# Separating distinct structures of multiple macromolecular assemblies from cryo-EM projections

Eric J. Verbeke[1-3], Yi Zhou[1-3], Andrew P. Horton[1-3], Anna L. Mallam[1-3], David W. Taylor[1–4*],

Edward M. Marcotte[1–3*]


**Affiliations:**

[1]Department of Molecular Biosciences

[2]Center for Systems and Synthetic Biology

[3]Institute for Cellular and Molecular Biology, University of Texas at Austin, Austin, TX 78712,

USA

[4]LIVESTRONG Cancer Institutes, Dell Medical School, Austin, TX 78712, USA

*Correspondence: marcotte@icmb.utexas.edu (E.M.M), dtaylor@utexas.edu (D.W.T),

1

1   **Abstract**

2   Cryo-electron microscopy is traditionally applied to samples purified to near homogeneity as

3   current reconstruction algorithms are unable to handle heterogeneous mixtures of structures from

4   many macromolecular complexes. We extend on long established methods and demonstrate that

5   relating two-dimensional projection images by their common lines in a graphical framework is

6   sufficient for partitioning distinct protein and multiprotein complexes within the same data set.

7   Using this approach, we first group a large set of synthetic reprojections from 35 unique

8   macromolecular structures ranging from ~30 – 3000 kDa into individual homogenous classes. We

9   then apply our algorithm on cryo-EM data collected from a mixture of five protein complexes and

10  use existing reconstruction methods to solve multiple three-dimensional structures *ab initio*.

11  Incorporating methods to sort cryo-EM data from heterogeneous mixtures will alleviate the need

12  for stringent purification and pave the way toward investigation of samples containing many

13  unique structures.

14 **Introduction**

15 Cryo-electron microscopy (cryo-EM) has undergone a revolutionary shift in the past few years.

16 Increased signal in electron micrographs, as a result of direct electron detectors, has allowed for

17 the near-atomic and atomic resolution structure determination of many macromolecules of various

18 shapes and sizes (Kühlbrandt, 2014). These new detectors combined with automated data

19 collection software and improvements in image processing suggest that cryo-EM could be utilized

20 as a high-throughput approach to structural biology. One major obstacle remains: sorting through

21 the immense heterogeneity present in a mixture of tens to hundreds to thousands of

22 macromolecular assemblies.

23 We and others have shown that cellular extract can be mined for identification of multiple

24 structures (Kastritis et al., 2017; Verbeke et al., 2018). More recently, we showed that it was

25 possible to reconstruct macromolecular machines from the lysate of a single *C. elegans* embryo

26 (Yi et al., 2018). These studies were limited to the identification of only the most abundant and

27 easily identifiable protein and protein–nucleic acid complexes due to a lack of methods to

28 efficiently categorize which two-dimensional (2D) projection images derive from which three-

29 dimensional (3D) assemblies on the basis of their structural features. While a number of 3D

30 classification schemes exist, all failed to produce reliable reconstructions for the majority of

31 particles in these complicated mixtures. This obstacle emphasizes the long-standing need to sort

32 mixtures of structures in addition to their conformational and compositional heterogeneity.

33 Several methods have been successfully implemented for sorting limited heterogeneity in

34 cryo-EM data. These approaches generally fall into three categories. Currently, the most popular

35 approach for sorting heterogeneity in cryo-EM data utilizes a maximum likelihood estimation to

36 optimize the correct classification of particles into multiple structures (Scheres, 2012; Sigworth,

37 1998; Sigworth et al., 2010). Another approach is to estimate the covariance in cryo-EM data to

38 search for regions of variability between the models and the data (Katsevich et al., 2015; Liao et

3

39  al., 2015; Penczek et al., 2006). The last approach involves computing similarities between

40  projection images in the data before applying clustering methods to separate the data into

41  homogenous subsets (Aizenbud and Shkolnisky, 2016; Herman and Kalinowski, 2008; Shatsky

42  et al., 2010). All of these approaches have been demonstrated on samples containing a primary

43  structure with multiple conformations or variable subunits. However, little work has been done for

44  sorting heterogeneous samples containing multiple distinct structures.

45  Here, we develop a pipeline for building 3D reconstructions from a mixture of distinct

46  particles by first grouping 2D projections into discrete, particle-specific classes using the principles

47  of common lines and a novel graphical clustering framework. We demonstrate our method by

48  partitioning reprojections from 35 previously solved X-ray crystal structures into their correct

49  groups. Furthermore, we applied this pipeline to an experimental set of cryo-EM micrographs

50  containing a mixture of several macromolecular complexes. We were able to reconstruct multiple

51  3D structures after our clustering, improving on classification of all particles simultaneously using

52  current 3D reconstruction software. These results are a necessary first step for moving cryo-EM

53  towards high-throughput structural biology.

54

55  **Results**

56  <u>Classifying projection images from multiple structures</u>

57  A major challenge facing "shotgun"-style cryo-EM is to reconstruct models from projection images

58  arising from multiple distinct structures present in a mixture. To overcome this obstacle, we sought

59  a method to computationally group heterogeneous projection images into discrete classes that

60  each derive from the same structure. Two-dimensional projection images from the same

61  asymmetric object can be related to each other if there is prior information of the three-

62  dimensional object (i.e. an initial model) using projection-matching algorithms. One approach that

63  circumvents the need for a starting model is to relate the 2D projection images based on common

4

64    lines (Van Heel, 1987), derived from the projection-slice theorem, which states that any two 2D

65    projections of the same 3D object must share a 1D line projection in common. In order to partition

66    projection images into homogenous subsets, we developed an algorithm for detecting **S**hared

67    **L**ines **I**n **C**ommon **E**lectron **M**aps (SLICEM). Using this algorithm, we score the similarity of 1D

68    line projections between sets of 2D projection images without knowledge of the underlying 3D

69    objects. Subsequently, these similarity scores can be put into a graphical framework and

70    clustering algorithms can be applied to group related 2D projection images for subsequent 3D

71    reconstructions (Figure 1).

72

73    <u>Synthetic data</u>

74    To test our approach using SLICEM, we generated synthetic reprojections from 35 previously

75    solved X-ray crystal structures (see Methods) (Figure S1). The structures ranged in molecular

76    weight from ~30 – 3000 kDa. Each structure was low-pass filtered to 9 Å and uniformly reprojected

77    to create 12 2D projection images (Ludtke et al., 1999). Next, we combined reprojections from all

78    models to simulate ideal 2D class averages from a heterogeneous cryo-EM dataset. The similarity

79    of 1D line projections from each image was scored using 6 different metrics (see Methods). The

80    precision and recall of correctly pairing 2D projection images from the same 3D structures was

81    computed in order to determine the performance of each metric, and cosine similarity was

82    determined to be the top performing metric (Figure 2A).

83         In order to identify sets of 2D projection images from the same 3D particles, we

84    constructed a network from the comparisons between projection images as follows: Each 2D

85    projection image was represented as a node in a directed graph, with each node connected by

86    edges to the nodes corresponding to the 5 most closely-related 2D projection images based on

87    the similarity of their 1D line projections. While the top-scoring metric in our precision/recall

88    analysis was cosine similarity, the network generated from the Euclidean distance similarity most

89    clearly showed communities (clusters of 2D projections) correctly partitioned by 3D structure

90    (Figure 2B). These results show that partitioning 2D projection images by scoring their common

91    lines is a powerful, unsupervised approach for sorting cryo-EM data from distinct 3D structures

92    within a heterogeneous mixture.

93

94    <u>Cryo-EM on a mixture of protein complexes</u>

95    After validating our SLICEM algorithm on a synthetic dataset, we performed cryo-EM on an

96    experimental mixture of structures and tested our approach as a proof-of-principle. Our

97    experimental mixture consisted of 40S, 60S and 80S ribosomes, apoferritin and β-galactosidase.

98    We collected ~2,400 images and used a template-based particle picking scheme to select

99    ~523,000 particles from the entire data set (Roseman, 2004). Raw micrographs showed a mixture

100   of disperse particles with varying size and shape (Figure S2). We then performed 2D classification

101   on the entire set of particles using RELION (Scheres, 2012). After 1 round of filtering junk

102   particles, the remaining ~203,000 particles were sorted into 100 classes using RELION. The class

103   averages contained many characteristic ribosome projections and had distinct structural features

104   (Figure S2). We were unable to identify any β-galactosidase particles in our collected images.

105   We then applied our SLICEM algorithm to the 100 2D class averages. The identity of each

106   class average was manually annotated, where it was easily recognizable, to assess whether our

107   algorithm was correctly separating the 2D projection images from our heterogeneous mixture

108   (Figure 3). Based on these manual annotations, we again tested the 6 different metrics in a

109   precision-recall framework to determine which metric performed better on experimental data

110   (Figure S3). Interestingly, the Euclidean distance and sum of the absolute difference scoring

111   metric significantly outperformed the cosine similarity. Using the sum of the absolute difference

112   scoring metric, the network naturally partitioned into 3 distinct communities, one for each

113   ribosome, prior to employing any community detection algorithms (Figure 3). As part of our

6

114     algorithm, we evaluated two community detection methods, edge betweenness and walktrap, to

115     determine if the network should be further subdivided (Latapy and Pons, 2004; Newman and

116     Girvan, 2004). We chose to use community detection algorithms to prevent biasing the data by

117     choosing the number of output clusters we expected. Using our SLICEM algorithm, we were able

118     to correctly separate 2D projection images from 3 large, asymmetric macromolecular complexes

119     from the same mixture.

120

121     <u>Relating summed pixel intensity to molecular weight</u>

122     Apart from partitioning 2D projection images into homogenous subsets for 3D reconstruction, one

123     additional goal was to determine the identity of each projection image. In previous studies, we

124     and others have leveraged mass spectrometry data to help identify electron microscopy

125     reconstructions from a heterogeneous mixture, such as cell lysate, where the architecture of every

126     protein or protein complex is not known (Kastritis et al., 2017; Verbeke et al., 2018). However,

127     this combined MS-EM approach was only useful for identifying highly abundant and easily

128     recognizable structures.

129         To provide evidence of macromolecular identity from the electron maps, we calculated the

130     sum of pixel intensities in each manually annotated 2D class average as a proxy for molecular

131     weight (Figure 4). We found that each of the three ribosomes and apoferritin had unique summed

132     pixel intensities that could be used to distinguish the class averages. A least-squares fit to the

133     mean of the summed pixel intensities showed a linear relationship between summed pixel

134     intensity and protein molecular weight. The summed pixel intensities were therefore used as an

135     additional filtering step by removing nodes in communities whose summed pixel intensities were

136     outliers in that community. Using this filtering step, the apoferritin class average was removed

137     from the community containing predominantly 40S ribosome reprojections. Our data suggest that,

138     given an appropriate set of standards, summed pixel intensity can be correlated to molecular

139    weight. Thus, summed pixel intensity could be useful in narrowing down the possible identities for

140    a set of electron densities, when combined with sequence information.

141

142    <u>3D classification of a mixture of protein complexes</u>

143    The ultimate goal of our pipeline is to reconstruct 3D models from our output of clustered 2D

144    projection images. We chose to use cryoSPARC for 3D reconstructions because it can perform

145    heterogeneous reconstruction without *a priori* information on structure or identity (Punjani et al.,

146    2017). We used the particles from each of our 3 distinct communities in addition to the isolated

147    apoferritin node for *ab initio* reconstruction in cryoSPARC (Figure 5). The cluster containing

148    primarily 40S ribosome particles was split into two classes to filter the additional junk particles

149    present in the community. Comparison of our models reconstructed after clustering to the models

150    produced using the entire data set as input for *ab initio* reconstruction in cryoSPARC with 4

151    classes (one for each protein complex in the mixture) showed our pre-sorting procedure improved

152    the resulting structures (Figure 5). In particular, we were able to build an apoferritin model that

153    was missed in the 3D classification of all particles from cryoSPARC. Our 80S model also shows

154    a more complete density for the small subunit than its counterpart in the model created without

155    clustering. We also observe that changing the number of classes using *ab initio* reconstruction in

156    cryoSPARC had a substantial impact on the quality of classification (Figure S4).

157        Each model was refined and evaluated using the gold-standard 0.143 Fourier shell

158    correlation criterion (Figure S5). We obtained easily identifiable 40S, 60S, and 80S ribosome

159    structures at 12, 4, and 5.4 Å resolution, respectively. We were also able to reconstruct the

160    smaller, more compact apoferritin at 19 Å resolution. Notably, the 40S and 80S models contain

161    streaks in one dimension, indicating that we are missing several orientations of the particles. We

162    attribute this to preferential orientation of the particles in ice, rather than an inability of our

163    algorithm to properly sort particles into correct communities. Together, these results demonstrate

8

164    a functioning pipeline for sorting 2D projection images from a heterogeneous mixture of 3D

165    structures, allowing for single particle EM to be applied to samples containing multiple proteins or

166    protein complexes. Importantly, aside from choosing the most appropriate similarity measure, our

167    approach is fully unsupervised, requiring no user defined estimate of the number of existing

168    classes.

169

170    **Discussion**

171    As cryo-EM continues to rapidly advance, one potential application would be to perform high-

172    throughput structural biology. The ability to sort and classify heterogeneous mixtures will become

173    a necessary feature. One advantage of this approach would be to study closer-to-native proteins

174    directly from cell lysate without the need to purify or alter the sample. Currently, handling

175    compositional and conformational heterogeneity is a major challenge for the EM field, usually

176    requiring expert, time-consuming steps. In this study, we present an unsupervised algorithm,

177    SLICEM, which extends on previous methods and demonstrates that sorting 2D projection images

178    based on the similarity of their common lines is capable of correctly clustering 2D projection

179    images from a mixture of protein and protein-nucleic acid complexes. We first demonstrate that

180    the algorithm successfully sorts a synthetic dataset of reprojections created from 35 unique

181    macromolecular structures. Next, we show the same algorithm can successfully partition 2D

182    projection images from an experimental data set containing multiple macromolecular complexes.

183    Pre-sorting 2D projection images prior to 3D classification allows current reconstruction

184    algorithms to be employed on datasets that would otherwise be too complex.

185         Although we demonstrated the feasibility of our approach on synthetic and experimental

186    data, we acknowledge that there are several limitations. In particular, our algorithm relies on the

187    quality of upstream 2D alignment, classification and averaging. As we observed during 2D

188    classification of our cryo-EM data, all apoferritin particles were grouped into a single class

9

189   average. However, during our network generation step, each class average is given multiple

190   edges to the most similar classes, forcing the single apoferritin class average to have multiple

191   spurious edges. This error will occur any time the number of class averages of a given structure

192   is less than the number of edges used in the graph. Future modifications to the algorithm could

193   include searching for symmetric class averages, where this error is more likely to occur, and

194   removing them prior to community detection.

195   As we move cryo-EM towards structural determination of heterogeneous mixtures, several

196   other technical challenges will emerge, such as universal freezing conditions. In our mixture of 5

197   macromolecular complexes, we were unable to easily find freezing conditions that accommodated

198   all proteins. The result was a mixture missing β-galactosidase and containing orientation

199   preferences for the 40S and 80S ribosome. However, previous work has produced e.g. high-

200   resolution structures of fatty-acid synthase from fractionated cell lysate, suggesting it is possible

201   to find suitable cryo-conditions for solutions containing many macromolecular species (Kastritis

202   et al., 2017). An additional challenge will be developing particle picking algorithms specifically for

203   mixtures, where the particle shape may be unknown and, perhaps more importantly, non-uniform.

204   While in this study we used a template picking scheme, future studies with mixtures of unknown

205   composition will require more sophisticated approaches.

206   An expert might be able to manually sort the class averages from our cryo-EM data set;

207   however, as mixtures grow in complexity, manual sorting will certainly become infeasible.

208   Introducing algorithms such as SLICEM will provide an unbiased way to group 2D projection

209   images and can be easily implemented in conjunction with a variety of image processing and 3D

210   reconstruction packages. One additional utility of this algorithm could be to remove junk class

211   averages from data in a semi-supervised manner by removal of communities of projection images

212   that do not appear to have structural features. Our approach for sorting mixtures of structures

213   combined with previous approaches for sorting conformational heterogeneity could be a powerful

214  tool for deep classification. Development of methods to sort mixtures of structures in single

215  particle cryo-EM will allow us to solve more structures in parallel and alleviate time-consuming

216  protein purification and sample preparation.

217

218  **Materials and Methods**

219  <u>Synthetic data generation</u>

220  The following list of PDB entries were used to create the dataset of synthetic reprojections (1A0I,

221  1HHO, 1NW9, 1WA5, 3JCK, 5A63, 1A36, 1HNW, 1PJR, 2FFL, 3JCR, 5GJQ, 1AON, 1I6H, 1RYP,

222  2MYS, 3VKH, 5VOX, 1FA0, 1JLB, 1S5L, 2NN6, 4F3T, 6B3R, 1FPY, 1MUH, 1SXJ, 2SRC, 4V6C,

223  6D6V, 1GFL, 1NJI, 1TAU, 3JB9, 5A1A). Each PDB entry was low-pass filtered to 9 Å and

224  converted to a 3D EM density using 'pdb2mrc' in EMAN (Ludtke et al., 1999). These densities

225  were then uniformly reprojected using 'project3d' in EMAN to create 12 2D reprojections for each

226  structure (Ludtke et al., 1999). Reprojections were centered in 350 Å boxes.

227

228  <u>Purification of apoferritin and β-galactosidase</u>

229  Size-exclusion chromatography was performed at 4 ºC on an AKTA FPLC (GE Healthcare).

230  Approximately 10 mg of apoferritin (Sigma A3660-1VL) and 5 mg of β-galactosidase G5635-5KU

231  were independently applied to a Superdex 200 10/300 GL analytical gel filtration column (GE

232  Healthcare) equilibrated in 20 mM HEPES KOH, 100 mM potassium acetate, 2.5 mM magnesium

233  acetate, pH 7.5 at a flow rate of 0.5 mL min-1. Fractions were collected every 0.5 mL.

234

235  <u>SLICEM Algorithm</u>

236  Our algorithm consists of five main steps: (1) Extracting 2D class average signal from background,

237  (2) Generating 1D line projections from the extracted 2D projection images, (3) Scoring the

11

238   similarity of all pairs of 1D line projections, (4) Building a nearest-neighbors graph of the 2D class

239   averages and (5) Partitioning communities within the graph.

240   (1) Extracting 2D class averages from background

241   The input to our algorithm is a set of centered and normalized 2D class averages. We then extract

242   the centered region of positive pixels values from the zero-mean normalized images to remove

243   background signal and extra densities that might be present in a class average.

244   (2) Generating 1D line projections from extracted 2D projection images

245   Each extracted class average is projected into 1D over 360 degrees in 5 degree intervals by

246   summing the pixel values along the projection axis. The 1D line projections are then

247   independently zero-mean normalized if the normalized cross-correlation or normalized Euclidean

248   scoring metric are selected.

249   (3) Scoring the similarity of all pairs of 1D line projections

250   To score the similarity of the 1D line projections we considered 6 different scoring metrics:

251   Euclidean distance, normalized Euclidean distance, cosine similarity, sum of the absolute

252   difference, cross-correlation and normalized cross-correlation. For the non-cross-correlation

253   metrics, the similarity of the 1D line projections is calculated for translations of the smaller 1D

254   projection across the larger 1D projection if there is a difference in projection size, analogous to

255   the 'sliding' feature of cross-correlations. The optimum score during the translations is kept for

256   each pair of 1D projections. After pairwise scoring of all 1D line projections, the similarity between

257   each pair of 2D class averages is defined by their respective highest scoring 1D line projections.

258   (4) Building a nearest-neighbors graph of the 2D class averages

259   SLICEM then constructs a directed graph using the similarity scores calculated for each pair of

260   2D class averages. Each node (2D class average) is connected to the 5 most similar (top scoring)

261   2D class averages. Each edge is assigned a weight computed as a z-score relative to all scores

262   for a given 2D class average.

12

263   (5) Partitioning communities within the graph.

264   The resulting graph is then subdivided using a community detection algorithm. Specifically, we

265   evaluated the edge-betweenness and walktrap algorithms to define clusters in the graph. Then,

266   the median absolute deviation of summed pixel intensities for each node is calculated to remove

267   outliers from the cluster. The final set of nodes in a cluster is then used as input for 3D

268   reconstruction in cryoSPARC.

269

270   Cryo-EM grid preparation and data collection

271   C-flat holey carbon grids (CF-1.2/1.3, Protochips Inc.) were pre-coated with a thin layer of freshly

272   prepared carbon film and glow-discharged for 30 seconds using a Gatan Solarus plasma cleaner

273   before addition of sample. 2.5 µl of a mixture of 75 nM 40S ribosome, 150 nM 60S ribosome, 50

274   nM 80S ribosome, 125 nM apoferritin and 125 nM β-galactosidase were placed onto grids, blotted

275   for 3 seconds with a blotting force of 5 and rapidly plunged into liquid ethane using a FEI Vitrobot

276   MarkIV operated at 4 °C and 100% humidity. Data were acquired using an FEI Titan Krios

277   transmission electron microscope (Sauer Structural Biology Laboratory, University of Texas at

278   Austin) operating at 300 keV at a nominal magnification of ×22,500 (1.1 Å pixel size) with defocus

279   ranging from -2.0 to -3.5 µm. The data were collected using a total exposure of 6 s fractionated

280   into 20 frames (300 ms per frame) with a dose rate of ~8 electrons per pixel per second and a

281   total exposure dose of ~40 e$^-$ Å$^{-2}$. A total of 2,423 micrographs were automatically recorded on a

282   Gatan K2 Summit direct electron detector operated in counting mode using the MSI Template

283   application within the automated macromolecular microscopy software LEGINON (Suloway et al.,

284   2005).

285

286

287

288    <u>Cryo-EM data processing</u>

289    All image pre-processing was performed in Appion (Lander et al., 2009). Individual movie frames

290    were aligned and averaged using 'MotionCor2' drift-correction software (Zheng et al., 2017).

291    These drift-corrected micrographs were binned by 8, and bad micrographs and/or regions of

292    micrographs were removed using the 'manual masking' command within Appion. A total of

293    522,653 particles were picked with a template-based particle picker using a reference-free 2D

294    class average from a small subset of manually picked particles as templates. The contrast transfer

295    function (CTF) of each micrograph was estimated using CTFFIND4 (Rohou and Grigorieff, 2015).

296    Selected particles were extracted from micrographs using particle extraction within RELION

297    (Scheres, 2012) and the EMAN2 coordinates exported from Appion. Two rounds of reference free

298    2D classification with 100 classes for each sample were performed in RELION to remove junk

299    particles, resulting in a clean stack of 202,611 particle images.

300

313     an Army Young Investigator supported by the Army Research Office (WW911NF-19-10021). This

314     work was supported in part by Welch Foundation Research Grants F-1938 (to D.W.T.) and F-

315     1515 (to E.M.M.), Army Research Office Grant (W911NF-15-0120) (to D.W.T.), Robert J. Kleberg,

316     Jr. and Helen C. Kleberg Foundation Medical Research Award (to D.W.T.) and grants from the

317     National Institutes of Health (GM122480, DK110520, HD085901) to E.M.M..

318

**Data Availability**

320     The cryo-EM reconstructions of the 40S, 60S, 80S, and apoferritin have been deposited in the

321     Electron Microscopy Databank with accession codes EMD-20109, EMD-20110, EMD-20111 and

322     EMD-20112, respectively. The motion-corrected sum micrographs have been deposited into

323     EMPIAR with accession code EMPIAR-10268. Computer code for SLICEM is available at

324     https://github.com/marcottelab/SLICEM.

325

**References**

327     Aizenbud, Y., and Shkolnisky, Y. (2016). A max-cut approach to heterogeneity in cryo-electron
328     microscopy. ArXiv160901100 Cs Math Q-Bio.

329     Cianfrocco, M.A., and Leschziner, A.E. (2015). Low cost, high performance processing of single
330     particle cryo-electron microscopy data in the cloud. ELife *4*.

331     Herman, G.T., and Kalinowski, M. (2008). Classification of heterogeneous electron microscopic
332     projections into homogeneous subsets. Ultramicroscopy *108*, 327–338.

333     Kastritis, P.L., O'Reilly, F.J., Bock, T., Li, Y., Rogon, M.Z., Buczak, K., Romanov, N., Betts, M.J.,
334     Bui, K.H., Hagen, W.J., et al. (2017). Capturing protein communities by structural proteomics in
335     a thermophilic eukaryote. Mol. Syst. Biol. *13*, 936.

336     Katsevich, E., Katsevich, A., and Singer, A. (2015). Covariance Matrix Estimation for the Cryo-
337     EM Heterogeneity Problem. SIAM J. Imaging Sci. *8*, 126–185.

338     Kühlbrandt, W. (2014). The resolution revolution. Science *343*, 1443–1444.

339     Lander, G.C., Stagg, S.M., Voss, N.R., Cheng, A., Fellmann, D., Pulokas, J., Yoshioka, C.,
340     Irving, C., Mulder, A., Lau, P.-W., et al. (2009). Appion: An integrated, database-driven pipeline
341     to facilitate EM image processing. J. Struct. Biol. *166*, 95–102.

342  Latapy, M., and Pons, P. (2004). Computing communities in large networks using random
343  walks. ArXivcond-Mat0412368.

344  Liao, H.Y., Hashem, Y., and Frank, J. (2015). Efficient Estimation of Three-Dimensional
345  Covariance and its Application in the Analysis of Heterogeneous Samples in Cryo-Electron
346  Microscopy. Structure *23*, 1129–1137.

347  Ludtke, S.J., Baldwin, P.R., and Chiu, W. (1999). EMAN: Semiautomated Software for High-
348  Resolution Single-Particle Reconstructions. J. Struct. Biol. *128*, 82–97.

349  Newman, M.E.J., and Girvan, M. (2004). Finding and evaluating community structure in
350  networks. Phys. Rev. E *69*.

351  Penczek, P.A., Frank, J., and Spahn, C.M.T. (2006). A method of focused classification, based
352  on the bootstrap 3D variance analysis, and its application to EF-G-dependent translocation. J.
353  Struct. Biol. *154*, 184–194.

354  Punjani, A., Rubinstein, J.L., Fleet, D.J., and Brubaker, M.A. (2017). cryoSPARC: algorithms for
355  rapid unsupervised cryo-EM structure determination. Nat. Methods *14*, 290–296.

356  Rohou, A., and Grigorieff, N. (2015). CTFFIND4: Fast and accurate defocus estimation from
357  electron micrographs. J. Struct. Biol. *192*, 216–221.

358  Roseman, A. (2004). FindEM—a fast, efficient program for automatic selection of particles from
359  electron micrographs. J. Struct. Biol. *145*, 91–99.

360  Russo, C.J., and Passmore, L.A. (2014). Ultrastable gold substrates for electron
361  cryomicroscopy. Science *346*, 1377–1380.

362  Scaiola, A., Peña, C., Weisser, M., Böhringer, D., Leibundgut, M., Klingauf-Nerurkar, P.,
363  Gerhardy, S., Panse, V.G., and Ban, N. (2018). Structure of a eukaryotic cytoplasmic pre-40S
364  ribosomal subunit. EMBO J. 13.

365  Scheres, S.H.W. (2012). RELION: Implementation of a Bayesian approach to cryo-EM structure
366  determination. J. Struct. Biol. *180*, 519–530.

367  Shatsky, M., Hall, R.J., Nogales, E., Malik, J., and Brenner, S.E. (2010). Automated multi-model
368  reconstruction from single-particle electron microscopy data. J. Struct. Biol. *170*, 98–108.

369  Shen, P.S., Park, J., Qin, Y., Li, X., Parsawar, K., Larson, M.H., Cox, J., Cheng, Y., Lambowitz,
370  A.M., Weissman, J.S., et al. (2015). Rqc2p and 60S ribosomal subunits mediate mRNA-
371  independent elongation of nascent chains. Science *347*, 75–78.

372  Sigworth, F.J. (1998). A Maximum-Likelihood Approach to Single-Particle Image Refinement. J.
373  Struct. Biol. *122*, 328–339.

374  Sigworth, F.J., Doerschuk, P.C., Carazo, J.-M., and Scheres, S.H.W. (2010). An Introduction to
375  Maximum-Likelihood Methods in Cryo-EM. In Methods in Enzymology, (Elsevier), pp. 263–294.

376    Suloway, C., Pulokas, J., Fellmann, D., Cheng, A., Guerra, F., Quispe, J., Stagg, S., Potter,
377    C.S., and Carragher, B. (2005). Automated molecular microscopy: The new Leginon system. J.
378    Struct. Biol. *151*, 41–60.

379    Van Heel, M. (1987). Angular reconstitution: A posteriori assignment of projection directions for
380    3D reconstruction. Ultramicroscopy *21*, 111–123.

381    Verbeke, E.J., Mallam, A.L., Drew, K., Marcotte, E.M., and Taylor, D.W. (2018). Classification of
382    Single Particles from Human Cell Extract Reveals Distinct Structures. Cell Rep. *24*, 259-268.e3.

383    Yi, X., Verbeke, E.J., Chang, Y., Dickinson, D.J., and Taylor, D.W. (2018). Electron microscopy
384    snapshots of single particles from single cells. J. Biol. Chem. jbc.RA118.006686.

385    Zheng, S.Q., Palovcak, E., Armache, J.-P., Verba, K.A., Cheng, Y., and Agard, D.A. (2017).
386    MotionCor2: anisotropic correction of beam-induced motion for improved cryo-electron
387    microscopy. Nat. Methods *14*, 331–332.

388

389

390   **Figure Legends**

391   **Figure 1. Computational pipeline for SLICEM**

392   Individual particle images are averaged after reference-free 2D alignment and classification.

393   Using a Radon transform, 1D line projections are created from the 2D class averages (referred

394   to as 2D projections). Each 1D line projection from every 2D projection is then scored for

395   similarity. The top scores between each projection are used to form edges connecting 2D

396   projections that have a similar 1D line projection to form a graph. 2D projection images are then

397   partitioned into groups belonging to the same putative structure using a community detection

398   algorithm. Individual particle images belonging to each 2D projection within a community are

399   subjected to *ab initio* 3D reconstruction.

400

401   **Figure 2. Separating mixtures of synthetic 2D reprojections**

402   Synthetic reprojections were generated from 35 distinct X-ray crystal structures low-pass filtered

403   to 9 Å from complexes ranging in molecular weight from ~30 – 3000 kDa, prior to separation

404   using SLICEM. (A) Precision-recall plot ranking 6 different metrics at scoring the similarity

405   between 1D line projections from each 2D reprojection. (B) Network output displaying

406   communities of 2D reprojection images determined using SLICEM. Each node represents a 2D

407   reprojection with 5 connecting edges to the most similar reprojections as scored using

408   Euclidean distance. The color of each node matches the structure from which it was reprojected

409   (shown as a surface).

410

411   **Figure 3. Experimental 2D class averages and resulting network**

412   Cryo-EM data was collected on a mixture of 5 protein and protein-nucleic acid complexes. (A)

413   Representative 2D class averages of the 4 complexes identified in the mixture. The identity of

18

414     each class average was manually annotated were it could be easily identified. The class average

415     corresponding to apoferritin was further subdivided into multiple classes for visualization. (B)

416     Network generated using SLICEM on the 100 2D class averages scored using the sum of the

417     absolute difference metric. Nodes representing each 2D class averages are colored by their

418     putative structural identity. The width of the box corresponds to 422 Å.

419

420     **Figure 4. Summed pixel intensities of 2D class averages correlate to molecular weight**

421     (A) 2D to 1D projections for representative 2D class averages of each structure present in the

422     mixture. 1D projection plots show the line profile for a single projection of each 2D class average.

423     Pixel heat maps show the intensity of the line profile at each pixel. (B) Distribution of the summed

424     1D projection pixel intensities, or integration of the 1D line profiles, calculated for each 2D class

425     average. Summed pixel intensities for each manually identified 2D class average are plotted

426     against their respective molecular weight. Black points are the mean summed pixel intensity for

427     each structure.

428

429     **Figure 5. *Ab initio* structures from an experimental mixture**

430     (A) High-resolution structures of the 80S ribosome EMD-2858 (Cianfrocco and Leschziner, 2015),

431     60S ribosome EMD-2811 (Shen et al., 2015), 40S ribosome EMD-4214 (Scaiola et al., 2018) and

432     apoferritin EMD-2788 (Russo and Passmore, 2014). (B) 3D models of the 80S ribosome, 60S

433     ribosome, 40S ribosome and apoferritin generated by sorting particles using SLICEM prior to *ab*

434     *initio* 3D reconstruction in cryoSPARC. (C) 3D models generated using *ab initio* reconstruction to

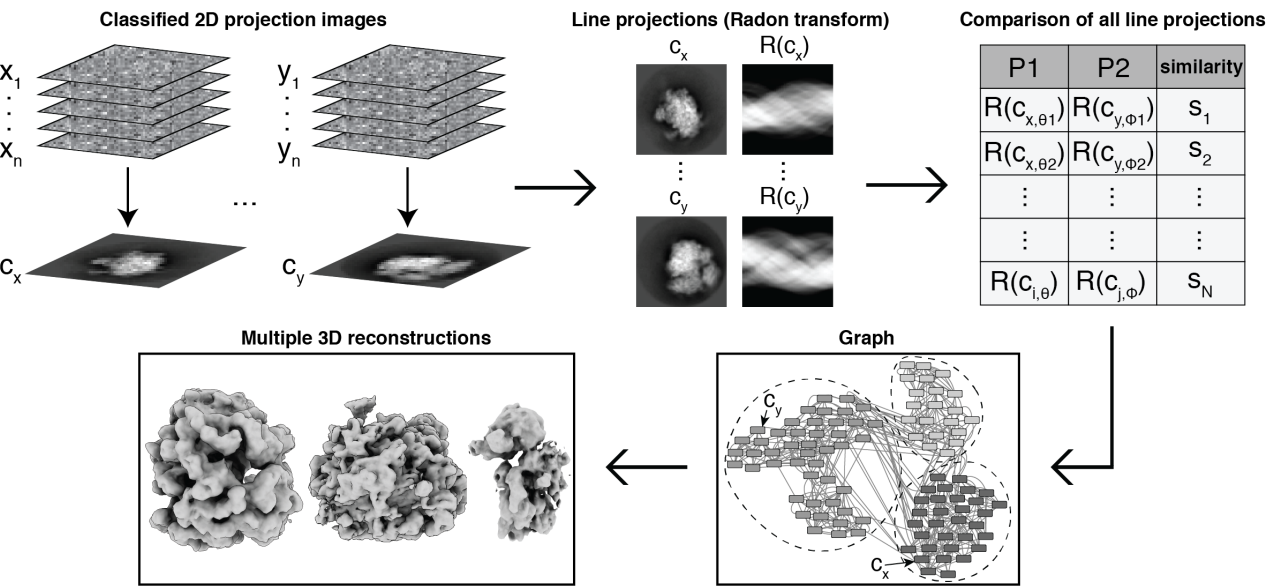435     generate 4 classes in cryoSPARC without pre-sorting particles using SLICEM.
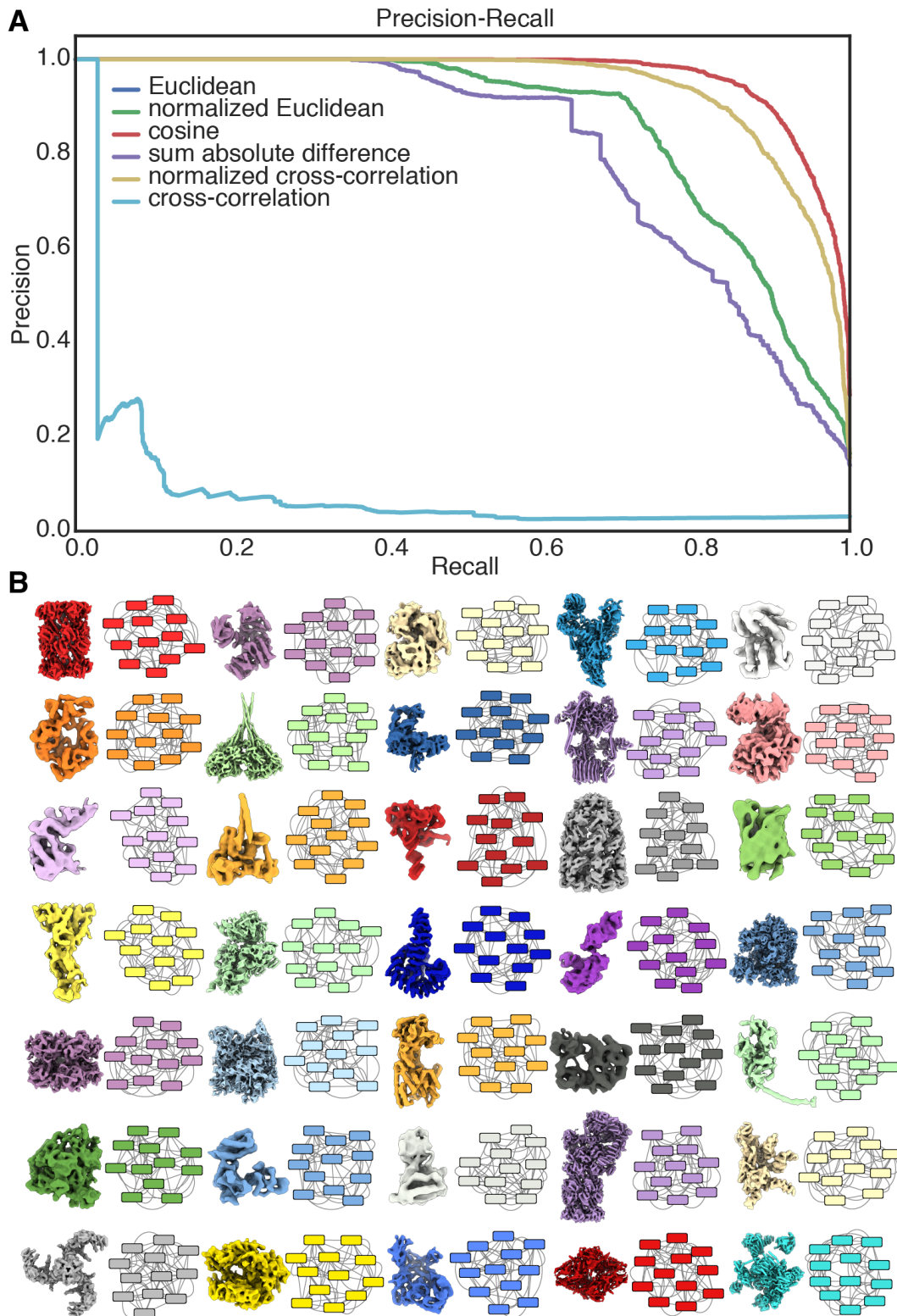
436

19

437     **Figures:**

438     **Figure 1**

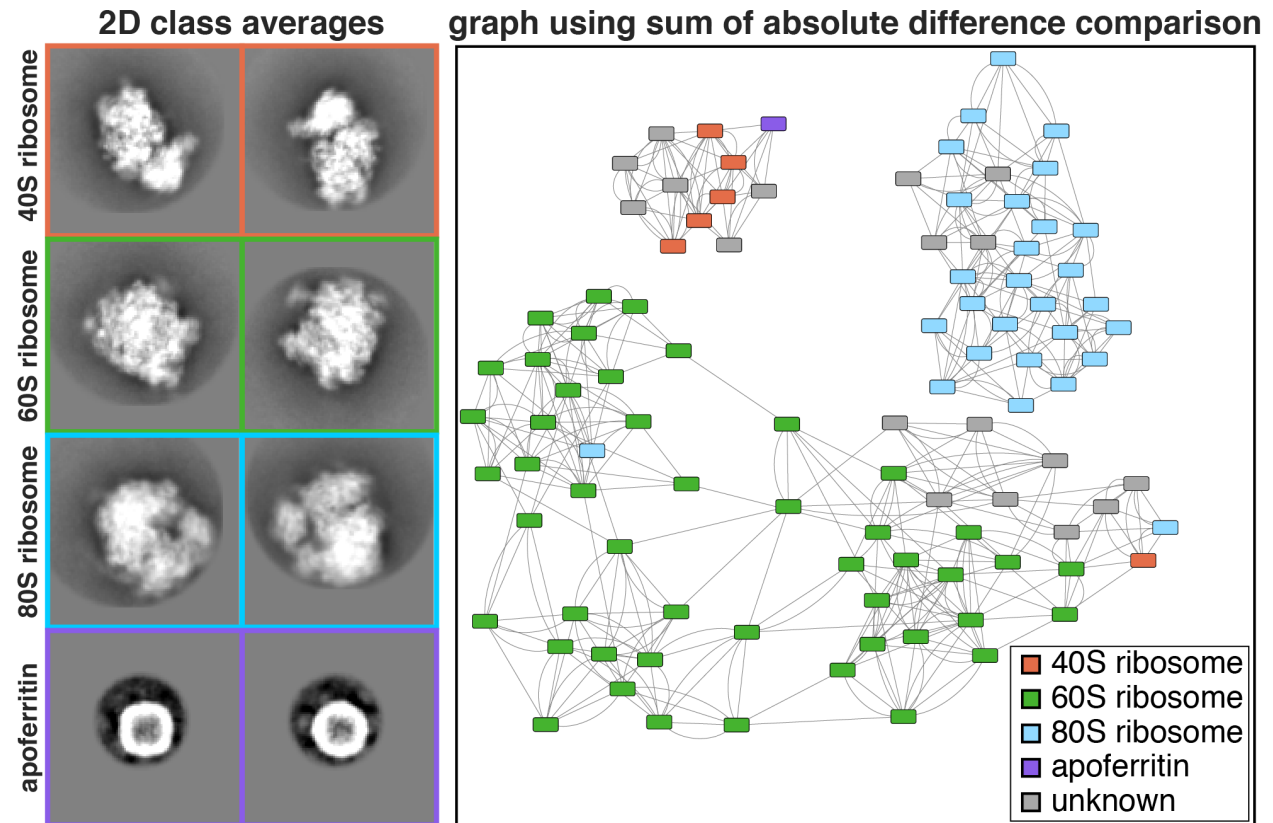439

**Figure 2**

**Figure 3**

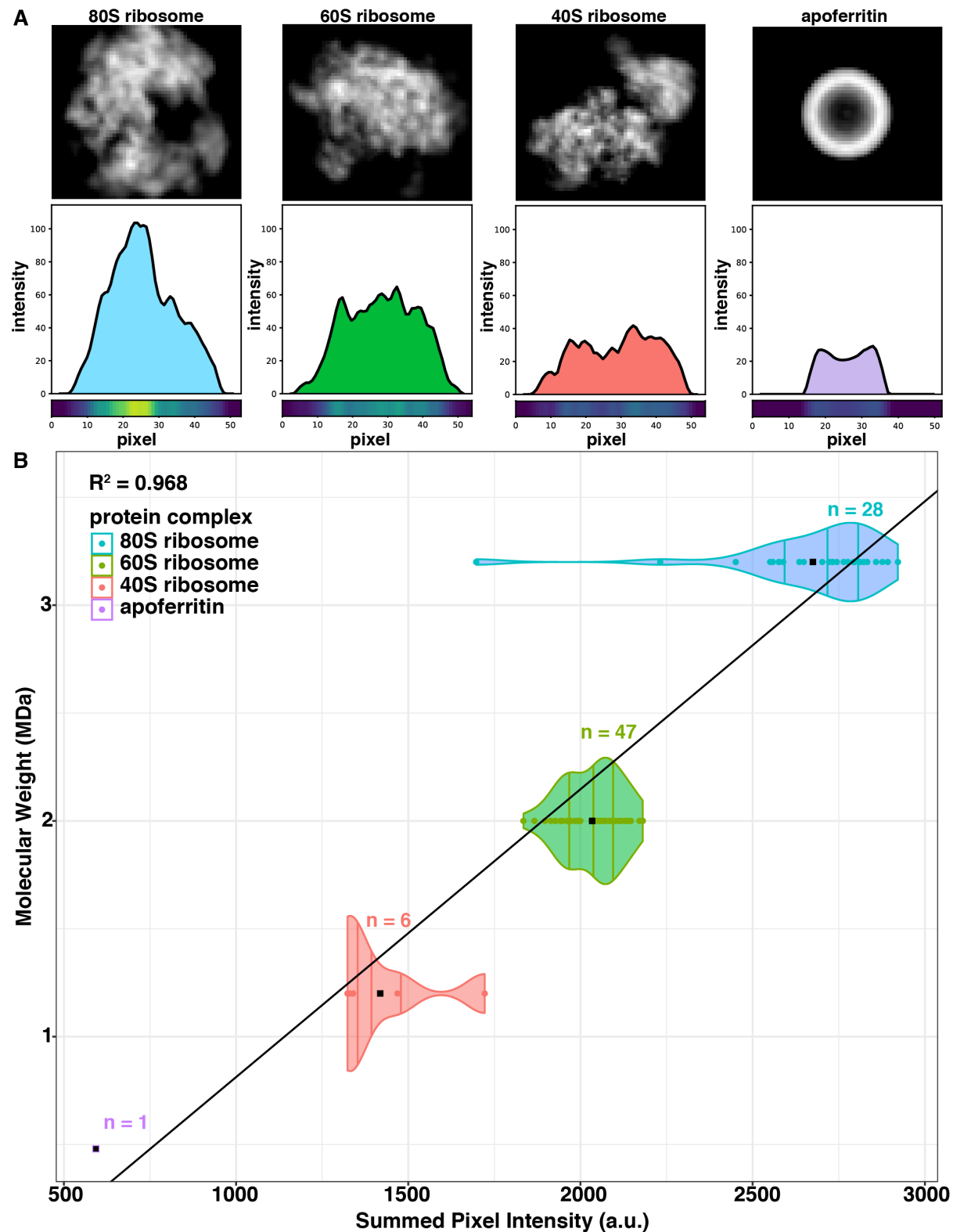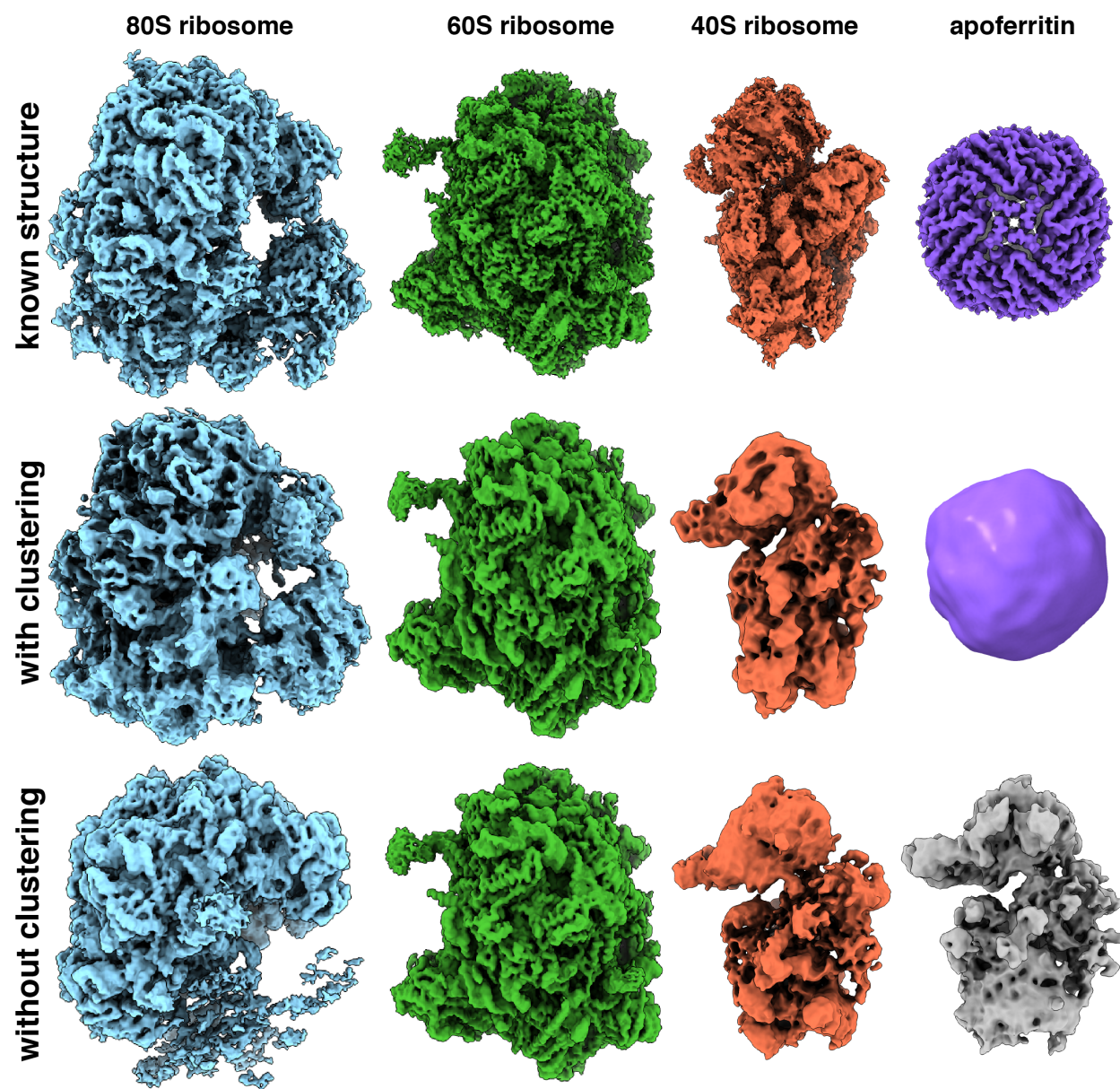**Figure 4**

**Figure 5**

562 **Supplemental Information**

563 Figure S1. 2D reprojections from synthetic dataset

564 Subset of 2D reprojections from 12 of the 35 structures in our synthetic dataset. Box size

565 corresponds to 300 Å.

566

567 Figure S2. 2D classification of particles using RELION

568 (A) Representative raw micrograph of a mixture containing 40S, 60S and 80S ribosomes,

569 apoferritin and β-galactosidase. (B) Reference-free 2D class averages generated using RELION

570 of ~203,000 template-picked particle images. Box size corresponds to 422 Å.

571

572 Figure S3. Precision-recall curves for experimental cryo-EM data

573 Precision-recall plot displaying 6 different metrics for scoring the similarity between 1D line

574 projections from the entire set of 2D class averages.

575

576 Figure S4. *Ab initio* reconstructions in cryoSPARC with varying class number

577 3D reconstructions using *ab initio* reconstruction in cryoSPARC from the entire data set with K =

578 3, 4, 5 and 6 classes, respectively.

579

580 Figure S5. Fourier shell correlations curves

581 FSC curves for our clustered 80S ribosome (blue), 60S ribosome (green), 40S ribosome (red)

582 and apoferritin (purple) shown in Figure 5B. Nominal resolutions were estimated to be 5.4, 4, 12

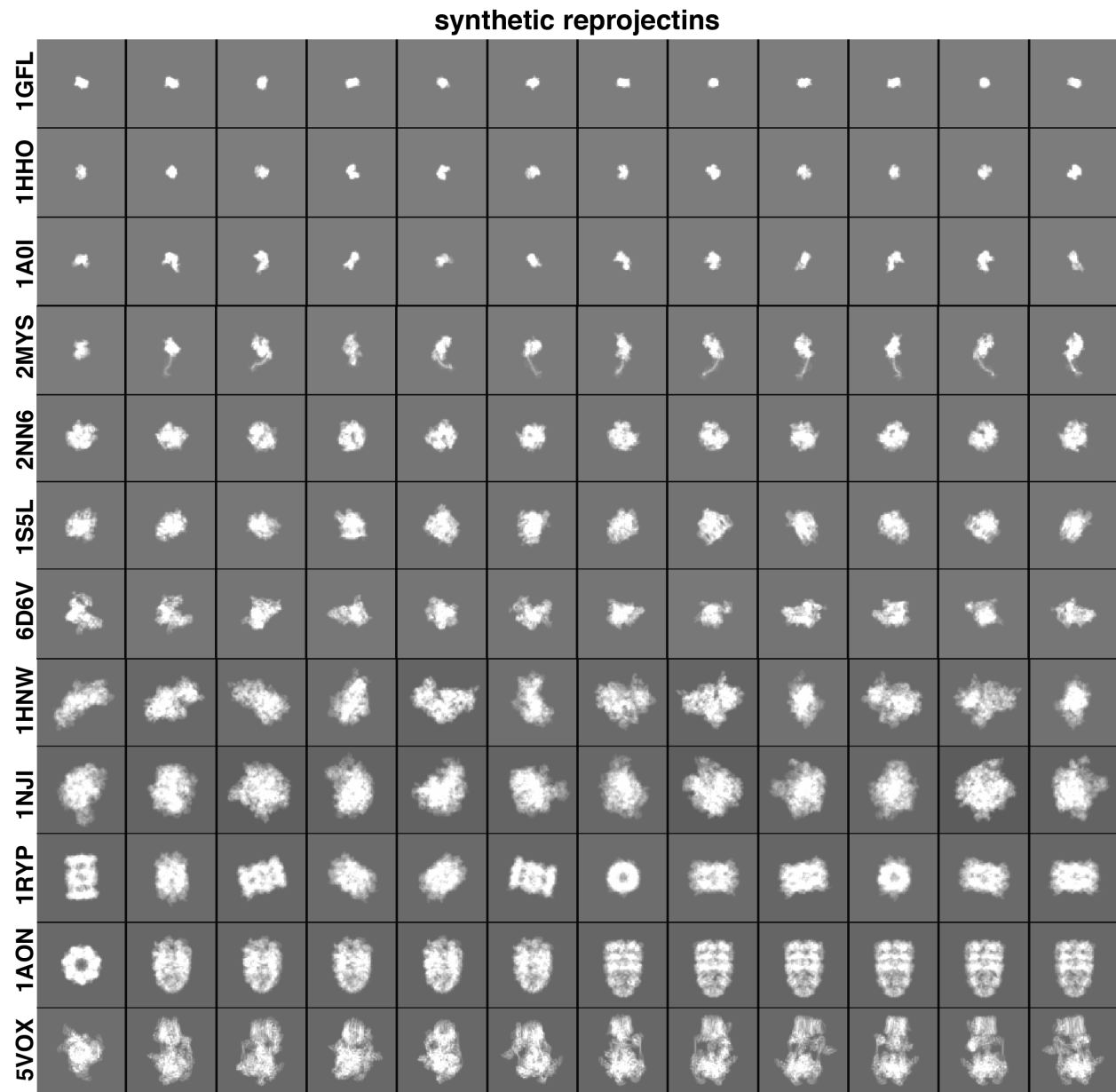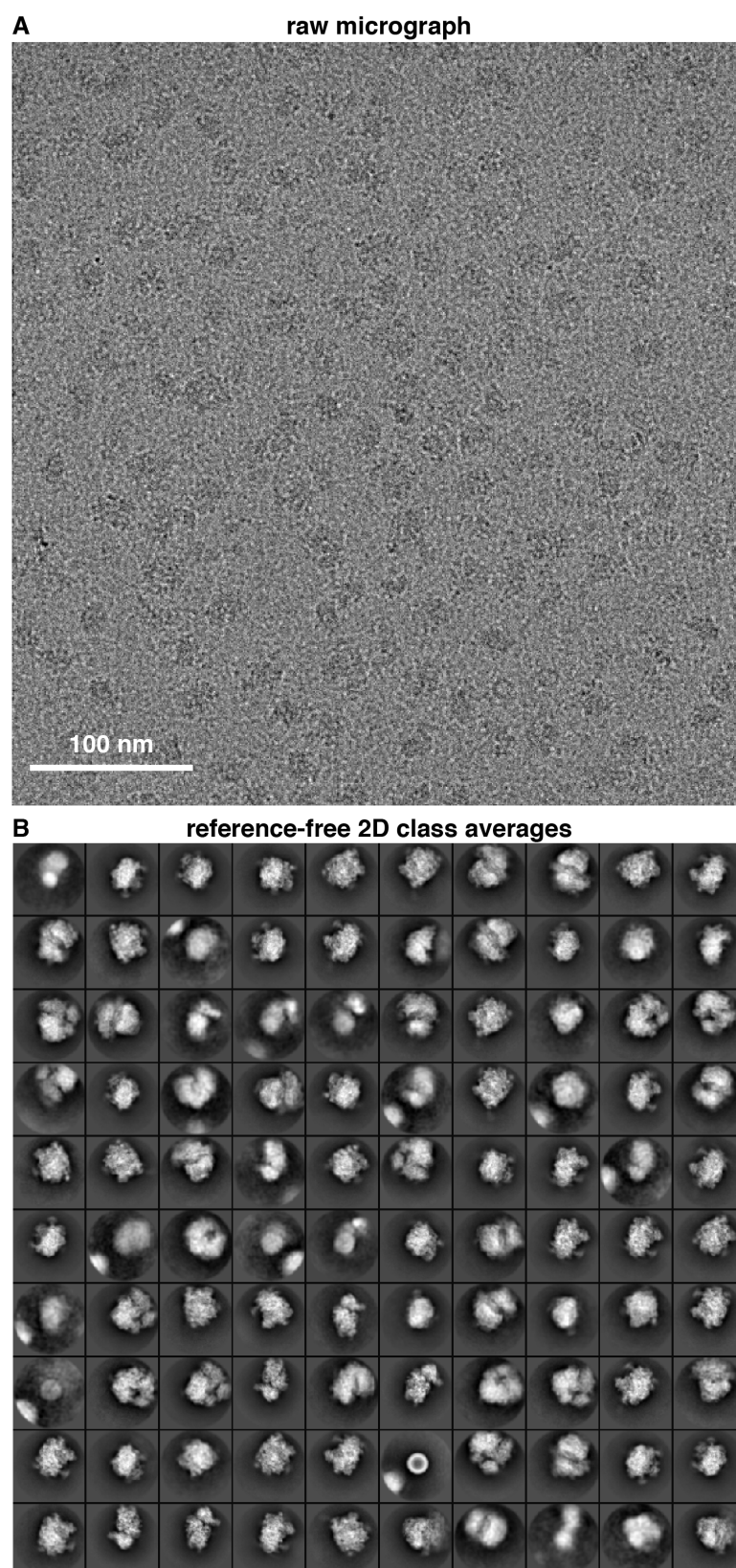583 and 19 Å, respectively, using the 0.143 gold-standard FSC criterion.
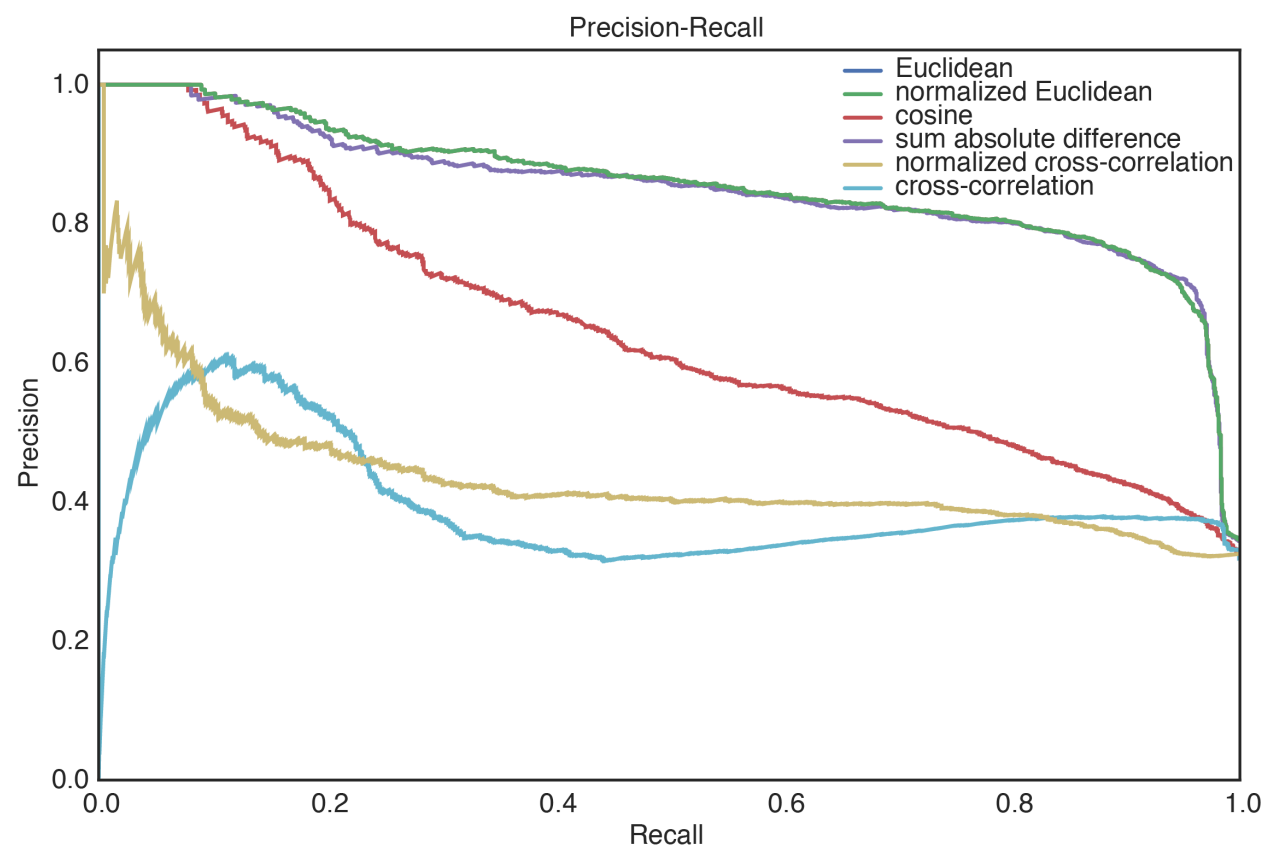
584

585

586

587    **Figure S1**

588



606

607

608

609

610

611

**Figure S2**

**A**  **raw micrograph**



100 nm

**B**  **reference-free 2D class averages**



27

**Figure S3**



Precision-Recall plot showing Precision (y-axis) versus Recall (x-axis) for six methods: Euclidean, normalized Euclidean, cosine, sum absolute difference, normalized cross-correlation, cross-correlation.

662 **Figure S4**



Fourier shell correlation curves

687    **Figure S5**



cryoSPARC models

footer: 30