

1 Map segmentation, automated model-building and their application to the Cryo-EM Model 2 Challenge

3
4

5 Thomas C. Terwilliger^{1,2}, Paul D. Adams^{3,4}, Pavel V. Afonine^{3,5}, Oleg V. Sobolev³

6
7

¹*Los Alamos National Laboratory, Los Alamos NM 87545 USA*

8

²*New Mexico Consortium, Los Alamos NM 87544 USA*

9

³*Molecular Biophysics & Integrated Bioimaging Division, Lawrence Berkeley National
10 Laboratory, Berkeley, CA 94720-8235, USA*

11

⁴*Department of Bioengineering, University of California Berkeley, Berkeley, CA, USA*

12

⁵*Department of Physics and International Centre for Quantum and Molecular Structures,
13 Shanghai University, Shanghai, 200444, People's Republic of China*

14

15

16 E-mail: tterwilliger@newmexicoconsortium.org

17

18

19 Abstract

20

21 A recently-developed method for identifying a compact, contiguous region representing the
22 unique part of a density map was applied to 218 cryo-EM maps with resolutions of 4.5 Å or
23 better. The key elements of the segmentation procedure are (1) identification of all regions of
24 density above a threshold and (2) choice of a unique set of these regions, taking symmetry into
25 consideration, that maximize connectivity and compactness. This segmentation approach was
26 then combined with tools for automated map sharpening and model-building to generate
27 models for the 12 maps in the 2016 cryo-EM model challenge in a fully automated manner. The
28 resulting models have completeness from 24% to 82% and RMS distances from reference
29 interpretations of 0.6 Å to 2.1 Å.

30

31 Introduction

32

33 In the 2016 Cryo-EM Modeling Challenge (see
34 http://challenges.emdatabank.org/?q=model_challenge; accessed 2017-11-19), a total of 12
35 maps were supplied to contestants along with reconstruction symmetry and the sequences of
36 the molecules present. One of the goals of the Challenge was to fully interpret such a map
37 given only the map, the symmetry and the sequence information. There are a number of tools
38 being developed by several groups for automated interpretation of cryo-EM maps (DiMaio and
39 Chiu, 2016). These include methods for identification of secondary structure (Jiang et al., 2001;
40 Kong and Ma, 2003; Kong et al., 2004; Baker, Ju and Chiu, 2007), methods for combination of
41 structure-modeling tools such as Rosetta with cryo-EM model-building (Lindert et al., 2012;
42 Wang et al., 2015; Frenz et al., 2017), semi-automated tools for full map interpretation (Baker
43 et al., 2011), and automated tools based on chain-tracing (Chen et al., 2016; Collins and Si,
44 2017) and template-matching approaches (Zhou, Wang and Wang, 2017).

45
46 Prior to the 2016 cryo-EM Model Challenge, we had begun development of software for
47 automatic map sharpening (*phenix.auto_sharpen*; Terwilliger et al., 2018a) and interpretation
48 of density in cryo-EM maps (*phenix.map_to_model*; Terwilliger et al., 2018b) as part of the
49 *Phenix* software package (Adams et al, 2010). It was possible in principle to apply these tools
50 directly to the 2016 Challenge, interpreting an entire map and ignoring the symmetry of the
51 map. It seemed however that it would be more efficient to work with just the unique part of a
52 map. We reasoned that this could be done by identifying a unique part of map that contained a
53 complete molecule, interpreting that part of the map, and then expanding the result using the
54 symmetry in the map to represent the entire map. Powerful tools existed for map
55 segmentation (e.g., Volkman, 2002; Baker, Chiu and Bajaj, 2006; Yu and Bajaj, 2008; Pintilie et
56 al., 2010), but we wanted to be able to integrate the segmentation and symmetry analysis with
57 automated model-building so that information from model-building could be used to make the
58 final choice of the regions of density representing a single molecular unit. We therefore
59 developed a new *Phenix* tool, *phenix.segment_and_split_map* (Terwilliger et al., 2018b) which
60 could be used for this purpose. Here we describe the application of
61 *phenix.segment_and_split_map* to a set of 218 cryo-EM maps selected to generally represent
62 the unique currently-available cryo-EM maps with resolution of 4.5 Å or better. We then
63 describe map segmentation, sharpening, and model-building (Terwilliger et al., 2018b) applied
64 to the 12 cryo-EM maps in the 2016 Cryo-EM Model Challenge.

65

66

67 Methods

68

69 Summary of map segmentation

70

71 The main goal of our segmentation procedure is to identify the density in a map that
72 corresponds to the unique part of that map. (Note that we use “density” to refer to map values.
73 They can be electron density, electric potential, or any other quantity that is being used to
74 describe the locations of atoms in the map). A secondary goal is to choose this density in such
75 a way that it corresponds as closely as possible to the unique biological unit in the map. Our
76 overall approach to map segmentation is (1) to identify all regions of density above an
77 automatically-determined threshold, and (2) to choose a unique set of density regions that
78 maximizes connectivity and compactness, taking into account the symmetry that is present. By
79 default, the process is repeated with a new threshold after removing the density that has been
80 used in the first iteration. The density threshold for consideration of a region of density is
81 chosen to yield a specific volume fraction (typically 20%) of the region of the macromolecule
82 above the threshold. The map is divided into regions of density above the threshold density,
83 where each region is composed of points above the threshold and that have at least one
84 neighbor above the threshold. A unique set of regions is chosen using the symmetry (if any)
85 supplied by the user and the criteria that the unique set should be as compact and connected
86 as feasible. The details of this segmentation procedure have recently been described
87 (Terwilliger et al., 2018b).

88

89 Symmetry present in a map

90

91 We identified symmetry relationships that were applied during map reconstruction using a
92 simplified version of approaches described by Zhang et al., (2012). In many cases the symmetry
93 applied during reconstruction is specified in the EM Data Bank (EMDB, Electron Microscopy
94 Data Bank; Lawson et al., 2016; as for example “I” for icosahedral reconstructions, “C6” for a 6-
95 fold symmetry axis). In others, the symmetry is specified in meta-data associated with the
96 deposited model in the Protein Data Bank (PDB; Bernstein et al., 1977; Berman et al., 2000). In
97 still others, the model deposited in the PDB contains symmetry-related copies which we
98 extracted with the *Phenix* tool *phenix.simple_ncs_from_pdb*. If present, we used symmetry
99 from the deposited models and their meta-data, and if not, we used the information from the
100 EMDB or literature specified in the deposition and the assumption that principal symmetry axes
101 (i.e., screw axes, rotational axes) are generally along the principal axes of the reconstruction to
102 find reconstruction symmetry in the density maps.

103

104

105 Map-model correlations

106

107 We calculated map-model correlations using the *Phenix* tool *phenix.map_model_cc*. This tool
108 identifies the region occupied by the model as all grid points in the target map within a
109 specified distance (typically 3 Å) of an atom in the model. Then it generates a model-based map
110 on the same grid and calculates the correlation of density values between the target map and
111 the model-based map inside the region occupied by the model.

112

113 Model-based maps were calculated in reciprocal space using elastic atomic scattering factors of
114 electrons for neutral atoms as described (Colliex et al., 2006, Afonine et al., 2018b). These
115 scattering factors are framed as the sum of Gaussian terms, represent electric potential, and
116 assume that all atoms are independent. These scattering factors do not include the effects of
117 charged residues and therefore they may be substantially incorrect for certain atoms, including
118 phosphates in RNA or DNA and side chains such as aspartate and glutamate. As improved
119 representations of electron scattering expressed as sums of Gaussian terms these become
120 available these can readily be incorporated in the *Phenix* framework.

121

122

123 Data used for map segmentation

124

125 We selected a group of 218 cryo-EM maps to test our segmentation algorithms. We started
126 with 492 maps we could extract from the EMDB in August of 2017 with simple *Phenix* tools and
127 that were reconstructed at resolutions of 4.5 Å or better. We excluded 91 maps where the
128 resolution in the EMDB and PDB differed by 0.2 Å or more or was not reported, 24 maps where
129 the map-model correlation was less than 0.3, and 16 maps for which the signal-to-noise in map
130 sharpening (Terwilliger et al., 2018a) was less than 3. We then removed map-model pairs that
131 were largely duplications by clustering based on sequence identity using a cutoff of 95%
132 identity and choosing the highest-resolution representative of each group. The sequence

133 identity of two structures was calculated after alignment of each chain in the first structure
134 with the closest-matching chain in the second structure. If either sequence was contained
135 within the other, the identity was considered to be 100%. Otherwise if the lengths of the
136 sequences differed by more than 5% , or the percentage of residues in all chains of the first
137 structure matching a corresponding residue in the second structure was less than 95%, the
138 sequences were considered to be different. Four of the maps in this set were associated with
139 two models, so one map-model pair was set aside for each of these, yielding 218 map-model
140 pairs that were analyzed in this work.

141

142 Evaluating the results of map segmentation by calculation of fraction of molecular unit within
143 the segmented region

144

145 We estimated the fraction of the molecular unit within the segmented region of a map from a
146 comparison of map-model correlations. Our method is related to the cross-correlation
147 variation metric described by Zhang et al. (2012) but it is extended to make an estimate of the
148 fraction of the molecular unit that matches the segmented map. The segmented map has
149 values of zero everywhere outside the segmented region of the map. The overall idea is that if
150 the segmented region contains a complete molecular unit, then the map-model correlation
151 between one complete molecular unit and the segmented map will be the same as the map-
152 model correlation with the original map. On the other hand, if the segmented region contains
153 part of one molecular unit and parts of symmetry-related ones, then the map-model
154 correlation between one intact molecular unit and the segmented map will be lower than the
155 correlation to the original map. We use this difference in map correlation to estimate the
156 fraction of a complete molecule that is within the segmented region.

157

158 We first calculated the map-model correlation between the original map and a map calculated
159 from single molecular unit extracted from the deposited model of the structure. Then we
160 calculated the map-model correlation between the segmented map and a single molecular unit.
161 The square of the ratio of these correlations is (see below) approximately equal to the fraction
162 of the molecular volume that is within the segmented map. The single molecular unit to
163 compare with the map was chosen to be a set of chains representing the unique part of the
164 deposited model. In cases with symmetry, each symmetry-related choice of molecular unit was
165 considered and the one with the highest map-model correlation was chosen.

166

167 The relationship between the map-model correlation for a single molecular unit and the
168 original map compared to the correlation for a molecular unit and a segmented map can be
169 calculated in a straightforward fashion with one assumption. This assumption is that the local
170 map-model correlation for the original map and this single molecular unit is approximately the
171 same everywhere in the region of the model. With this assumption, we can readily calculate
172 the effect of setting all but a fraction f of the map density in the region of the model to zero.
173 This corresponds to calculating the map-model correlation of the segmented map to one full
174 molecular unit, where a fraction f of the molecular unit is present in the segmented map.

175

176 The correlation coefficient CC between two maps with density values represented by D_1 and D_2
177 can be written (after adjusting each map to set the mean density for each to zero so that $\langle D_1 \rangle$
178 $= \langle D_2 \rangle = 0$) as,

$$180 \quad CC = \langle D_1 D_2 \rangle / \sqrt{\langle D_1^2 \rangle \langle D_2^2 \rangle} , \quad (1)$$

181
182 where the calculation in this case is carried out over all the grid points near the model. Now
183 suppose we create a new map D_1' in which we set D_1 to zero at a fraction $(1-f)$ of these grid
184 points. Referencing Eq. (1), this means that the values of $D_1' D_2$ and $D_1'^2$ will be zero at all these
185 grid points, but D_2^2 will be the same. Assuming then that the values of $\langle D_1 D_2 \rangle$, $\langle D_1^2 \rangle$, and
186 $\langle D_2^2 \rangle$ are approximately the same everywhere near the model, we can write that the map-
187 model correlation for the segmented map (CC') with all but $(1-f)$ of the map set to zero is
188 related to the map-model correlation for the original map (CC) by,

$$190 \quad CC' = CC \sqrt{f} , \quad (2)$$

191
192 so that f , the fraction inside the mask, is approximately given by $f = CC'^2 / CC^2$.

193
194
195 Automated model-building

196
197 The *Phenix* tool *phenix.map_to_model* has recently been described in detail (Terwilliger et al.,
198 2018c). The inputs required are a map file (CCP4/MRC format, Cheng et al., 2015), a sequence
199 file with the sequences of residues or nucleotides in each unique chain in the structure, and the
200 nominal resolution of the map. If symmetry was used in the reconstruction process, then the
201 symmetry operators can be supplied as well. All other parameters are fully optional and it is
202 normally not necessary for a user to adjust them.

203
204 For the model-building described here, the maps, sequence files, symmetry operators, and
205 resolution were all obtained from the 2016 Model Challenge web site at
206 http://challenges.emdatabank.org/?q=model_challenge.

207
208 The first step carried out by the *map_to_model* tool is to automatically sharpen the map with
209 the *phenix.auto_sharpen* tool (Terwilliger et al., 2018a). In this approach the map is sharpened
210 (or blurred) to attempt to simultaneously maximize the level of detail in the map and the
211 connectivity of the map.

212
213 The second step is to carry out automatic map segmentation as described above, yielding one
214 map that represents the unique part of the sharpened map along with a set of small maps each
215 representing one small region of connected density (all above a contour level determined
216 automatically during the segmentation process).

217
218 The third step is to carry out automatic model-building for each chain type that is represented
219 in the sequence file. This is done for the map representing the unique part of the sharpened

220 map and for each small map. Model-building is done using tools available in *Phenix* that include
221 placement of helices and strands in density of corresponding shapes (Terwilliger, 2010a;
222 Terwilliger, 2010b), tracing density along a chain and replacement with main-chain atoms
223 (Terwilliger, 2010c), placement of short fragments by convolution-based searches followed by
224 extension with 3-residue fragments from structures in the PDB (Terwilliger, 2003; Terwilliger et
225 al., 2018c), and recently-described methods for model-building of RNA that are extensions of
226 these procedures for protein (Terwilliger et al., 2018c).

227
228 The fourth step is to combine all the models. The principal method for combining models is to
229 rank all segments (fragments of a model that have no chain breaks) based on map-model
230 correlation, segment length, and secondary structure, then to go through this ranked list and
231 place whatever part of each model does not overlap with a higher-scoring model (Terwilliger et
232 al., 2018c).

233
234 After each model is built, after models are combined, and after application of reconstruction
235 symmetry to the final model, each working model is refined with real-space refinement
236 (Afonine et al., 2018a).

237

238

239

240 Data used from the Cryo-EM Model Challenge

241

242 The maps and reconstruction symmetry used for the 12 cryo-EM maps in the 2016 Cryo-EM
243 Model Challenge were taken from the Model Challenge site at
244 http://challenges.emdatabank.org/?q=model_challenge (accessed 2017-11-19). The Challenge
245 consisted of 8 unique molecules, four of which were associated with two maps at different
246 resolutions, leading to 12 different maps (Table I). Of these maps, most were associated with
247 previously-deposited models that were likely to be more accurate than the ones we built
248 automatically and that were therefore suitable for use as references for the accuracy of our
249 models. For one additional map (groEL, EMD entry 6422) there was no deposited model,
250 however there is a model for a related structure in the PDB (1ss8) which we offset
251 superimposed on this map and used as a reference. One final structure was recently
252 interpreted (the proteasome structure; Veessler, D., unpublished) and we used that structure as
253 a reference model. We checked the map-model agreement with *phenix.map_model_cc* and
254 these map-model correlations ranged from 0.34 (rather low, supporting only low confidence in
255 the model), to 0.85, suggesting that the model is in good agreement with the map.

256

257 Results and Discussion

258

259 Fig. 1 illustrates the application of our map segmentation procedure (Terwilliger et al., 2018b)
260 to the 2.9 Å cryo-EM reconstruction of the anthrax protective antigen pore (Jiang et al., 2015).
261 The map has C7 symmetry (a 7-fold symmetry axis). Fig. 1A shows the 7-fold symmetry of the
262 pore and illustrates one of the 7 chains in purple. Fig. 1B shows the density map with 7-fold
263 symmetry. It can be seen that the density is much stronger for the extracellular region of the

264 molecule than for the transmembrane part below. The 7-fold symmetry was used along with
265 the map to identify symmetry-related regions of density in the map. Then a compact and
266 connected unique set of density regions was chosen to represent the molecule. Fig. 1C shows
267 the individual segmented regions of the map, and Fig. 1D shows the segmented region,
268 augmented by neighboring regions of density.

269

270 We then applied our segmentation procedure to a large set of cryo-EM maps from the EMDB
271 (Fig. 2). As expected, using the reconstruction symmetry of the maps in segmentation often
272 resulted in a very large reduction in the volume that needed to be considered to include the
273 unique part of each map (Fig. 2A). The average volume after segmentation and placing the
274 unique segmented region in a new box was 8% of the starting volume of the maps. In most
275 (206 of 218) of the cases illustrated in Fig. 2 we used the *add_neighbors* keyword to add a layer
276 of regions around the unique molecular volume in order to increase the chance of finding a
277 complete molecule. The 12 cases (EMD_2807, EMD_3137, EMD_5185, EMD_5600, EMD_6346,
278 EMD_6630, EMD_6637, EMD_6688, EMD_8598, EMD_8605, EMD_8644, EMD_9518) where
279 this was not done are those where the map was large (maps with 16M to 134M elements) and
280 we attempted to keep the size of the region to be worked on to the minimum possible.

281

282 Fig. 2B illustrates the fraction of the unique molecular unit that is within the unique segmented
283 region used for each map in Fig. 2A. This fraction of the molecule contained within the
284 segmented region is estimated from the map-model correlation between the molecule and a
285 map which is set to zero everywhere outside the segmented region, normalized to the map-
286 model correlation without setting any of the map to zero. If the molecule is within the
287 segmented region this normalized correlation will be unity, while if the molecule is split
288 between different segmented regions it will be smaller. As shown in Fig 2B, the fraction within
289 a single segmented region varies considerably among the 218 maps analyzed here, but the
290 mean fraction was 0.72, indicating that typically a large fraction, but not all, of the molecular
291 unit was contained within the segmented region.

292

293 We examined whether the fraction of the molecular unit contained within the segmented
294 region (Fig. 2B) depended on the number of symmetry copies or the resolution of the map. The
295 number of symmetry copies had only a small effect: maps with a single copy had an average
296 fraction of 0.73 and maps with 60 copies had an average of 0.71. On the other hand, resolution
297 had quite a substantial impact on the fraction within the segmented region: maps with
298 resolution of 3.5 Å or better had a mean fraction of 0.82; maps with resolution of 4 Å or worse
299 had a mean of 0.63.

300

301 We applied the combination of map sharpening, segmentation, and model-building as
302 implemented in the *Phenix* tool *phenix.map_to_model* (Terwilliger et al., 2018b) to the 12 maps
303 in the 2016 cryo-EM Challenge. The maps and corresponding reference models are listed in
304 Table I along with the CPU hours required for the analysis, which ranged from 7 to 422 hours.
305 Table II lists the number of residues that were built with C_α or P atoms within 3 Å of the
306 corresponding atoms in the reference model by the *phenix.map_to_model* procedure, along
307 with the fraction of the reference model represented by the model that was built and the

308 fraction of residues that were assigned the correct residue identity. The number of residues
309 built more than 3 Å from any residue in the reference model is also listed.

310
311 Overall, from 35% to 82% of the protein portions of the 12 structures were built within 3 Å of
312 the corresponding reference models. For the two RNA structures, 24% and 54% of the RNA
313 portions were built within 3 Å of the corresponding reference models. From 8% to 75% of the
314 protein and RNA sequences were correctly assigned. For the non-ribosome structures, only a
315 small proportion of the models built did not correspond at all to the deposited models. On the
316 other hand, for the ribosome structures, a large fraction (over half for the 3.6 Å map) of the
317 protein residues built did not correspond to the deposited models. Most of these incorrectly-
318 built residues are located in regions that are RNA in the deposited models (recently we have
319 developed a tool, *phenix.remove_poor_fragments* that can remove some of these incorrectly-
320 built residues, but it was not available at the time of this work, TT, OS, PDA and PVA,
321 unpublished). The models built by *phenix.map_to_model* have RMS distances for C α /P atoms
322 from reference interpretations of 0.6 Å to 2.1 Å.

323
324
325 The procedures developed here for map segmentation could be applied automatically to all of
326 the 218 maps that we examined in the tests shown in Fig. 2. Further, all 12 of the maps in the
327 2016 Cryo-EM Model Challenge could be automatically sharpened, segmented and partially
328 interpreted by the *phenix.map_to_model* procedure. It seems likely that combining the
329 techniques developed here with other approaches for automatic model-building might lead to
330 procedures that can automatically interpret an even larger part of cryo-EM maps.

331
332
333 Acknowledgements

334
335 This work was supported by the NIH (grant GM063210 to PDA and TT) and the *Phenix* Industrial
336 Consortium. This work was supported in part by the US Department of Energy under Contract
337 No. DE-AC02-05CH11231 at Lawrence Berkeley National Laboratory. This research used
338 resources provided by the Los Alamos National Laboratory Institutional Computing Program,
339 which is supported by the U.S. Department of Energy National Nuclear Security Administration
340 under Contract No. DE-AC52-06NA25396. Molecular graphics and analyses were performed with
341 the UCSF *Chimera* package (Pettersen et al., 2005) and with *Coot* (Emsley et al., 2010).

342
343 Figure Legends

344
345 Figure 1. Segmentation of density for the anthrax protective antigen pore. A. Deposited
346 structure of anthrax protective antigen pore with one of the 7 chains in purple. B. Density map
347 illustrating the 7-fold symmetry used in the reconstruction. C. Individual segmented regions of
348 the map superimposed on a single chain from the deposited structure. Note that the deposited
349 structure was not used in the segmentation process. D. Illustration of the segmented region,
350 augmented by neighboring regions of density.

351

352 Figure 2. Histograms showing the results of application of the segmentation procedure to cryo-
353 EM maps from the EMDB. Datasets are grouped according to (panel A), the fraction of original
354 map required to represent the segmented region of each map, or (panel B), the fraction of each
355 molecular unit contained within the segmented region of each map. In each panel, the label
356 corresponds to the lower bound of each grouping. The values are grouped in increments of
357 0.05, so for example the number of datasets with values from 0.00 to 0.05 is shown over the
358 ordinate of "0".

359
360

361

362 References

363

364 Adams, P. D., Afonine, P.V., Bunkóczy, G., Chen, V.B., Davis, I.W., Echols, N., Headd, J.J., Hung, L.-
365 W. Kapral, G.J., Grosse-Kunstleve, R.W., McCoy, A.J., Moriarty, N.W., Oeffner, R., Read, R.J.,
366 Richardson, D.C., Richardson, J.S., Terwilliger, T.C. & Zwart, P.H., 2010. PHENIX: a
367 comprehensive Python-based system for macromolecular structure solution. *Acta Cryst.*, D66,
368 213-221.

369

370 Afonine, P.V., Poon, B.K., Read, R.J., Sobolev, O.V., Terwilliger, T.C., Urzhumtsev, A., Adams, P.D.
371 (2018a). Real-space refinement in *Phenix* for cryo-EM and crystallography. *Acta Cryst D.*, in
372 press.

373

374 Afonine, P.V., Klaholz, B.P., Moriarty, N.W., Poon, B.K., Sobolev, O.V., Terwilliger, T.C., Adams,
375 P.D., Urzhumtsev, A. (2018b). New tools for the analysis and validation of Cryo-EM maps and
376 atomic models. *BioRxiv* doi.org/10.1101/279844.

377

378 Baker, M.L., Yu, Z., Chiu, W., Bajaj, C., 2006. Automated segmentation of molecular subunits in
379 electron cryomicroscopy density maps. *J. Structural Biology* 156, 432-441.

380

381 Baker, M.L., Ju, T., Chiu, W., 2007. Identification of secondary structure elements in
382 intermediate-resolution density maps. *Structure* 15, 7–19.

383

384 Baker, M.L., Abeyasinghe, S.S., Schuh, S., Coleman, R.A., Abrams, A., March, M.P., Hryc, C.F.,
385 Ruths, T., Chiu, W., Ju, T. (2011). Modeling protein structure at near atomic resolutions with
386 Gorgon. *J. Struct. Biol.* 174, 360-373.

387

388 Bai, X.C., Yan, C.Y., Yang, G.H., Lu, P.L., Ma, D., Sun, L.F., Zhou, R., Scheres, S.H.W., Shi, Y.G.
389 (2015). An atomic structure of human γ -secretase. *Nature* 525, 212-217.

390

391 Bartesaghi, A., Merk, A., Banerjee, S., Mathies, D., Wu, X, Milne, J.L., Subramaniam, S. (2015).
392 2.2 Å resolution cryo-EM structure of β -galactosidase in complex with a cell-permeant inhibitor.
393 *Science* 348, 1147-1151.

394

395 Bartesaghi, A., Matthies, D., Banerjee, S., Merk, A., Subramaniam, S. Structure of β -
396 galactosidase at 3.2- Å resolution obtained by cryo-electron microscopy. Proc. Natl. Acad. Sci.
397 USA 111, 11709-11714.
398
399 Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N. and
400 Bourne, P.E., 2000. The Protein Data Bank. Nucleic Acids Research, 28, 235-242.
401
402 Bernstein, F.C., Koetzle, T.F., Williams, G.J., Meyer, E.F., Jr., Brice, M.D., Rodgers, J.R., Kennard,
403 O., Shimanouchi, T. and Tasumi, M., 1977. The Protein Data Bank: a computer-based archival
404 file for macromolecular structures. J. Mol. Biol. 112, 535-542.
405
406 Campbell, M.G., Veessler, D., Cheng, A., Potter, C.S., Carragher, B. (2015). 2.8 Angstrom
407 resolution reconstruction of the Thermoplasma acidophilum 20S proteasome using cryo-
408 electron microscopy. eLife 4, e06380-e06380.
409
410 Cheng, A., Henderson R., Mastrorarde, D., Ludtke, S., Schoenmakers, R.H.M., Short, J.,
411 Marabini, R., Dallakyan, S., Agard, D., Winn, M. (2015). J. Struct. Biol. 192, 146-150.
412
413 C. Colliex, C., Cowley, J. M., Dudarev, S. L., Fink, M., Gjønnnes, J., Hilderbrandt, R., Howie, A.,
414 Lynch, D. F. , Peng, L. M. , Ren, G., Ross, A. W. , Smith, V. H., Jr, Spence, J. C. H. , Steeds, J. W.,
415 Wang, J., Whelan, M. J., Zvyagin, B. B. International Tables for Crystallography (2006). Vol. C, ch.
416 4.3, pp. 259-429.
417
418 Collins, P., Si, D. (2017). A graph based method for the prediction of backbone trace from cryo-
419 EM density maps. ACM-BCB '17 Proceedings of the 8th ACM International Conference on
420 Bioinformatics, Computational Biology, and Health Informatics.
421
422 Chen, M., Baldwin, P.R., Ludtke, S.J., Baker, M.L., 2016. De Novo modeling in cryo-EM density
423 maps with Pathwalking. J. Structural Biol. 196, 289-298.
424
425 DiMaio, F., Chiu, W. (2016). Tools for model building and optimization into near-atomic
426 resolution electron cryo-microscopy density maps. Methods. Enzymol. 679, 255-276.
427
428 Emsley, P., Lohkamp, B., Scott, W. G. & Cowtan, K. (2010). Features and development of Coot.
429 Acta Cryst. D66, 486-501.
430
431 Fischer, N., Neumann, P., Konevega, A.L., Bock, L.V., Ficner, R., Rodnina, M.V., Stark, H. (2015).
432 Structure of the E. coli ribosome-EF-Tu complex at <3 Å resolution by Cs-corrected cryo-EM.
433 Nature 520, 567-570.
434
435 Frenz, B., Walls, A.C., Egelman, E.H., Veessler, D., DiMaio, F., 2017. RosettaES: a sampling
436 strategy enabling automated interpretation of difficult cryo-EM maps. Nature Methods 14, 797-
437 803.
438

- 439 Fromm, S.A., Bharat, T.A., Jakobi, A.J., Hagen, W.J., Sachse, C. (2015). Seeing tobacco mosaic
440 virus through direct electron detectors. *J. Struct. Biol.* 189, 87-97.
441
- 442 Jiang, W., Baker, M.L., Ludtke, S.J., Chiu, W., 2001. Bridging the information gap: computational
443 tools for intermediate resolution structure interpretation. *J. Mol. Biol.* 308, 1033–1044.
444
- 445 Jiang J., Pentelute B.L., Collier R.J., Zhou, Z.H., 2015. Atomic structure of anthrax protective
446 antigen pore elucidates toxin translocation. *Nature* 521, 545-549.
447
- 448 Kong, Y., Ma, J., 2003. A structural-informatics approach for mining beta-sheets: locating sheets
449 in intermediate-resolution density maps. *J. Mol. Biol.* 332, 399–413.
450
- 451 Kong, Y., Zhang, X., Baker, T.S., Ma, J., 2004. A Structural-informatics approach for tracing beta-
452 sheets: building pseudo-C(alpha) traces for beta-strands in intermediate-resolution density
453 maps. *J. Mol. Biol.* 339, 117–130.
454
- 455 Lawson C.L., Patwardhan, A., Baker, M.L., Hryc, C., Garcia, E.S., Hudson, B.P., Lagerstedt, I.,
456 Ludtke, S.J., Pintilie, G., Sala, R., Westbrook, J.D., Berman, H.M., Kleywegt, G.J., Chiu, W. 2016.
457 EMDatabank unified data resource for 3DEM. *Nucleic Acids Res.* 44, D396–D403.
458
- 459 Li, X., Mooney P, Zheng, S., Booth, C., Braunfeld ,M.B., Gubbens, S., Agard, DA., Cheng, Y.
460 (2013). Electron counting and beam-induced motion correction enable near-atomic-resolution
461 single-particle cryoEM. *Nat. Methods* 10, 584-590.
462
- 463 Li, W., Liu, Z., Koripella, R.K., Langlois, R., Sanyal, S., Frank, J. (2015). Activation of GTP hydrolysis
464 in mRNA-tRNA translocation by elongation factor G. *Sci. Adv.* 1, e1500169.
465
466
- 467 Liao, M., Cao, E., Julius, D., Cheng, Y. (2013). Structure of the TRPV1 ion channel determined by
468 electron cryo-microscopy *Nature* 504, 107-112.
469
- 470 Lindert, S., Alenxander, N., Wotzel, N., Karakas, M., Stewart, P.L., Meiler, J., 2012. Em-Fold: De
471 novo atomic-detail protein structure determination from medium-resolution density maps.
472 *Structure* 20, 464-478.
473
- 474 Lu ,P.L., Bai, X.C., Ma, D., Xie, T., Yan, C.Y., Sun, L.F., Yang, G.H., Zhao, Y.Y., Zhou, R., Scheres,
475 S.H.W., Shi, Y.G. (2014). Three-dimensional structure of human gamma-secretase. *Nature* 512,
476 166-170.
477
478
- 479 Pettersen E.F., Goddard T.D., Huang C.C., Couch G.S., Greenblatt D.M., Meng E.C., Ferrin, T.E.
480 2004. UCSF Chimera--a visualization system for exploratory research and analysis. *J Comput.*
481 *Chem.* 25, 1605-1612.
482

- 483 Pintilie, G. D., Zhang, J., Goddard, T. D., Chiu, W., & Gossard, D. C., 2010. Quantitative analysis
484 of cryo-EM density map segmentation by watershed and scale-space filtering, and fitting of
485 structures by alignment to regions. *Journal of Structural Biology* 170, 427–438.
486
- 487 Roh, S.H., Hryc, C.F., Jeong, H.H., Fei, X., Jakana, J., Lorimer, G.H., Chiu, W. (2017). Subunit
488 conformational variation within individual GroEL oligomers resolved by Cryo-EM. *Proc. Natl.*
489 *Acad. Sci. U.S.A.* 114, 8259-8264.
490
- 491 Terwilliger, T. C. (2003). [Automated main-chain model-building by template-matching and](#)
492 [iterative fragment extension](#). *Acta Cryst.* D59, 38-44.
493
- 494 Terwilliger, T. C. (2010a). Rapid model-building of β -sheets in electron density maps. [Acta Cryst.](#)
495 [D66, 276-284](#).
496 Terwilliger, T. C. (2010b). Rapid model-building of α -helices in electron density maps. [Acta](#)
497 [Cryst. D66, 268-275](#).
498
- 499 Terwilliger, T. C. (2010c). Rapid chain-tracing of polypeptide backbones in electron density
500 maps. [Acta Cryst. D66, 285-294](#).
501
- 502 Terwilliger, T.C., Sobolev, O., Afonine, P.V., Adams, P.D. (2018a). Automated map sharpening by
503 maximization of detail and connectivity. bioRxiv doi: <https://doi.org/10.1101/247049>.
504
- 505 Terwilliger, T.C., Adams, P.D., Afonine, P.V., Sobolev, O.V. (2018b). A fully automatic method
506 yielding initial models from high-resolution electron cryo-microscopy maps. bioRxiv doi:
507 <https://doi.org/10.1101/267138>
508
- 509 Volkmann, N. (2002). A novel three-dimensional variant of the watershed transform for
510 segmentation of electron density maps. *J. Struct. Biol.* 138, 123–129.
511
- 512 Wang, Z., Hryc, C.F., Bammes, B., Afonine, P.V., Jakana, J., Chen, D.H., Liu, X, Baker, M.L., Kao,
513 C., Ludtke, S.J., Schmid, M.F., Adams, P.D., Chiu W. (2014). An atomic model of brome mosaic
514 virus using direct electron detection and real-space optimization. *Nature Commun.* 5 4808.
515
- 516 Wang, R. Y-R., Kudryashev, M., Li, X., Egelman, E.H., Basler, M., Cheng, Y., Baker, D., DiMaio, F.,
517 2015. De novo protein structure determination from near-atomic-resolution cryo-em maps.
518 *Nature Methods* 12, 335-341.
519
- 520
- 521 Yu, Z., Bajaj, C. (2008). Computational approaches for automatic structural analysis of large
522 biomolecular complexes. *IEEE/AC Transactions on computational biology and bioinformatics* 5
523 568-582.
524
- 525 Zhang, Q., Bettadapura, R., Bajaj, C. (2012). Macromolecular structure modeling from 3D EM
526 using VolRover 2.0. *Biopolymers* 97, 709-731.

527

528 Zhou N., Wang, H., Wang, J., 2017. EMBuilder: A Template Matching- based automatic model-
529 building program for high-resolution cryo-electron microscopy maps. Scientific Rep. 7, 2664.

530

531

532 Table I. Cryo-EM Model Challenge structures analyzed with *phenix.map_to_model*

Structure	Map	Reference model	Resolution	Reference map-model correlation	Cpu (h)	Reference
	(EMDB entry)	(PDB entry ¹)	(Å)			
Beta galactosidase	2984	5A1A	2.2	0.73	32	Bartesaghi et al. (2015)
Proteasome	6287	undeposited ²	2.8	0.81	16	Campbell et al. (2015)
<i>E. coli</i> 70S ribosome	2847	5AFi	2.9	0.85	422	Fischer et al. (2015)
Beta galactosidase	5995	3J7H	3.2	0.76	39	Bartesaghi et al. (2014)
Proteasome	5623	3J9i	3.3	0.77	14	Li et al. (2013)
trpV1	5778	3J5P	3.3	0.56	14	Liao et al. (2013)
TMV	2842	4UDV	3.4	0.73	7	Fromm et al. (2015)
Gamma secretase	3061	5A63	3.4	0.41	34	Bai et al. (2015)
<i>E. coli</i> 70S ribosome	6316	3JA1	3.6	0.38	322	Li et al. (2015)
Brome mosaic virus	6000	3J7L	3.8	0.76	16	Wang et al. (2014)
GroEL	6422	1SS8	4.1	0.83	153	Unpublished data ³
Gamma secretase	2677	5A63	4.5	0.34	23	Lu et al. (2014)

¹ The PDB codes are written following the convention outlined in the editor's notes in the Computation Crystallography Newsletter (Comput. Cryst. Newsl. 2015:6; https://www.phenix-online.org/newsletter/CCN_2015_07.pdf).

² The recently-determined proteasome structure (Veesler, D., unpublished) was used as a reference model.

³ The PDB entry 1ss8 was used as a model for entry EMD_6422, as used in the related entry EMD_8750 (Roh et al., 2017)

533 Table II. Results of Cryo-EM Challenge analysis with *phenix.map_to_model*
534

Structure	Resolution (Å)		Residues in reference model	Built within 3 Å (Residues)	RMS distance (Å) ¹	Built (%)	Matched to sequence ² (%)	Built further than 3 Å (Residues)
Beta galactosidase	2.2		1022	842	0.6	82	75	10
Proteasome	2.8		422	327	0.6	78	59	3
<i>E. coli</i> 70S ribosome	2.9	Protein	6322	3212	0.7	51	23	2112
		RNA	4748	2566	1.0	54	49	280
Beta galactosidase	3.2		1022	676	1.2	66	23	26
Proteasome	3.3		427	246	1.2	58	48	4
trpV1	3.3		592	330	0.8	56	30	18
TMV	3.4		153	76	1.3	50	34	8
Gamma secretase	3.4		1223	832	1.0	68	24	80
<i>E. coli</i> 70S ribosome	3.6	Protein	7125	2479	2.0	35	8	3979
		RNA	4685	1140	1.9	24	26	440
Brome mosaic virus	3.8		479	198	1.5	41	11	14
GroEL	4.1		524	309	1.3	59	16	10
Gamma secretase	4.5		1223	619	2.1	51	8	185

¹ The RMS distances for C α /P atoms from the reference interpretations, only including the residues built within 3 Å of the reference model.

² The percentage of matched to sequence is the total number of residues in the automatically-built model correctly matched to sequence divided by the total number of residues in the reference model.

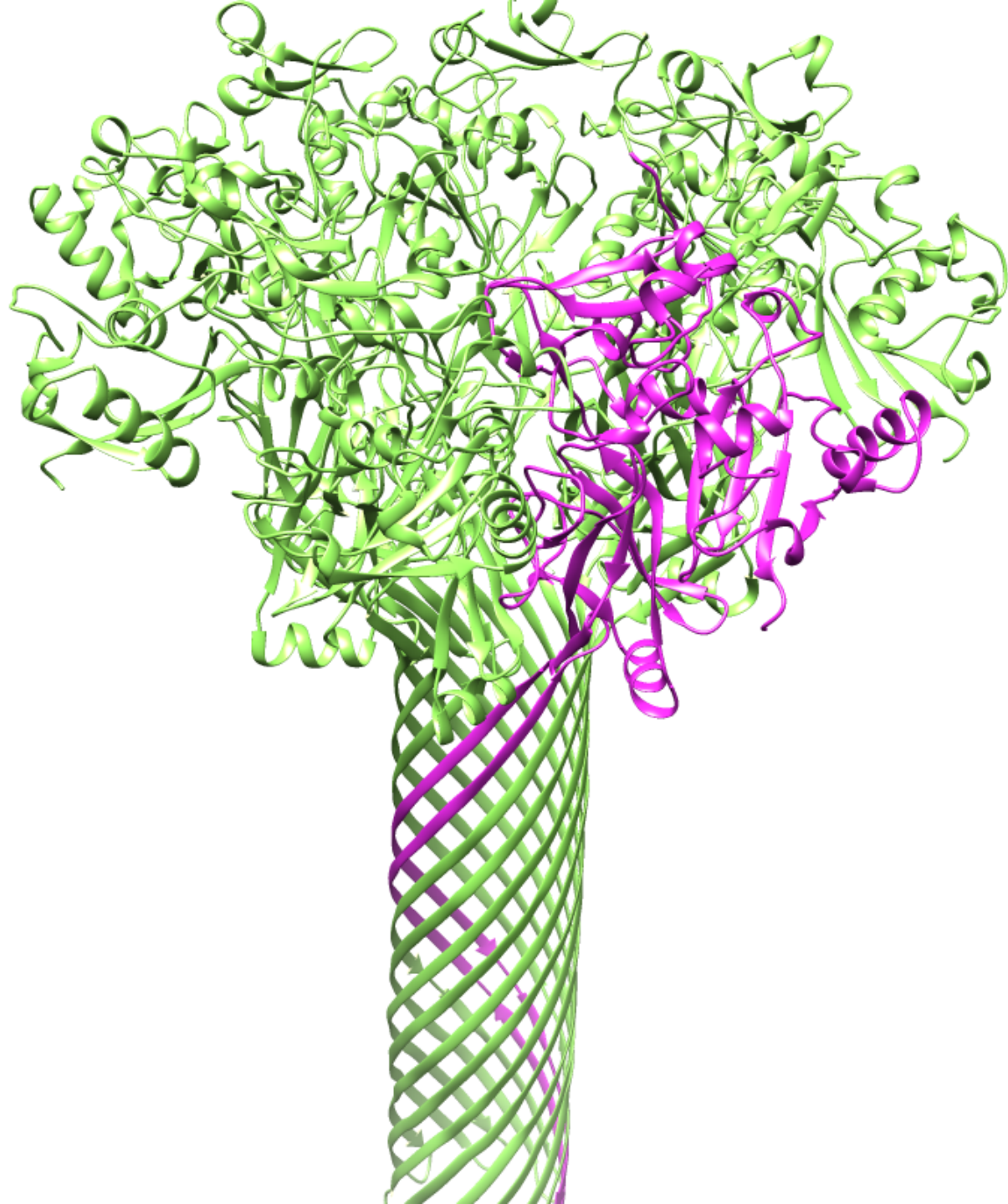


Fig 1A

Fig 1B

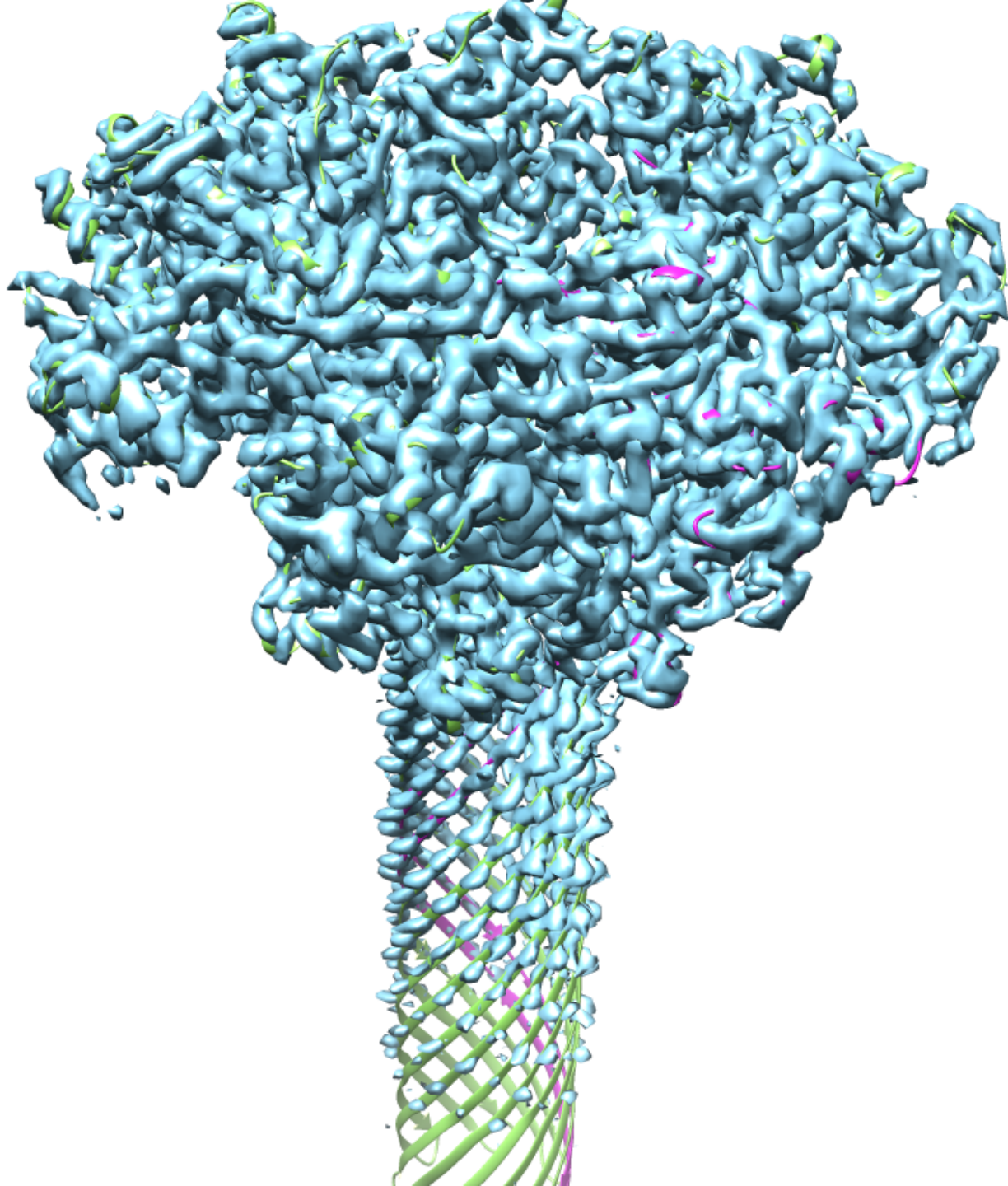


Fig 1C

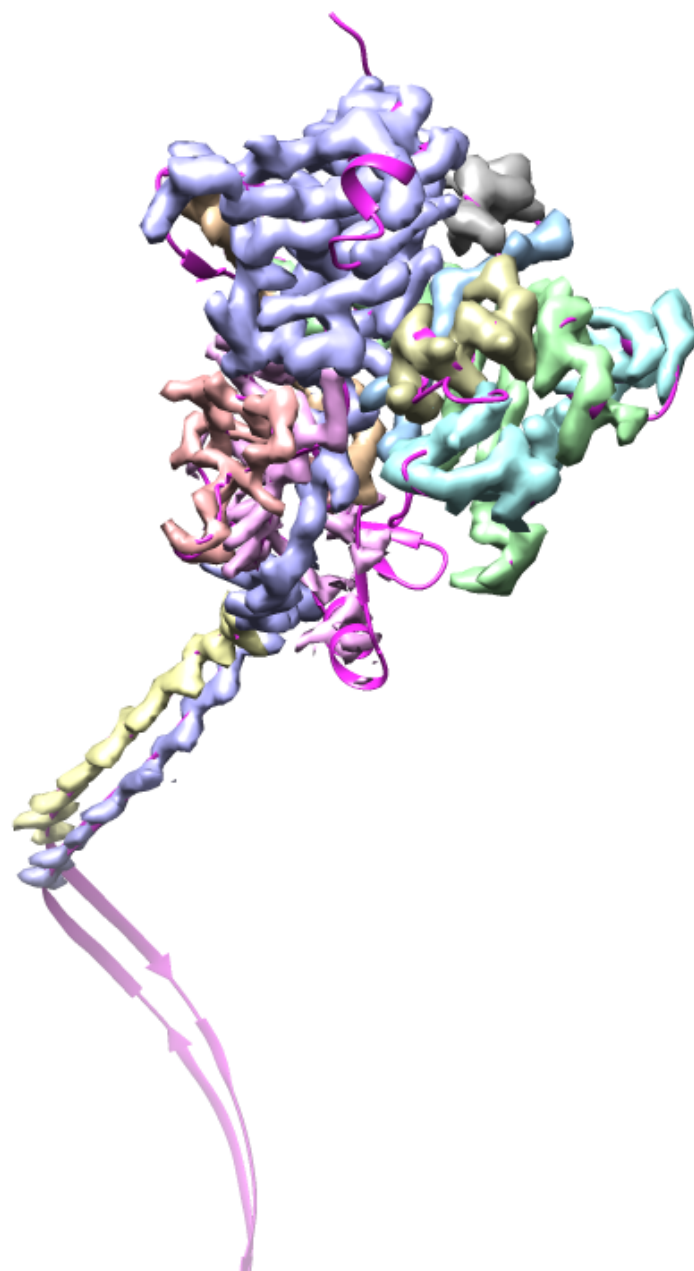


Fig 1D

