

1 cryoSPARC: Algorithms for rapid unsupervised cryo-EM structure determination

2
3 Ali Punjani¹, John L. Rubinstein^{2,3,4}, David J. Fleet⁵, and Marcus A. Brubaker⁶

4
5 1. Department of Computer Science, The University of Toronto, Toronto, Ontario, Canada

6 2. Molecular Structure and Function Program, The Hospital for Sick Children Research Institute,
7 Toronto, Ontario, Canada

8 3. Department of Biochemistry, The University of Toronto, Toronto, Ontario, Canada

9 4. Department of Medical Biophysics, The University of Toronto, Toronto, Ontario, Canada

10 5. Department of Computer Science, The University of Toronto, Toronto, Ontario, Canada

11 6. Department of Electrical Engineering and Computer Science, York University, Toronto,
12 Ontario, Canada

13
14 Correspondence to:

15 Ali Punjani (alipunjani@cs.toronto.edu) and Marcus Brubaker (mab@eecs.yorku.ca)

16 17 18 Abstract

19
20 Single particle electron cryomicroscopy (cryo-EM) is a powerful method for determining
21 the structures of biological macromolecules. With automated microscopes, cryo-EM data
22 can often be obtained in a few days. However, processing cryo-EM image data to reveal
23 heterogeneity in the protein structure and to refine 3-D maps to high resolution frequently
24 becomes a severe bottleneck, requiring expert intervention, prior structural knowledge,
25 and weeks of calculations on expensive computer clusters. Here we show that stochastic
26 gradient descent (SGD) and branch and bound maximum likelihood optimization
27 algorithms permit the major steps in cryo-EM structure determination to be performed in
28 hours or minutes on an inexpensive desktop computer. Furthermore, SGD with Bayesian
29 marginalization allows *ab initio* 3-D classification, enabling automated analysis and
30 discovery of unexpected structures without bias from a reference map. These algorithms
31 are combined in a user-friendly computer program named cryoSPARC.

32 Introduction

33

34 Scientific approaches can be transformed by innovations that decrease the cost and improve non-
35 expert usability of technology, as seen with DNA sequencing and synthesis, microarray
36 technology, and even the use of computers themselves. These changes can occur both
37 quantitatively, by allowing more experiments to be done in a shorter time by experts and non-
38 specialists, and qualitatively by changing the type and scope of experiments that are feasible.
39 Recent advances in single particle cryo-EM^{1,2} have enabled near-atomic resolution structure
40 determination of biomedically important protein complexes³⁻⁵, bringing the technique to the
41 attention of the general biological research community and pharmaceutical companies. The
42 throughput and automation of cryo-EM becomes increasingly important as the technique is used
43 for structure-based drug design⁶ and time-critical structural studies of pathogens⁷. Automated
44 electron microscopes can collect datasets for atomic resolution structure determination in as little
45 as 24 or 48 hours given appropriately prepared specimens, and centralized cryo-EM facilities are
46 now providing instrument access to non-specialist investigators. Calculation of 3-D maps from
47 cryo-EM images, however, can require weeks of computational analysis by an expert user. With
48 routine collection of cryo-EM datasets that contain millions of single particle images
49 corresponding to different 3-D conformations of the sample⁸, the cost of image analysis can
50 exceed 500,000 CPU hours on large, expensive computer clusters⁹. Further, without significant
51 user expertise, there are a variety of ways in which incorrect and misleading 3-D maps can be
52 generated at various stages in the image analysis pipeline^{10,11}. The computational cost and the
53 requirement for user input are bottlenecks both for automation and widespread use of cryo-EM.

54

55 To address these issues, we developed two new algorithms. The first of these algorithms, for the
56 first time, makes it possible to perform unsupervised *ab initio* 3-D classification, whereby
57 multiple 3-D states of a protein can be discovered from a single sample without user input of
58 prior structural knowledge, and without the assumption that all 3-D states resemble each other. In
59 contrast, existing techniques for 3-D refinement of cryo-EM maps require an initial structure that
60 is close to the correct target structure^{12,13}. The second algorithmic development radically speeds
61 up high-resolution refinement of cryo-EM maps by exploiting characteristics of image alignment
62 to achieve massive computational savings by removing redundant computation. These two
63 abilities are combined in a standalone Graphics Processing Unit (GPU) accelerated software
64 package that we have named cryoSPARC (cryo-EM single particle *ab initio* reconstruction and
65 classification). CryoSPARC can refine multiple high-resolution 3D structures directly from
66 single particle images, with no user input or expertise required. These combined steps are done
67 in a matter of hours on a single consumer-grade desktop computer. GPU hardware has been used
68 previously to accelerate cryo-EM contrast transfer function estimation¹⁴ and identification of
69 particles within images¹⁵. Related work has shown that exploiting GPU hardware in the popular
70 program RELION can significantly speed up existing algorithms for reference-based 3-D
71 classification and refinement⁹. The algorithms presented here provide a further order-of-
72 magnitude reduction in computational cost compared to GPU acceleration, which would require
73 at least an additional ~7 years if driven by hardware advances alone¹⁶. Based on the combination
74 of algorithms, inexpensive hardware, and an easy-to-use graphical user interface, cryoSPARC
75 can allow new non-specialist cryo-EM users to process data rapidly without needing to purchase
76 or set up their own computer clusters and with minimal user input and expertise.

77

78 Results

79

80 Formally, structure determination by cryo-EM is an optimization problem and may be described
81 in a Bayesian likelihood framework^{12,17}:

$$\arg \max_{V_{1...K}} \log p(V_{1...K} | X_{1...N}) = \arg \max_{V_{1...K}} \sum_{i=1}^N \log \sum_{j=1}^K \frac{1}{K} \int p(x_i, \phi_i | V_j) d\phi_i + \log p(V_{1...K}) \quad (1)$$

82

83 The aim of the optimization is to find the 3-D structures (V_1 to V_K) that best explain the observed
84 images (X_1 to X_N), by marginalizing over class assignment (j) and the unknown pose variable
85 (ϕ_i), which describes a 3-D rotation and a 2-D translation for each particle image.

86 Numerical optimization problems have been studied extensively in computer science¹⁸.
87 Traditionally, optimization is formulated as the maximization of a single, monolithic objective
88 function. With this approach, the variables of a function are iteratively altered until the ‘best’
89 values, which give an optimum value to the function, are identified. Sophisticated algorithms for
90 iterative optimization have been developed¹⁹ and are central to a myriad of problems in data
91 modeling and engineering. In the case of cryo-EM map calculation, the objective function
92 (Equation 1) quantifies how well cryo-EM maps explain the collected experimental images, and
93 the variables in the function include the 3-D maps represented as density voxels on a 3-D grid.

94 We use a stochastic gradient descent (SGD) optimization scheme to quickly identify one or
95 several low-resolution 3-D structures that are consistent with a set of observed images. This
96 algorithm allows for *ab initio* heterogeneous structure determination with no prior model of the
97 molecule’s structure. Once approximate structures are determined, a branch and bound algorithm
98 for image alignment helps rapidly refine structures to high resolution. The speed and robustness
99 of these approaches allow structure determination in a matter of minutes or hours on a single
100 inexpensive desktop workstation.

101

102 **Stochastic Gradient Descent: Discovery of protein structure from random initialization**

103

104 Cryo-EM map calculation is a non-convex optimization problem. These problems are among the
105 most computationally challenging optimization problems known and are characterized by the
106 presence of multiple locally-optimal settings of variables, each of which forms an attractor where
107 typical iterative optimization algorithms can become stuck if poorly initialized¹⁹ (Figure 1A).
108 Sensitivity to local optima is seen in most optimization algorithms, including those used in cryo-
109 EM^{12,13} and as a result, refinement programs require a reasonably accurate initial model for the
110 structure that initializes the search near the global optimum. However, recent methods have been
111 discovered that perform well on non-convex problems. One such method is stochastic gradient
112 descent (SGD)²⁰ (Figure 1). SGD was popularized as a key tool in Deep Learning for the
113 optimization of non-convex functions, resulting in near-human level performance in tasks like
114 image and speech recognition^{21,22}.

115

116 In Equation 1, each term of the outer sum represents the contribution of a single particle image to
117 the overall likelihood of the 3-D map. SGD repeatedly approximates this objective function by
118 selecting a different random subset of terms (i.e., single particle images) at each iteration, and

119 computes the sum of those terms (Figure 1B). In a single iteration, the optimization variables
120 (i.e., the 3-D map) are updated based on the gradient of this approximate objective
121 (Supplementary Note 1). Each iteration requires analyzing only a small subset of single particle
122 images. As a consequence, a single iteration is inexpensive and hundreds or thousands of
123 iterative changes can be made during each pass through the full dataset. It is commonly believed
124 that it is because of these many noisy changes that SGD is insensitive to local optima and often
125 finds effective solutions to non-convex problems (Figure 1C).

126
127 We implemented an SGD method for *ab initio* structure determination and 3-D classification.
128 Applied to several different datasets, the use of SGD enables convergence to correct low-
129 resolution structures from arbitrary random initialization, allowing both *ab initio* structure
130 determination and *ab initio* 3-D classification (Figure 2). With 35,645 TRPV1 particle images³
131 SGD optimization resulted in a low-resolution 3-D map in 75 minutes from random initialization
132 (Figure 2A) using a single inexpensive desktop workstation with an Intel i7-5820K Processor
133 and a single NVIDIA Tesla K40 GPU. When applied to a dataset of conformationally
134 heterogeneous *Thermus thermophilus* V/A-ATPase particle images²³, the algorithm was able to
135 discern three different conformational states for the enzyme, again from random initializations
136 (Figure 2B). These three states correspond to the three different rotational positions of the central
137 rotor of the enzyme²⁴. This finding is particularly notable as previous analysis with reference-
138 based classification¹² and the same dataset of images was only able to detect two of the three
139 states²³. The newly identified third rotational state is the conformation of the enzyme that differs
140 the most from the other two. This observation illustrates the importance of reference-free *ab*
141 *initio* classification for unbiased identification of states that differ from the expected structures
142 present in the dataset.

143

144 **Branch and bound: rapid refinement of maps to high resolution**

145

146 The primary computational burden in map refinement is the search for orientation parameters
147 that best align each 2-D single particle image to a 3-D density map. The branch and bound
148 algorithm design paradigm²⁵ can accelerate this search by quickly and inexpensively ruling out
149 large regions of the search space that cannot contain the optimum of the objective function
150 (Figure 3A).

151

152 In cryo-EM map refinement, the optimal pose for a particle image minimizes the error between
153 the observed image and a projection of the 3-D map. To find this optimal pose using the branch
154 and bound approach (Figure 3B), an inexpensive lower bound on the error is first computed
155 across the entire space of poses. At the pose that minimizes this lower bound, the
156 computationally expensive true error function is evaluated. All regions of the search space where
157 the lower bound exceeds this computed value of the true error function cannot contain the
158 optimal pose and can be excluded from further search. A new lower bound is then calculated that
159 fits more tightly to the true error function but is more expensive to calculate. The process of
160 evaluating the error function at the optimum of the lower bound, discarding regions of search
161 space where the true error is above the lower bound, and recalculating a tighter-fitting lower
162 bound, is repeated until only the optimal pose remains.

163

164 Although conceptually straightforward, application of the branch and bound strategy requires an
165 informative and inexpensive lower bound for the objective function. Suitable lower bounds are
166 well known for other problems^{26,27} but use of the method for determining the orientations of
167 single particle cryo-EM images required derivation of an appropriate bound (Supplementary
168 Note 2). The derivation we describe was based on the signal-to-noise ratio of single particle
169 images over a range of resolutions. It is worth emphasizing that the branch and bound approach
170 is a global pose search that requires no prior estimate of an optimal pose. In contrast, strategies to
171 accelerate orientation determination based solely on local search risk selection of a pose that is
172 not the global optimum^{12,13}. In practice, an approximation to this branch and bound search is
173 used (Supplementary Note 2) that was found to be equally effective but even more efficient.

174
175 We implemented the branch and bound approach and applied it to high-resolution structure
176 determination from several published datasets: the 20S proteasome from *Thermoplasma*
177 *acidophilum*²⁸, the 80S ribosome from *Plasmodium falciparum*²⁹, amphipol-solubilized rat
178 TRPV1³, as well as the *T. thermophilus* V/A-ATPase²³. Computations were carried out with the
179 same desktop workstation and single NVIDIA Tesla K40 GPU used for *ab initio* SGD
180 calculations. Applied to 35,645 TRPV1 particle images, branch and bound orientation
181 determination produced a 3.3 Å resolution map in 66 minutes with C4 symmetry enforced using
182 a gold-standard refinement procedure³⁰, the FSC=0.143 resolution criterion³¹, and correction for
183 effects of masking on the FSC by high-resolution noise substitution³² (Figure 2C). This
184 resolution slightly exceeds the previously published resolution of 3.4 Å from the same dataset³.
185 With *T. thermophilus* V/A-ATPase particle images sorted into three classes by SGD, the branch
186 and bound search produced maps of all three states in a total of 2.4 hours (Figure 2D). The
187 resolutions estimated for the states were 6.4 Å, 7.6 Å, and 7.9 Å, compared to 6.4 Å and 9.5 Å
188 for the two states identified in the previously published analysis²³.

189
190 Following SGD *ab initio* structure determination, the application of the branch and bound
191 method allowed high-resolution refinement of the 80S ribosome to 3.2 Å resolution, equivalent
192 to the published resolution²⁹, in 2.2 hours (Figure 4A), demonstrating the capability of the
193 method to deal with large and asymmetric protein complexes. Notably, on the same computer
194 hardware (desktop computer with one GPU), this dataset of particle images would take
195 approximately 20 hours for refinement using the GPU accelerated program RELION⁹. Similarly,
196 the 20S proteasome structure was refined to 2.8 Å with D7 symmetry enforced, matching the
197 published sub-3 Å resolution from the dataset²⁸ (Figure 4B) but in only 70 minutes. These
198 refined maps show clear high-resolution detail and side-chain density, illustrating the
199 performance of the method at near-atomic resolution.

200
201

202 Discussion

203

204 *Ab initio* reconstruction of 3-D maps from cryo-EM images has long been known as a significant
205 problem. While random initialization can be successful for highly-symmetric particles³³, this has
206 not been the case for asymmetric or low-degree of symmetry structures where incorrect
207 structures have been published³⁴. Previous approaches for determining low-resolution initial
208 maps often involve collecting image tilt pairs^{35,36}. In that method, the need to switch to a
209 different experimental procedure to generate an initial map is unwieldy and presents a barrier to
210 automated structure determination. Other investigators have proposed algorithms to generate
211 initial maps from images obtained under standard conditions. The approaches have included
212 evolutionary algorithms³⁷, a statistical weighted least squares approach³⁸, complex annealing
213 procedures³⁹, matching of common lines⁴⁰ and statistical weighting⁴¹. However, all of these
214 algorithms rely on analyzing all images in batch, making them intrinsically slower than our
215 approach, particularly as the number of particle images in datasets grow. In contrast, SGD
216 processes random subsets of data at each iteration, making it efficient, even in the face of large
217 datasets. We previously showed that SGD could produce a reasonable low-resolution map *ab*
218 *initio* for a homogenous dataset⁴². Here we have demonstrated that SGD, unlike other
219 approaches, is sufficiently robust to enable reconstruction of multiple 3-D classes from
220 independent arbitrary initializations. To our knowledge, all existing techniques for 3-D
221 classification use a single initial reference from which analysis of heterogeneity proceeds.
222 Removal of the assumption that all 3-D classes are similar to the single input reference is
223 particularly advantageous for discovering 3-D states that are unexpected and different from the
224 consensus structure. It is important to note that, like other algorithms, SGD will fail when the
225 particle images do not contain a sufficient series of views to define the 3-D structure of the
226 molecule. It can also fail if there are sufficient views, but strongly preferred orientations for
227 particles. Other pathological situations may include analysis of datasets with little contrast at
228 low-resolution. This situation can occur when insufficient defocus is used with a cryo-EM
229 microscope that does not possess a phase plate or when imaging low molecular weight
230 complexes.

231

232 Combination of the SGD approach and branch and bound refinement provides a complete
233 framework for rapid *ab initio* calculation of multiple high-resolution maps from a heterogeneous
234 dataset on inexpensive computer hardware. The bound derived and used in this work is based on,
235 and provides a mathematical basis for, the common intuition that high-resolution features in an
236 image contribute less to alignment than low-resolution features. This intuition has previously
237 been used in heuristic methods that perform alignment and reconstruction at iteratively
238 increasing resolution levels¹² or decompose the space of particle images into basis vectors that
239 contain low-resolution features⁴³. A number of heuristic methods have also been employed to
240 accelerate the alignment of particle images to a structure at a fixed resolution. Most commonly,
241 locally restricted high-resolution searches are used in later iterations of refinement, after
242 exhaustive search at early iterations provides a guess for the optimal pose of each image^{12,13}.
243 These approaches can still be computationally expensive, require extra tunable parameters for
244 when to start and how much to restrict local search, and run the risk of missing the optimal
245 alignment. Branch and bound optimization provides a risk-free, parameter-free approach to
246 accelerating computationally expensive search problems, is significantly faster than heuristic
247 methods, and will likely find other applications in cryo-EM image analysis.

248
249 With the recent push to re-implement existing algorithms on new hardware (e.g., GPUs),
250 attempts have also been made to simplify the task of accessing and using computer clusters
251 through cloud computing service providers, notably Amazon EC2⁴⁴. However, even with
252 computer clusters available for rent, existing software methods do not scale well, providing
253 diminishing returns with larger clusters. As the pace of cryo-EM data collection grows, and
254 studies aim to distinguish increasingly subtle structural differences between 3-D classes^{8,45},
255 improved computational efficiency through algorithm development will be a critical enabler for
256 both academic and industrial researchers using cryo-EM.

257
258 The new cryoSPARC software is available as a standalone program that can run on either
259 commodity desktop workstations or rackmount servers. CryoSPARC is also available as a web-
260 service, for new users to try prior to installing locally. Once particle images are selected and
261 corrected for anisotropic beam-induced movement⁴⁶ and the effects of radiation damage^{46,47} they
262 may be processed through the program's web-browser graphical user interface (GUI). At
263 minimum, a single consumer or professional grade NVIDIA GPU is required. The easy-to-use
264 GUI (Supplementary Video 1) provides the same interface through both the web-service and in
265 local installations. This GUI allows for multiple users within a laboratory to have separate
266 accounts, access the program remotely, upload and share datasets, manage experimental results,
267 launch computational tasks, and view results streaming in real-time as they are computed. A
268 protocol detailing use of the software package has been prepared⁴⁹ (Supplementary Protocol).

269 270 **Software availability**

271 The software package, including source code, is available for non-commercial use as a download
272 and as a web service at www.cryosparc.com. Results reported in this work were computed using
273 cryoSPARC version 0.2.36.

274 275 Accession Codes

276 Data Availability Statement

277 The cryo-EM images used to experimentally demonstrate the effectiveness of algorithms were
278 taken from previously published studies. Several datasets were downloaded directly from
279 EMPIAR⁴⁸, including the 80S Ribosome (EMPIAR-10028), 20S proteasome (EMPIAR-10025),
280 TRVP1 channel (EMPIAR-10005). Images of the *T. thermophilus* V/A-ATPase are available
281 from the authors upon request. In all cases, the single particle images that were used in the
282 original studies were input directly into cryoSPARC, with no further preprocessing.

283 284 Acknowledgements

285 We thank Suhail Dawood for construction of the GUI front end and members of the Rubinstein
286 laboratory for testing cryoSPARC. AP was supported by a scholarship from the Natural Sciences
287 and Engineering Research Council (NSERC), JLR was supported by the Canada Research Chairs
288 program, and DJF was supported in part by the Learning in Machines and Brains program of the
289 Canadian Institute for Advanced Research. This research was also supported by NSERC
290 Discovery Grants RGPIN 2015-05630 (DJF) and 401724-12 (JLR), and an NVIDIA Academic
291 Hardware Grant (MAB, AP). Part of this work was performed while MAB was a postdoctoral
292 fellow at the University of Toronto.

293 Author Contributions

294 AP and MAB designed algorithms and implemented software. AP, MAB and JLR performed
295 experimental work. JLR, DJF and MAB contributed expertise and supervision. All authors
296 contributed to manuscript preparation.

297

298 Competing Financial Interests

299 All authors are engaged in a venture to commercially support cryoSPARC for industrial use.

300

301

302 References for main text

- 303 1. Kühlbrandt, W. *et al.* Biochemistry. The resolution revolution. *Science* **343**, 1443–4
304 (2014).
- 305 2. Smith, M. T. J. & Rubinstein, J. L. Beyond blob-ology. *Science (80-.)*. **345**, 617–619
306 (2014).
- 307 3. Liao, M., Cao, E., Julius, D. & Cheng, Y. Structure of the TRPV1 ion channel determined
308 by electron cryo-microscopy. *Nature* **504**, 107–12 (2013).
- 309 4. Bai, X. C., Fernandez, I. S., McMullan, G. & Scheres, S. H. W. Ribosome structures to
310 near-atomic resolution from thirty thousand cryo-EM particles. *Elife* **2013**, 2–13 (2013).
- 311 5. Yan, C. *et al.* Structure of a yeast spliceosome at 3.6-angstrom resolution. *Science* **349**,
312 1182–1191 (2015).
- 313 6. Banerjee, S. *et al.* 2.3 Å resolution cryo-EM structure of human p97 and mechanism of
314 allosteric inhibition. *Science* **351**, 871–5 (2016).
- 315 7. Sirohi, D. *et al.* The 3.8 Å resolution cryo-EM structure of Zika virus. *Science* **352**, 467–
316 70 (2016).
- 317 8. Abeyrathne, P. D., Koh, C. S., Grant, T., Grigorieff, N. & Korostelev, A. A. Ensemble
318 cryo-EM uncovers inchworm-like translocation of a viral IRES through the ribosome.
319 *Elife* **5**, e14874 (2016).
- 320 9. Kimanius, D., Forsberg, B. O., Scheres, S. & Lindahl, E. *Accelerated cryo-EM structure*
321 *determination with parallelisation using GPUs in RELION-2*. *bioRxiv* (Cold Spring
322 Harbor Labs Journals, 2016). doi:10.1101/059717
- 323 10. Henderson, R. Avoiding the pitfalls of single particle cryo-electron microscopy: Einstein
324 from noise. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 18037–41 (2013).
- 325 11. Henderson, R. *et al.* Outcome of the first electron microscopy validation task force
326 meeting. *Structure* **20**, 205–214 (2012).
- 327 12. Scheres, S. H. W. A bayesian view on cryo-EM structure determination. *J. Mol. Biol.* **415**,
328 406–418 (2012).
- 329 13. Grigorieff, N. FREALIGN: High-resolution refinement of single particle structures. *J.*
330 *Struct. Biol.* **157**, 117–125 (2007).
- 331 14. Zhang, K. Gctf: Real-time CTF determination and correction. *J. Struct. Biol.* **193**, 1–12
332 (2016).
- 333 15. Hoang, T. V, Cavin, X., Schultz, P. & Ritchie, D. W. gEMpicker: a highly parallel GPU-
334 accelerated particle picking tool for cryo-electron microscopy. *BMC Struct. Biol.* **13**, 25
335 (2013).
- 336 16. Moore, G. E. Progress in digital integrated electronics. *1975 Int. Electron Devices Meet.*
337 **21**, 11–13 (1975).
- 338 17. Sigworth, F. J. A maximum-likelihood approach to single-particle image refinement. *J.*
339 *Struct. Biol.* **122**, 328–39 (1998).
- 340 18. Nocedal, J. & Wright, S. J. *Numerical Optimization*. (Springer New York, 2000).
341 doi:10.1007/BF01068601
- 342 19. Calafiore, G. C. & El Ghaoui, L. *Optimization Models*. (Cambridge University Press,
343 2014).
- 344 20. Bottou, L. Large-Scale Machine Learning with Stochastic Gradient Descent. *Proc.*
345 *COMPSTAT'2010* 177–186 (2010). doi:10.1007/978-3-7908-2604-3_16
- 346 21. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet Classification with Deep
347 Convolutional Neural Networks. *Adv. Neural Inf. Process. Syst.* 1–9 (2012).

- 348 doi:<http://dx.doi.org/10.1016/j.protcy.2014.09.007>
- 349 22. Taigman, Y., Yang, M., Ranzato, M. & Wolf, L. DeepFace: Closing the gap to human-
350 level performance in face verification. *Proc. IEEE Comput. Soc. Conf. Comput. Vis.*
351 *Pattern Recognit.* 1701–1708 (2014). doi:10.1109/CVPR.2014.220
- 352 23. Schep, D. G., Zhao, J. & Rubinstein, J. L. Models for the a subunits of the Thermus
353 thermophilus V/A-ATPase and Saccharomyces cerevisiae V-ATPase enzymes by cryo-
354 EM and evolutionary covariance. *Proc. Natl. Acad. Sci. U. S. A.* **113**, 3245–3250 (2016).
- 355 24. Zhao, J., Benlekbir, S. & Rubinstein, J. Electron cryomicroscopy observation of rotational
356 states in a eukaryotic V-ATPase. *Nature* **521**, 241–245 (2015).
- 357 25. Kearfott, R. B. *Rigorous global search: Continuous problems. Igarss 2014* (Springer US,
358 2014). doi:10.1007/s13398-014-0173-7.2
- 359 26. Little, J. D. C., Karel, C., Murty, K. G. & Sweeney, D. W. An algorithm for the traveling
360 salesman problem. *Oper. Res.* **11**, 972–989 (1963).
- 361 27. Yang, J., Li, H. & Jia, Y. Go-ICP: Solving 3D Registration Efficiently and Globally
362 Optimally. *2013 IEEE Int. Conf. Comput. Vis.* 1457–1464 (2013).
363 doi:10.1109/ICCV.2013.184
- 364 28. Campbell, M. G., Veessler, D., Cheng, A., Potter, C. S. & Carragher, B. 2.8 ?? Resolution
365 Reconstruction of the Thermoplasma Acidophilum 20 S Proteasome Using Cryo-Electron
366 Microscopy. *Elife* **2015**, 1–22 (2015).
- 367 29. Wong, W. *et al.* Cryo-EM structure of the Plasmodium falciparum 80S ribosome bound to
368 the anti-protozoan drug emetine. *Elife* **2014**, 1–20 (2014).
- 369 30. Scheres, S. H. W. & Chen, S. Prevention of overfitting in cryo-EM structure
370 determination. *Nat. Methods* **9**, 853–854 (2012).
- 371 31. Rosenthal, P. B. & Henderson, R. Optimal determination of particle orientation, absolute
372 hand, and contrast loss in single-particle electron cryomicroscopy. *J. Mol. Biol.* **333**, 721–
373 745 (2003).
- 374 32. Chen, S. *et al.* High-resolution noise substitution to measure overfitting and validate
375 resolution in 3D structure determination by single particle electron cryomicroscopy.
376 *Ultramicroscopy* **135**, 24–35 (2013).
- 377 33. Yan, X., Cardone, G., Zhang, X., Zhou, Z. H. & Baker, T. S. Single particle analysis
378 integrated with microscopy: A high-throughput approach for reconstructing icosahedral
379 particles. *J. Struct. Biol.* **186**, 8–18 (2014).
- 380 34. Murray, S. C. *et al.* Validation of cryo-EM structure of IP3R1 channel. *Structure* **21**, 900–
381 909 (2013).
- 382 35. Radermacher, M., Wagenknecht, T., Verschoor, A. & Frank, J. Three-dimensional
383 reconstruction from a single-exposure, random conical tilt series applied to the 50S
384 ribosomal subunit of Escherichia coli. *J. Microsc.* **146**, 113–136 (1987).
- 385 36. Leschziner, A. E. & Nogales, E. The orthogonal tilt reconstruction method: An approach
386 to generating single-class volumes with no missing cone for ab initio reconstruction of
387 asymmetric particles. *J. Struct. Biol.* **153**, 284–299 (2006).
- 388 37. Penczek, P. A. & Asturias, F. J. Ab initio cryo-EM structure determination as a validation
389 problem. in *2014 IEEE International Conference on Image Processing (ICIP)* 2090–2094
390 (2014). doi:10.1109/ICIP.2014.7025419
- 391 38. Sorzano, C. O. S. *et al.* A statistical approach to the initial volume problem in Single
392 particle analysis by electron microscopy. *J. Struct. Biol.* **189**, 213–219 (2015).
- 393 39. Jaitly, N., Brubaker, M. A., Rubinstein, J. L. & Lilien, R. H. A Bayesian method for 3D

394 macromolecular structure inference using class average images from single particle
395 electron microscopy. *Bioinformatics* **26**, 2406–2415 (2010).

396 40. Elmlund, D. & Elmlund, H. SIMPLE: Software for ab initio reconstruction of
397 heterogeneous single-particles. *J. Struct. Biol.* **180**, 420–427 (2012).

398 41. Elmlund, H., Elmlund, D. & Bengio, S. PRIME: Probabilistic initial 3D model generation
399 for single-particle cryo-electron microscopy. *Structure* **21**, 1299–1306 (2013).

400 42. Brubaker, M. A., Punjani, A. & Fleet, D. J. Building proteins in a day: Efficient 3D
401 molecular reconstruction. in *Proceedings of the IEEE Computer Society Conference on*
402 *Computer Vision and Pattern Recognition* **07–12–June**, (2015).

403 43. Dvornek, N. C., Sigworth, F. J. & Tagare, H. D. SubspaceEM: A fast maximum-a-
404 posteriori algorithm for cryo-EM single particle reconstruction. *J. Struct. Biol.* **190**, 200–
405 214 (2015).

406 44. Cianfrocco, M. A. & Leschziner, A. E. Low cost, high performance processing of single
407 particle cryo-electron microscopy data in the cloud. *Elife* **4**, e06664 (2015).

408 45. Bai, X.-C., Rajendra, E., Yang, G., Shi, Y. & Scheres, S. H. Sampling the conformational
409 space of the catalytic subunit of human γ -secretase. *Elife* **4**, e11182 (2015).

410 46. Rubinstein, J. L. & Brubaker, M. A. Alignment of cryo-EM movies of individual particles
411 by optimization of image translations. *J. Struct. Biol.* **192**, 188–195 (2015).

412 47. Grant, T. & Grigorieff, N. Measuring the optimal exposure for single particle cryo-EM
413 using a 2.6 Å reconstruction of rotavirus VP6. *Elife* **4**, e06980 (2015).

414 48. Iudin, A., Korir, P. K., Salavert-Torres, J., Kleywegt, G. J. & Patwardhan, A. EMPIAR: a
415 public archive for raw electron microscopy image data. *Nat. Methods* **13**, 387–388 (2016).

416 49. Punjani, A., Rubinstein, J., Fleet, D. and Brubaker, M. “Protocol for rapid unsupervised
417 cryo-EM structure determination using cryoSPARC software” *Protocol*
418 *Exchange* (2016) DOI: [to-be-provided].
419
420
421

422 Figure Legends for main text

423

424 **Figure 1. Stochastic gradient descent for cryo-EM map calculation.** **A**, Iterative refinement
425 methods are sensitive to initialization. An arbitrary initialization far from the correct 3-D map
426 will be refined into an incorrect structure that attains a locally optimal probability within the
427 space of all 3-D maps. An accurate initialization will be refined to the correct structure. Iterative
428 refinement uses all single particle images in a dataset to compute each step. **B**, Random selection
429 of particle images in the SGD algorithm. At each iteration, a different small random selection of
430 images is used to approximate the true optimization objective. Each iteration may use a different
431 number of images. **C**, Stochastic Gradient Descent (SGD) algorithm enables *ab initio* structure
432 determination through insensitivity to initialization. An arbitrary computer generated random
433 initialization is incrementally improved by many noisy steps. Each step is based on the gradient
434 of the approximated objective function obtained by random selection in (B). These approximate
435 gradients do not exactly match the overall optimization objective. The success of SGD is
436 commonly explained by the noisy sampling approximation allowing the algorithm to widely
437 explore the space of all 3-D maps to arrive finally near the correct structure.

438

439 **Figure 2. Evolution of 3-D cryo-EM maps as computation progresses.** **A**. Low-resolution
440 map of the TRPV1 channel calculated in 75 minutes from random initialization. **B**. Multiple
441 conformations of the *Thermus thermophilus* V/A-ATPase calculated simultaneously from
442 separate random initializations. **C**. Refinement of TRPV1 to 3.3 Å resolution on a single GPU
443 desktop workstation in 66 minutes with C4 symmetry enforced. Density is apparent that
444 corresponds to amino acid side chains. **D**. Refinement of each of three V/A-ATPase rotational
445 states. The rotational state of the central rotor (indicated by red circles) is seen in cross sections
446 through the 3-D maps. All computations were performed on a single desktop computer with a
447 single NVIDIA Tesla K40 GPU. Scale bars, 25 Å.

448

449 **Figure 3. The branch and bound approach to high-resolution cryo-EM map refinement.** **A**,
450 Two iterations of a simplified 1-D representation of the branch and bound approach. Candidate
451 poses are iteratively eliminated by evaluation of an inexpensive lower bound over all poses, and
452 the true error function at the minimum of the lower bound. **B**, For cryo-EM images, the true error
453 function over all poses (top left) for an individual particle (top right) is never evaluated. Instead,
454 the entire lower bound is computed (middle left), the true error is calculated at the minimum of
455 the bound, and all poses where the lower bound exceeds this calculated error are eliminated
456 (middle right). A tighter lower bound is evaluated and the process repeated until the optimum
457 pose is identified (bottom left and right).

458

459 **Figure 4. High-resolution structures from branch and bound refinement.** **A**, 80S ribosome
460 structure refined to 3.2 Å resolution in 2.2 h with 105,247 particle images. Amino acids side
461 chain and RNA base densities are clearly visible in α -helices, β -strands, and rRNA (inset). **B**, A
462 20S proteasome map refined to 2.8 Å in 70 min with 49,954 particle images and D7 symmetry
463 enforced. Well-resolved densities are apparent for small and large residues (inset). Branch and
464 bound refinement of both structures was initialized with *ab initio* maps from SGD. Scale bars, 25
465 Å.

466

467 Tables: None.

468 Online Methods

469

470 **Statistics**

471

472 In all 3-D map refinement experiments, the Fourier shell correlation (FSC) between two
473 independently refined half-maps (the “gold standard”) was used to assess resolution³⁰, along with
474 the FSC=0.143 resolution criterion³¹ and correction of the FSC for effects of masking by high-
475 resolution noise substitution³².

476

477 **Computational Hardware**

478

479 All experiments were carried out on a single desktop workstation, equipped with an Intel i7-
480 5820K 6-core CPU, NVIDIA Tesla K40 GPU, 64GB of CPU RAM, and a 512GB SSD for file
481 storage. Tests were also run and equivalent running times were achieved using the consumer-
482 grade NVIDIA Titan Z GPU. It should be noted that at the time of writing, the Tesla K40 GPU is
483 over two years old, and more recent GPU cards will perform significantly faster.

484

485 **Implementation**

486

487 CryoSPARC is a software package written in a mixture of Python, CUDA C, and Javascript.
488 Algorithms are implemented in Python and the GPU computation routines are written in CUDA
489 C. Computations are parallelized over images, pixels, and search parameters. Two CPU threads
490 are used for the GPU to improve utilization, and images are loaded from SSD and prepared by
491 the CPU simultaneously with GPU processing of a different batch of images.

492

493 **Stochastic Gradient Descent**

494

495 SGD is initialized from a computer generated random initialization for each 3D class
496 (Supplementary Note 1). The number of images used in each iteration of SGD is automatically
497 determined based on the current resolution. A model of the noise level in single particle images
498 is initialized with an over-estimate relative to measured noise levels. Approximate gradients of
499 Equation 1 are computed along with second-order curvature information to enable estimation of
500 an optimal step size for descent at each iteration. Step directions are averaged over iterations
501 using a classical momentum method⁵⁰. Resulting iterative steps are applied to the 3-D maps and
502 a projection operation is used to enforce non-negativity of 3-D map density. The noise model is
503 refined based on errors between the images and projections of the 3-D map at each iteration,
504 converging to the optimal noise model over several iterations. The descent step size is decreased
505 monotonically over iterations to improve convergence once an approximately correct structure is
506 found. Further details can be found in Supplementary Note 1.

507

508 **Branch and Bound Image Alignment**

509

510 The branch and bound method is applied to each image individually at each iteration of high-
511 resolution map refinement. A space partitioning tree-structure is used to segment the space of
512 orientation parameters, which are represented using axis-angle coordinates. A coarse initial
513 sampling of the orientation space forms the first level of the tree, and each stage of branch and

514 bound subdivides and prunes branches in the tree until only the optimal pose remains to within a
515 specified angular precision of 0.18° . A similar tree structure is used to segment and subdivide the
516 2-D space of in-plane shifts for each image, resulting in a specified translational precision of
517 0.04 pixels. Further details including the derived lower bound and approximations can be found
518 in Supplementary Note 2.

519

520 **Program Settings**

521

522 Default cryoSPARC settings were used in all refinement experiments, and the number of classes
523 was set in each *ab initio* reconstruction experiment. Symmetry was enforced in refinement
524 experiments where noted, but not in *ab initio* reconstruction.

525

526

527

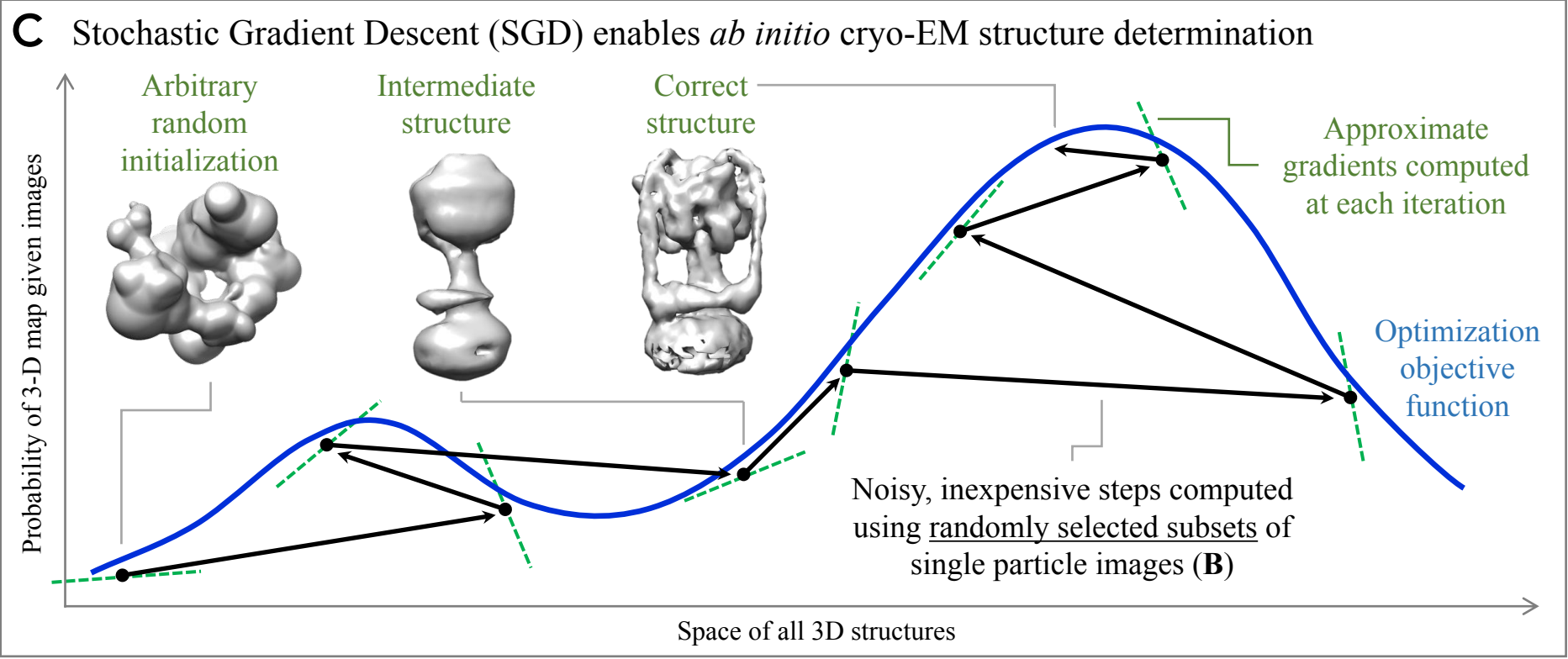
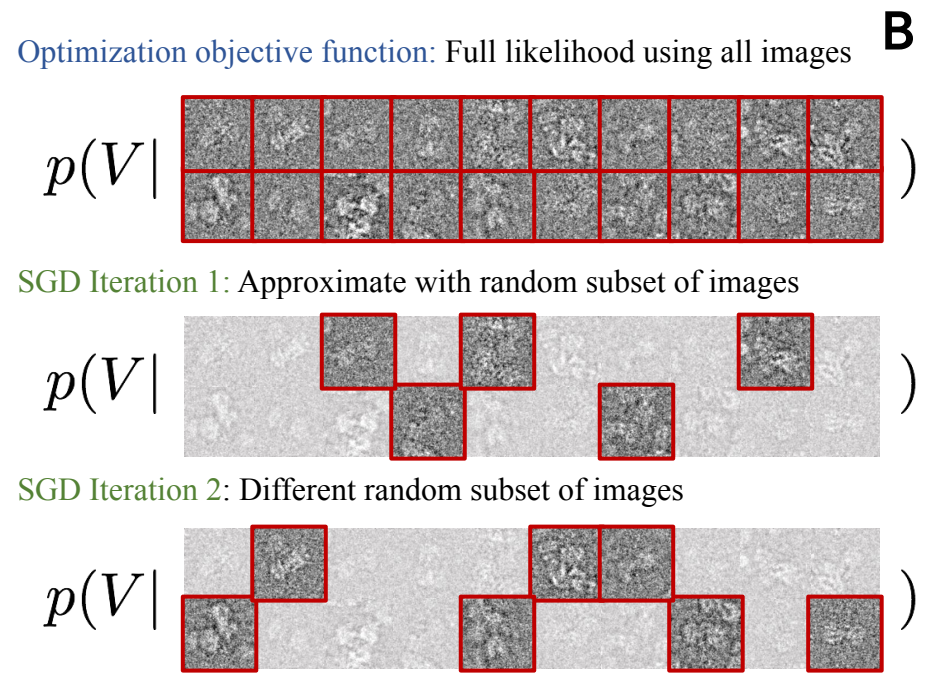
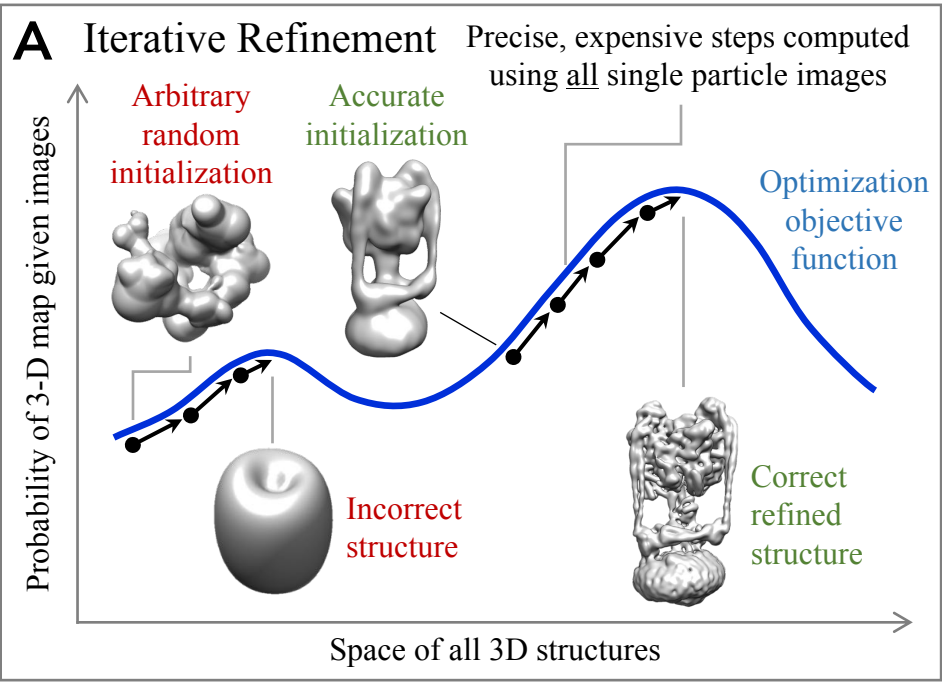
528 Methods-only References

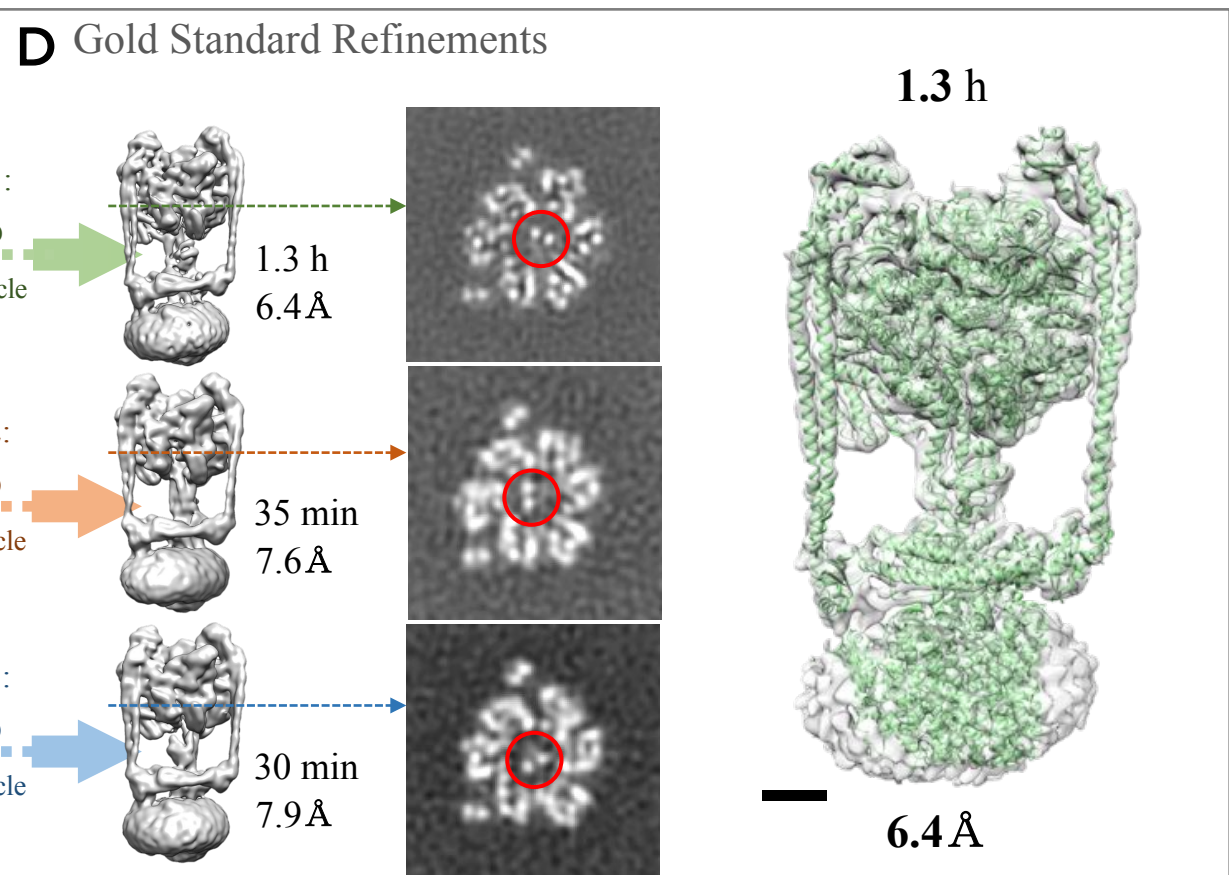
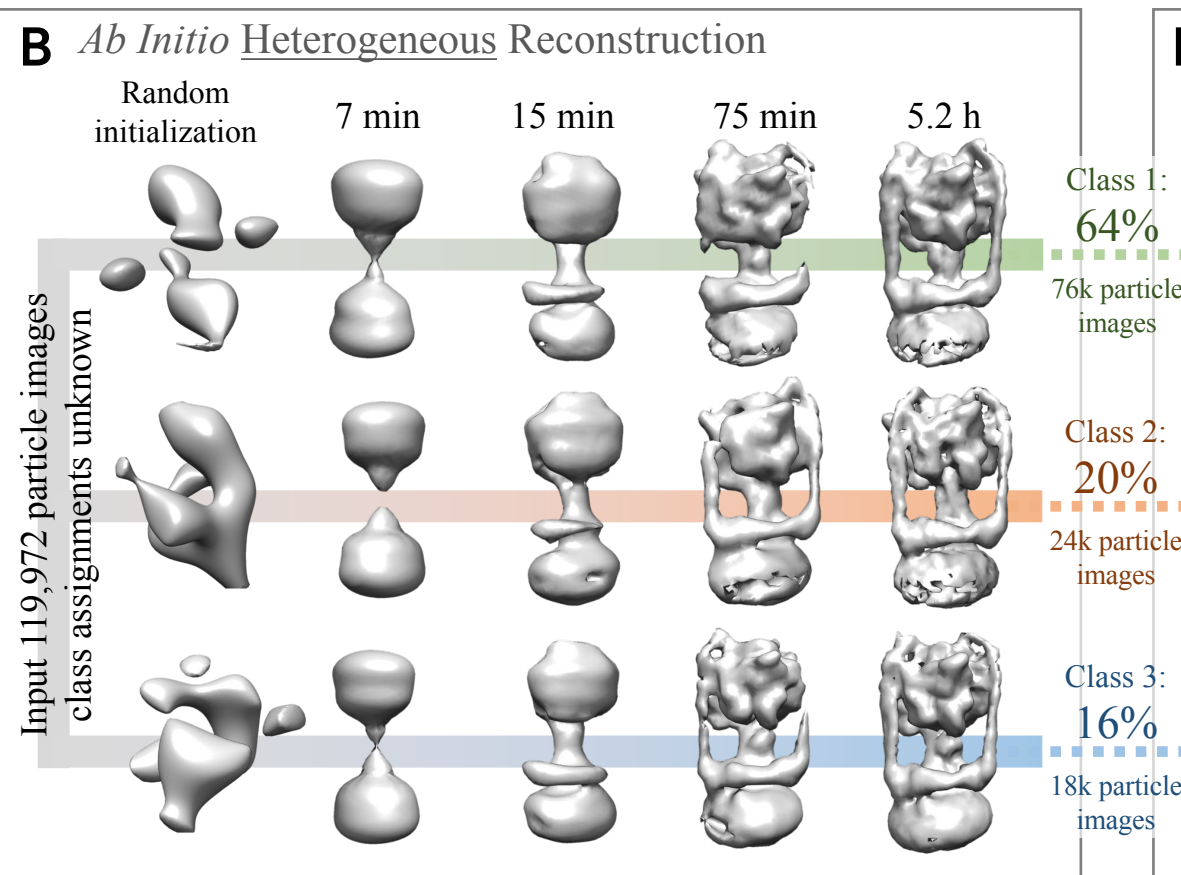
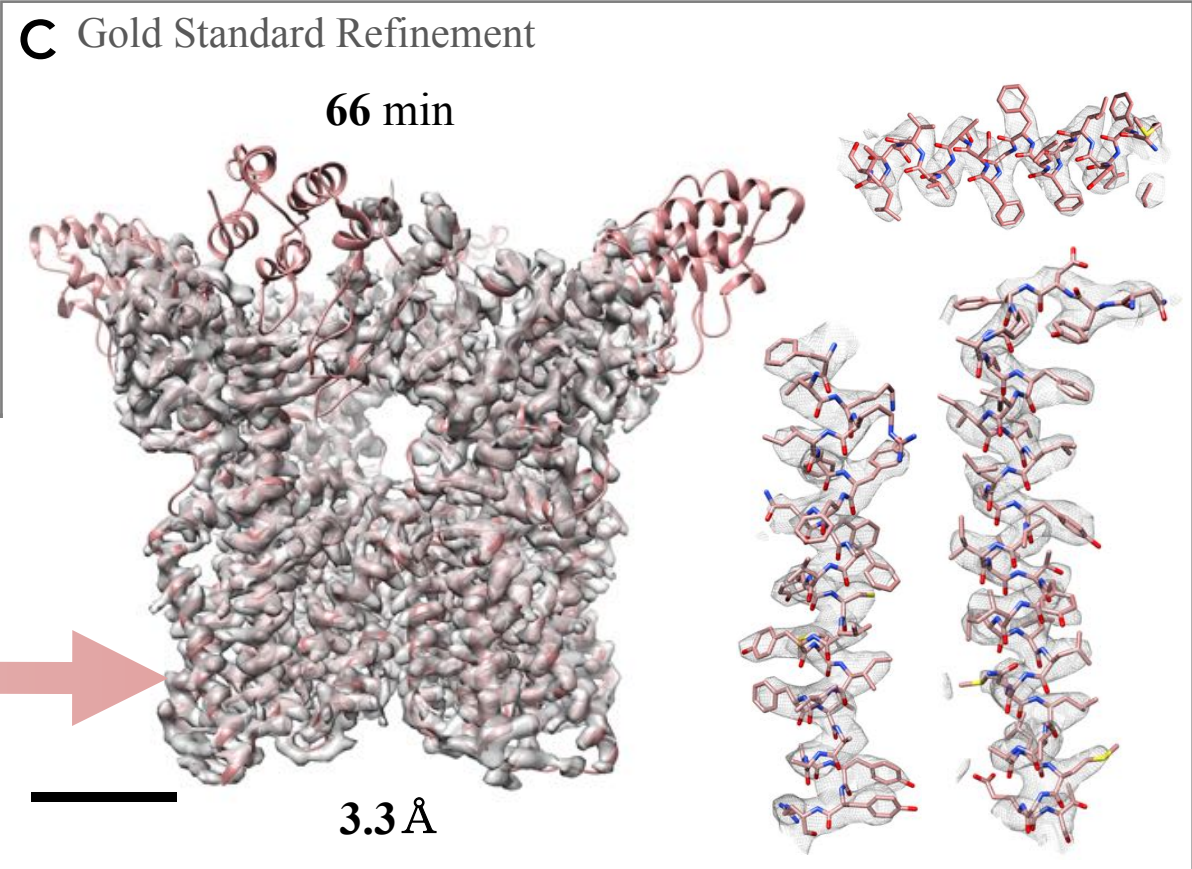
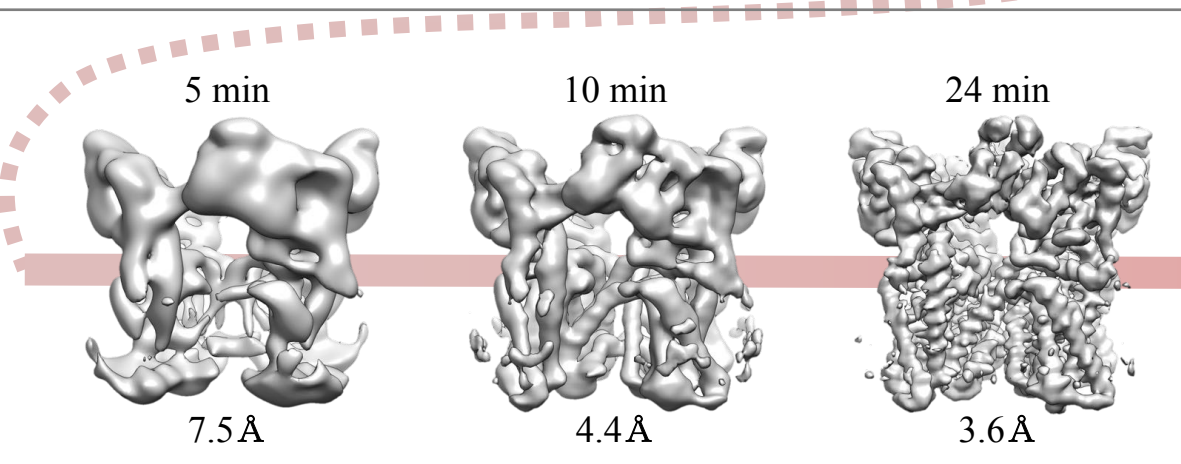
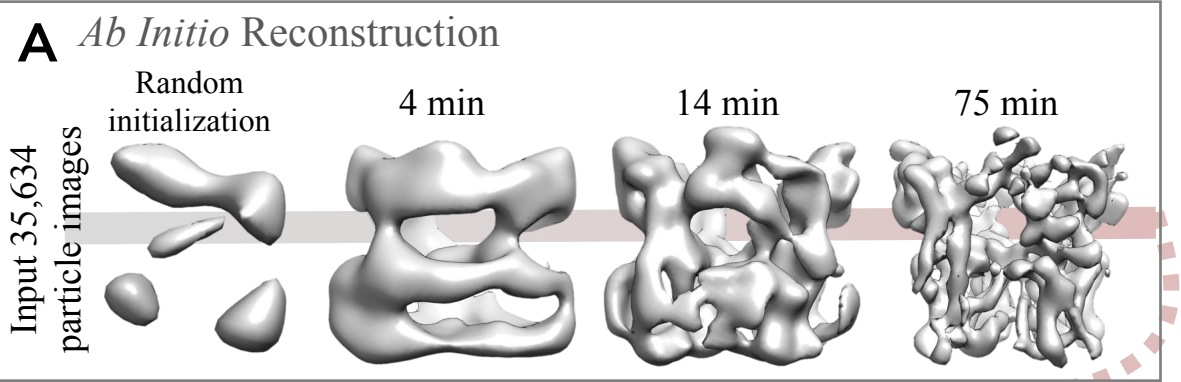
529

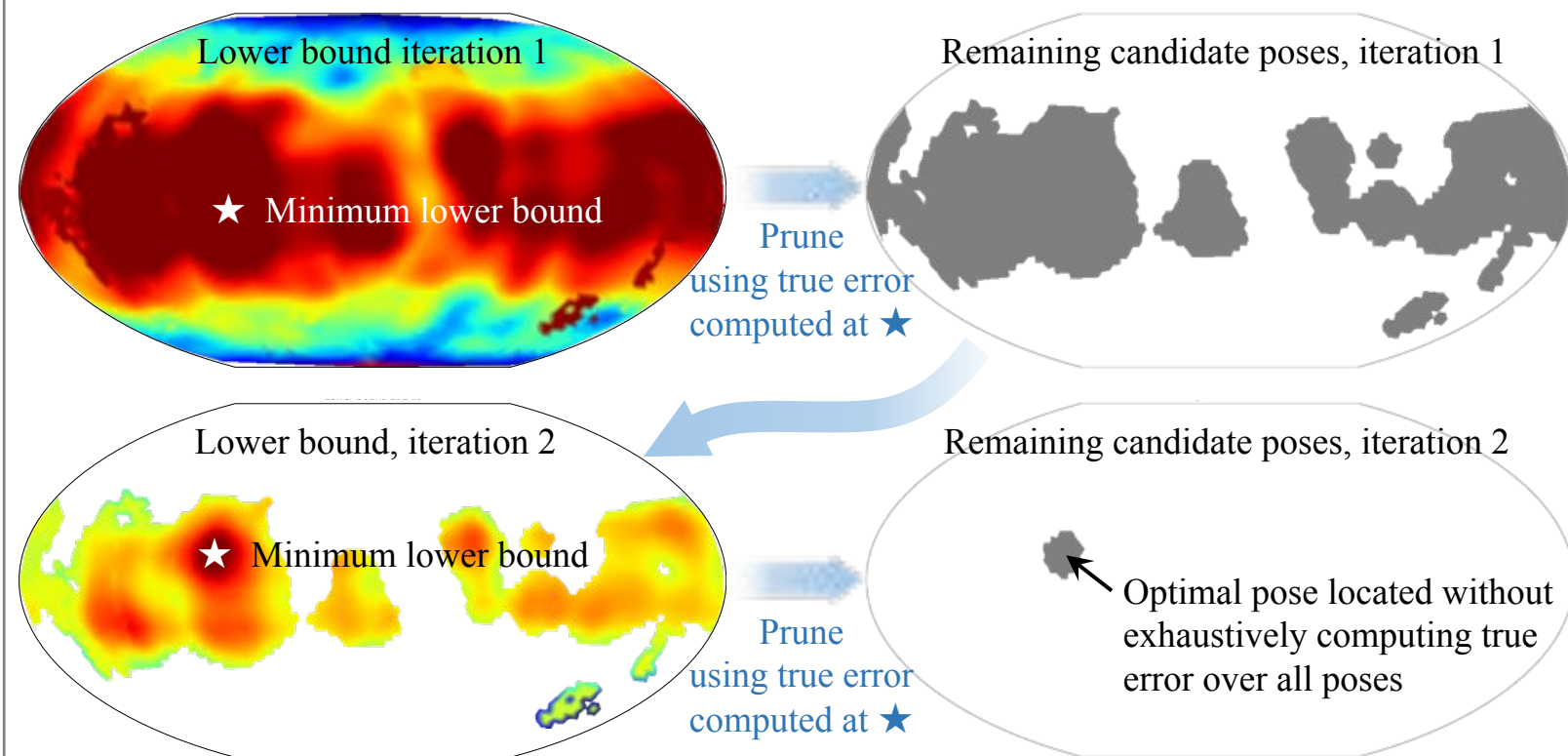
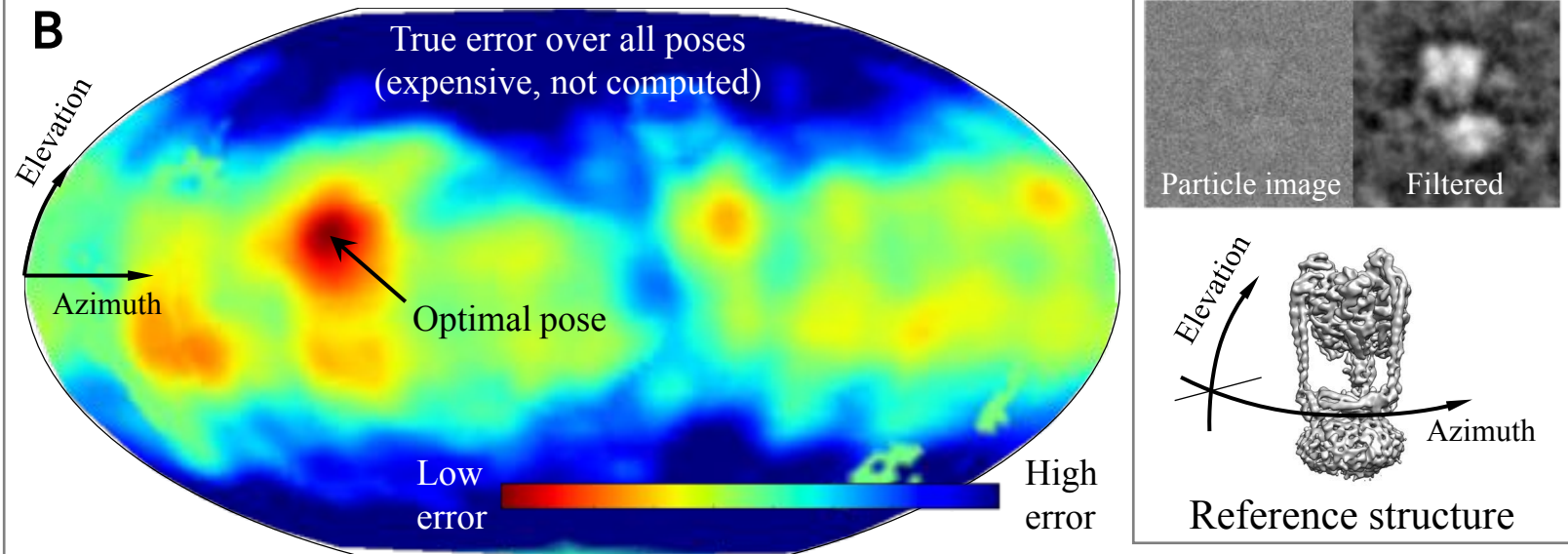
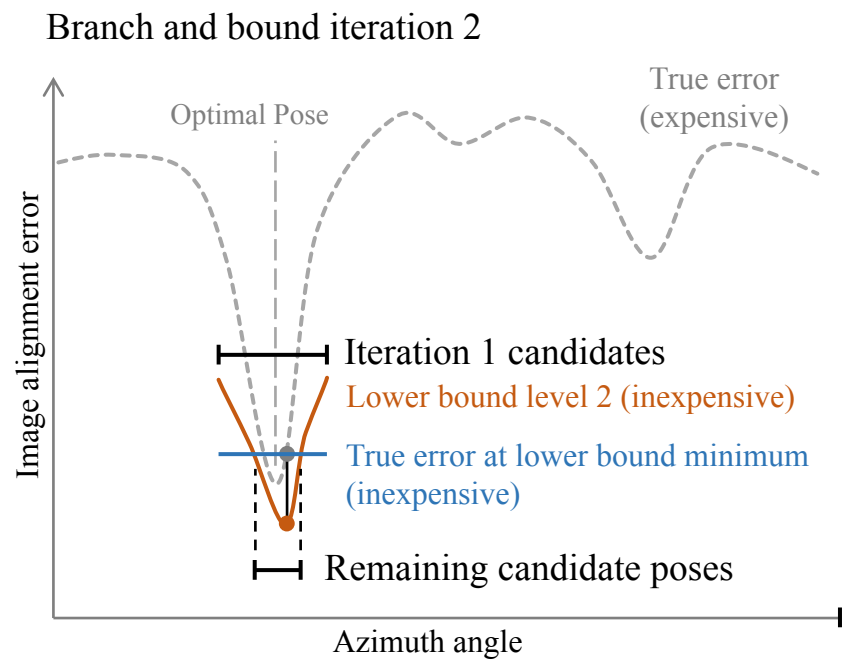
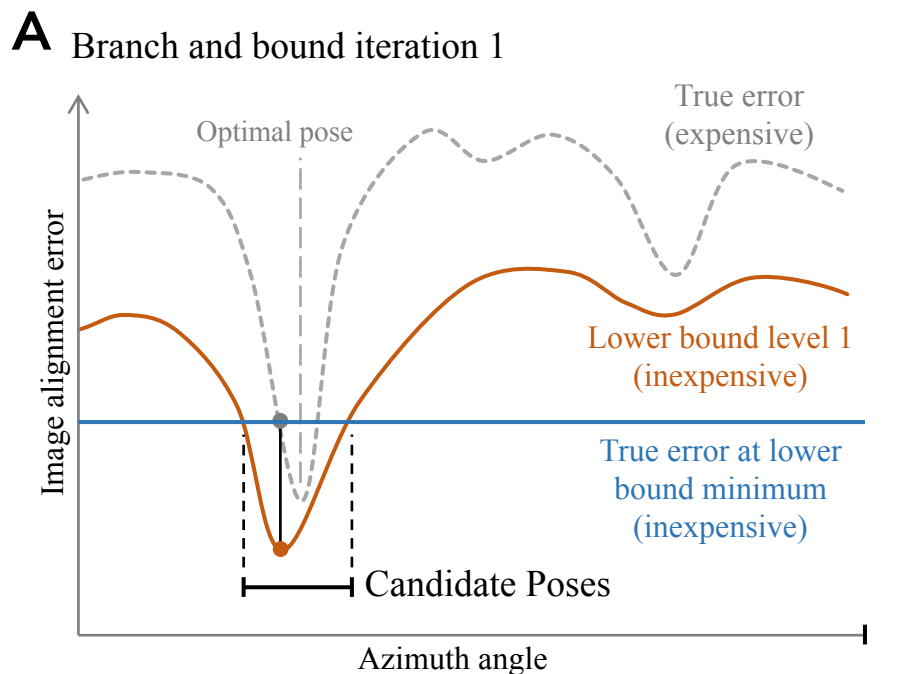
- 530 50. Sutskever, I., Martens, J., Dahl, G. E. & Hinton, G. E. On the importance of initialization
531 and momentum in deep learning. *ICML* **28**, 1139–1147 (2013).

532

533

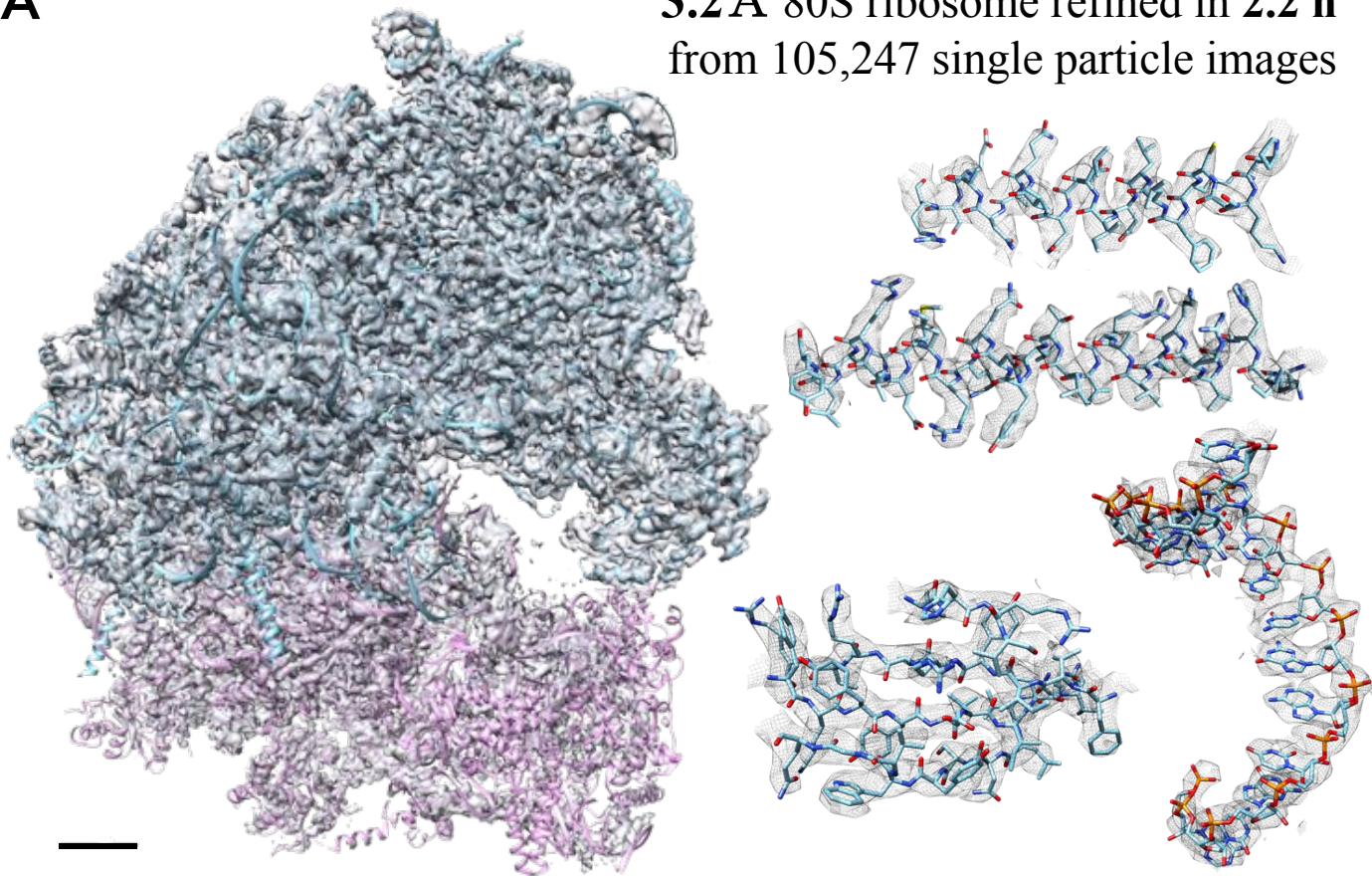






A

3.2 Å 80S ribosome refined in **2.2 h**
from 105,247 single particle images

**B**

2.8 Å T20S proteasome refined in **70 min**
from 49,954 single particle images

