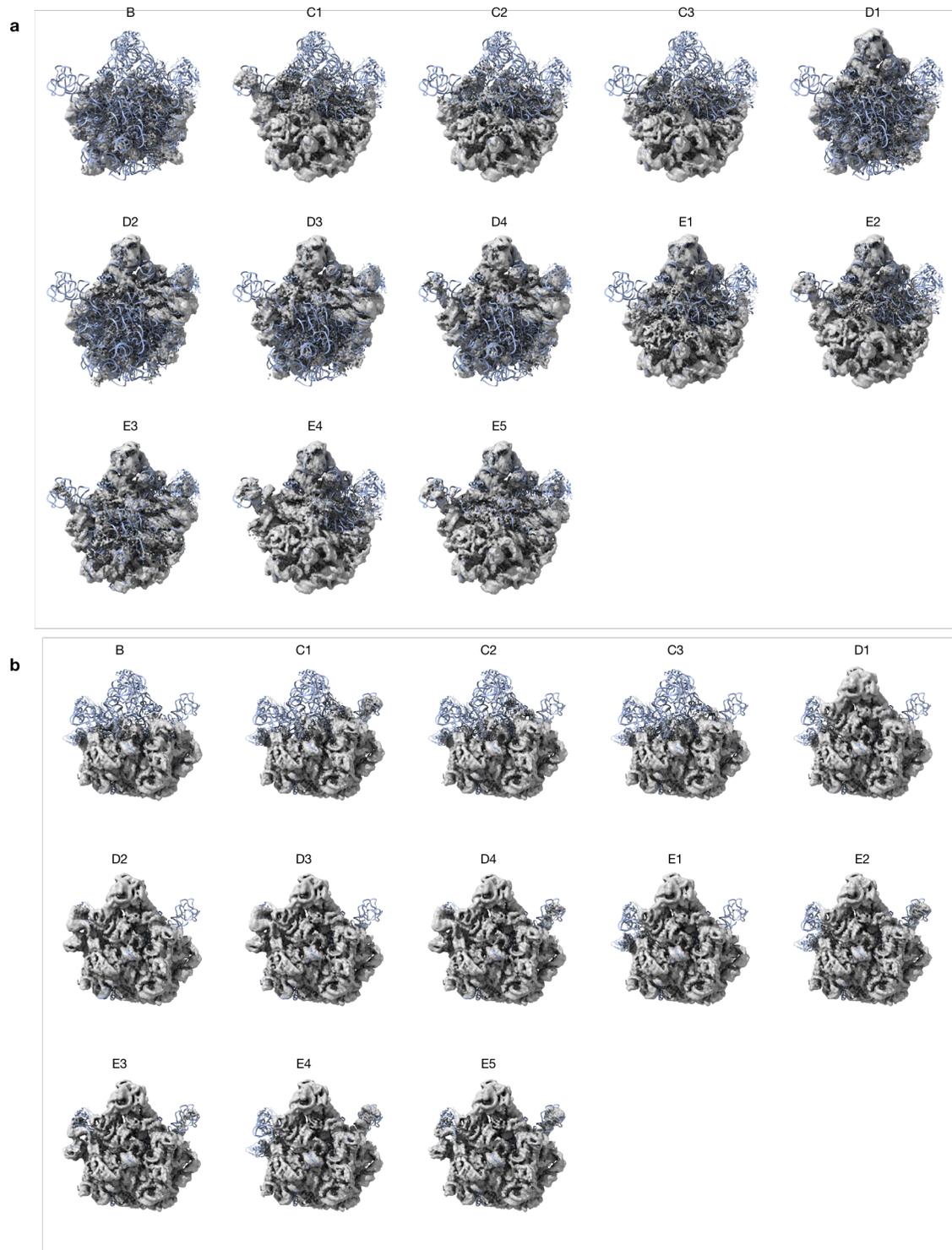


Supplementary Figure 3. Filtering of particles from the assembling ribosome dataset. a) 2D class averages of discarded particles with $z \leq -1$ from Figure 5a. **b)** 2D class averages of kept particles with $z > 1$ from Figure 5a. **c)** CryoSPARC *ab-initio* reconstruction of kept particles ($z > -1$) and **d)** of discarded particles ($z \leq -1$) from Figure 5a.



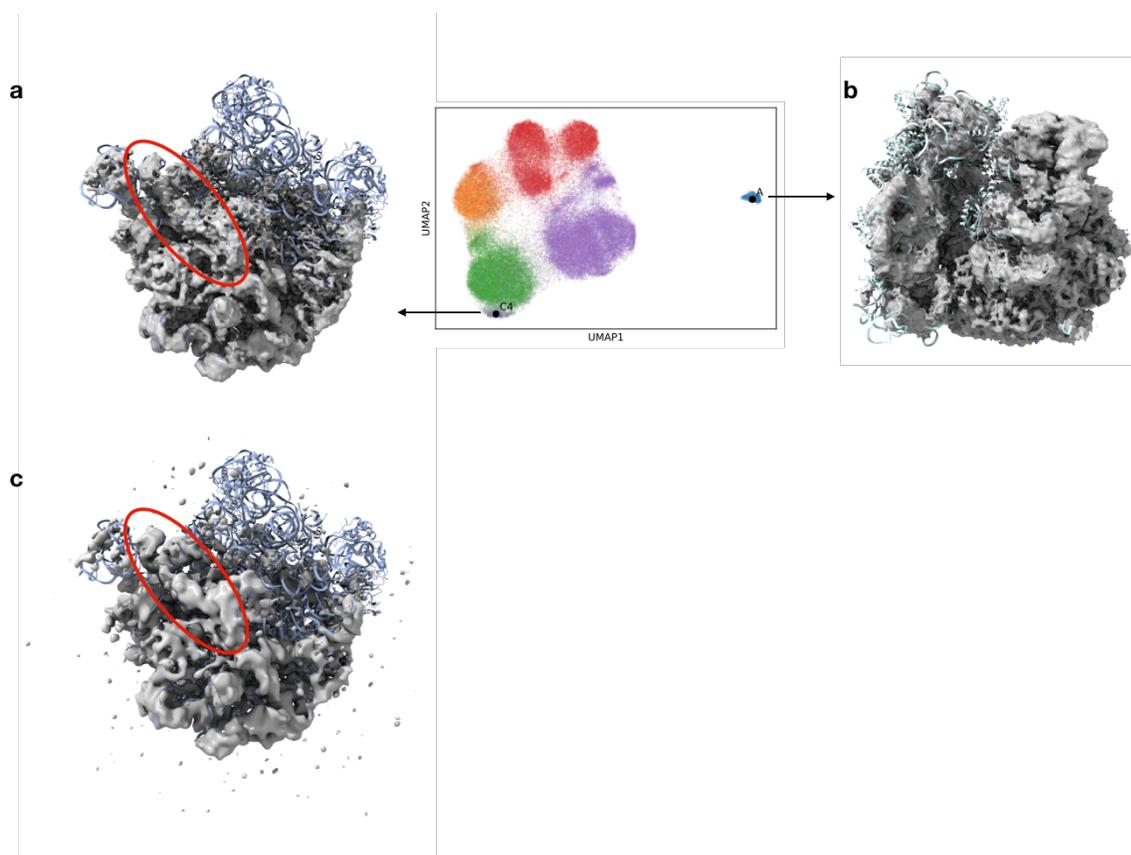
Supplementary Figure 4. CryoDRGN latent encodings trained on the assembling ribosome.

UMAP embedding of the latent space encodings of particle images after training a cryoDRGN 10D latent variable model on EMPIAR 10076. Points are colored by the 3D classification labels corresponding to the minor states of LSU assembly from *Davis et al.*

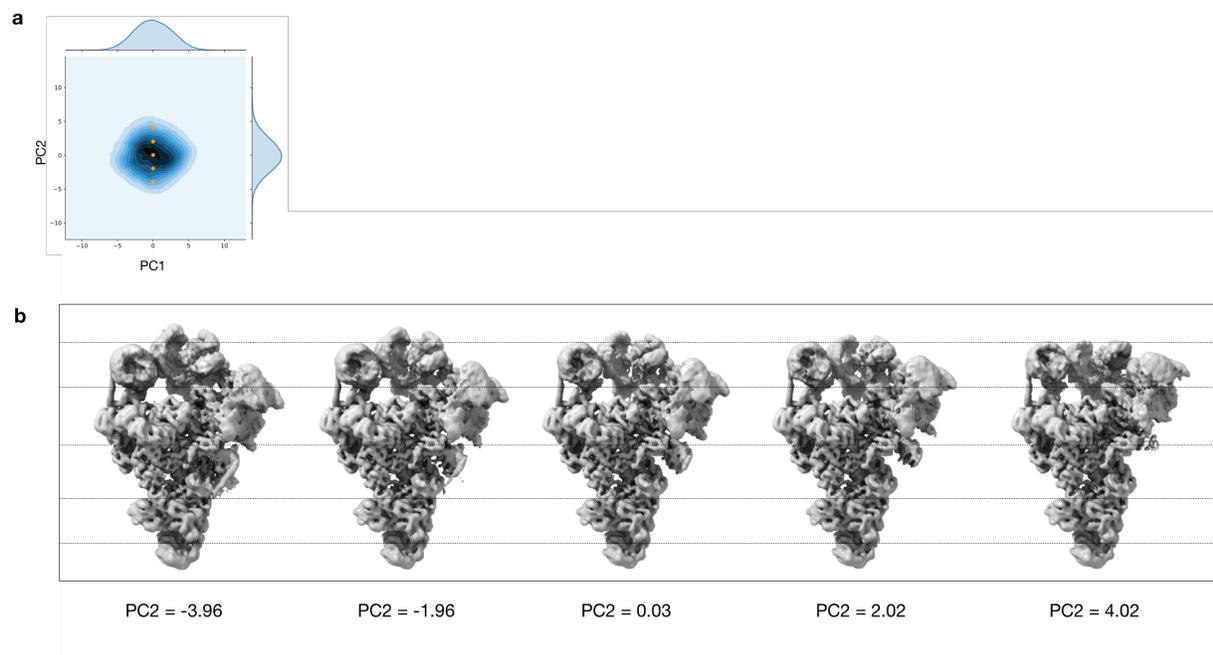


Supplementary Figure 5. Minor LSU assembly states reconstructed from cryoDRGN trained on the assembling ribosome dataset. a) Front view and b) back view of minor state density maps after training a cryoDRGN 10D latent variable model on particle images from EMPIAR 10076

with the 50S crystal structure docked (PDB 4YBB). Each cryoDRGN structure is generated at the latent variable values shown in Supplementary Figure 4, which are computed from the mean latent encoding of particles with the corresponding class assignment from *Davis et al.* Views match perspectives from Figure 5d.



Supplementary Figure 6. Additional structures reconstructed from cryoDRGN trained on the assembling ribosome dataset. a) Density map of a new assembly state, class C4, produced by cryoDRGN. Helix 68 (red oval) was exclusively associated with mature classes E4 and E5 in Davis *et al.* The structure is generated from point C4 in latent space, which belongs to a small cluster proximal to class C that was classified into class E4 and E5 by Davis *et al.* **b)** The density map of the 70S ribosome reconstructed by cryoDRGN from point A in latent space. **c)** Voxel-array backprojection of the 1,211 particles contained in the latent space cluster corresponding to the new assembly state with atomic model docked and helix 68 highlighted (red oval).



Supplementary Figure 7. Additional structures of the pre-catalytic spliceosome reconstructed by cryoDRGN. a) PCA projections of the 10D latent encodings from cryoDRGN with 5 points along PC2 shown in orange. **b)** Density maps produced by cryoDRGN at the 5 highlighted points from (a).

Supplemental Movie 1. Trajectory of the pre-catalytic spliceosome. Traversal in latent space (left) and corresponding structures generated from cryoDRGN (right).

Supplemental Movie 2. Trajectory of the assembling ribosome along the D-class assembly pathway described in Davis *et al.* Traversal in latent space (left) and corresponding structures generated from cryoDRGN (right).

Supplemental Movie 3. Trajectory of the assembling ribosome along the C-class assembly pathway described in Davis *et al.* Traversal in latent space (left) and corresponding structures generated from cryoDRGN (right).

Supplemental Movie 4. Trajectory of the 80S ribosome. Traversal in latent space (left) and corresponding structures generated from cryoDRGN (right).