

1 ***cis*TEM, User-friendly software for single-particle image processing**

2
3 Tim Grant ¹, Alexis Rohou ^{1,2} and Nikolaus Grigorieff ¹

4 ¹ Janelia Research Campus, Howard Hughes Medical Institute, Ashburn, Virginia, USA

5 ² Present address: Department of Structural Biology, Genentech, South San Francisco, California, USA

6
7 **Abstract**

8 We have developed new open-source software called *cis*TEM (computational imaging system
9 for transmission electron microscopy) for the processing of data for high-resolution electron
10 cryo-microscopy and single-particle averaging. *cis*TEM features a graphical user interface that is
11 used to submit jobs, monitor their progress, and display results. It implements a full processing
12 pipeline including movie processing, image defocus determination, automatic particle picking,
13 2D classification, ab-initio 3D map generation from random parameters, 3D classification, and
14 high-resolution refinement and reconstruction. Some of these steps implement newly-developed
15 algorithms; others were adapted from previously published algorithms. The software is
16 optimized to enable processing of typical datasets (2000 micrographs, 200k – 300k particles) on
17 a high-end, CPU-based workstation in half a day or less, comparable to GPU-accelerated
18 processing. Jobs can also be scheduled on large computer clusters using flexible run profiles that
19 can be adapted for most computing environments. *cis*TEM is available for download from
20 cistem.org.

22 **Introduction**

23 The three-dimensional (3D) visualization of biological macromolecules and their assemblies by
24 single-particle electron cryo-microscopy (cryo-EM) has become a prominent approach in the
25 study of molecular mechanisms (Cheng et al., 2015; Subramaniam et al., 2016). Recent advances
26 have been primarily due to the introduction of direct electron detectors (McMullan et al., 2016).
27 With the improved data quality, there is increasing demand for advanced computational
28 algorithms to extract signal from the noisy image data and reconstruct 3D density maps from
29 them at the highest possible resolution. The promise of near-atomic resolution (3 – 4 Å), where
30 densities can be interpreted reliably with atomic models, has been realized by many software
31 tools and suites (Frank et al., 1996; Hohn et al., 2007; Lyumkis et al., 2013; Punjani et al., 2017;
32 Scheres, 2012; Tang et al., 2007; van Heel et al., 1996). Many of these tools implement a
33 standard set of image processing steps that are now routinely performed in a single particle
34 project. These typically include movie frame alignment, contrast transfer function (CTF)
35 determination, particle picking, two-dimensional (2D) classification, 3D reconstruction,
36 refinement and classification, and sharpening of the final reconstructions.

37 We have written new software called *cis*TEM to implement a complete image processing
38 pipeline for single-particle cryo-EM, including all these steps, accessible through an easy-to-use
39 graphical user interface (GUI). Some of these steps implement newly-developed algorithms
40 described below; others were adapted from previously published algorithms. *cis*TEM consists of
41 a set of compiled programs and tools, as well as a wxWidgets-based GUI. The GUI launches
42 programs and controls them by sending specific commands and receiving results via TCP/IP
43 sockets. Each program can also be run manually, in which case it solicits user input on the
44 command line. The design of *cis*TEM, therefore, allows users who would like to have more

45 control over the different processing steps to design their own procedures outside the GUI. To
46 adopt this new architecture, a number of previously existing Fortran-based programs were
47 rewritten in C++, including Unblur and Summovie (Grant and Grigorieff, 2015b),
48 mag_distortion_estimate and mag_distortion_correct (Grant and Grigorieff, 2015a), CTFFIND4
49 (Rohou and Grigorieff, 2015), and Frealign (Lyumkis et al., 2013). Additionally, algorithms
50 described previously were added for particle picking (Sigworth, 2004), 2D classification
51 (Scheres et al., 2005) and ab-initio 3D reconstruction (Grigorieff, 2016), sometimes with
52 modifications to optimize their performance. *cisTEM* is open-source and distributed under the
53 Janelia Research Campus Software License
54 (http://license.janelia.org/license/janelia_license_1_2.html).

55 *cisTEM* currently does not support computation on graphical processing units (GPUs).

56 Benchmarking of a hotspot identified in the global orientational search to determine particle
57 alignment parameters showed that an NVIDIA K40 GPU performs approximately as well as 16
58 Xeon E5-2687W CPU cores after the code was carefully optimized for the respective hardware
59 in both cases. Since CPU code is more easily maintained and more generally compatible with
60 existing computer hardware, the potential benefit of GPU-adapted code is primarily the lower
61 cost of a high-end GPU compared with a high-end CPU. We chose to focus on optimizing our
62 code for CPU.

63

64 **Results**

65 *Movie alignment and CTF determination*

66 Movie alignment and CTF determination are based on published algorithms previously
67 implemented in Unblur and Summovie (Grant and Grigorieff, 2015b), and CTFFIND4 (Rohou
68 and Grigorieff, 2015), respectively, and these are therefore only briefly described here. Unblur
69 determines the translations of individual movie frames necessary to bring features (particles)
70 visible in the frames into register. Each frame is aligned against a sum of all other frames that is
71 iteratively updated until there is no further change in the translations. The trajectories along the
72 x- and y-axes are smoothed using a Savitzky–Golay filter to reduce the possibility of spurious
73 translations. Summovie uses the translations to calculate a final frame average with optional
74 exposure filtering to take into account radiation damage of protein and maximize its signal in the
75 final average. *cis*TEM combines the functionality of Unblur and Summovie into a single panel
76 and exposes all relevant parameters to the user (Figure 1). Both programs were originally written
77 in Fortran and have been rewritten entirely in C++.

78 CTFFIND4 fits a calculated two-dimensional CTF to Thon rings (Thon, 1966) visible in the
79 power spectrum calculated from either images or movies. The fitted parameters include
80 astigmatism and, optionally, phase shifts generated by phase plates. When computed from
81 movies, the Thon rings are often more clearly visible compared to Thon rings calculated from
82 images (Figure 2; (Bartesaghi et al., 2014)). When selecting movies as inputs, the user can
83 specify how many frames should be summed to calculate power spectra. An optimal value to
84 amplify Thon rings would be to sum the number of frames that correspond to an exposure of
85 about 4 electrons/Å² (McMullan et al., 2015).

86 Since our original description of the CTFFIND4 algorithm (Rohou and Grigorieff, 2015),
87 several significant changes were introduced. (1) The initial exhaustive search over defocus
88 values can now be performed using a one-dimensional version of the CTF (i.e. with only two

89 parameters: defocus and phase shift) against a radial average of the amplitude spectrum. This
90 search is much faster than the equivalent search over the 2D CTF parameters (i.e., four
91 parameters: two for defocus, one for astigmatism angle and one for phase shift) and can be
92 expected to perform well except in cases of very large astigmatism (Zhang, 2016). Once an
93 initial estimate of the defocus parameter has been obtained, it is refined by a conjugate gradient
94 minimizer against the 2D amplitude spectrum, as done previously. In *cis*TEM, the default
95 behavior is to perform the initial search over the 1D amplitude spectrum, but the user can revert
96 to previous behavior by setting a flag in the “Expert Options” of the “Find CTF” Action panel.

97 (2) If the input micrograph’s pixel size is smaller than 1.4 Å, the resampling and clipping of its
98 2D amplitude spectrum will be adjusted so as to give a final spectrum for fitting with an edge
99 corresponding to $1/2.8 \text{ \AA}^{-1}$, to avoid all of the Thon rings being located near the origin of the
100 spectrum, where they can be very poorly sampled. (3) The computation of the quality of fit
101 (CC_{fit} in (Rohou and Grigorieff, 2015)) is now computed over a moving window, similar to
102 (Sheth et al., 2015), rather than at intervals delimited by nodes in the CTF. (4) Following
103 background subtraction as described in (Mindell and Grigorieff, 2003), a radial, cosine-edged
104 mask is applied to the spectrum, and this masked version is used during search and refinement of
105 defocus, astigmatism and phase shift parameters. The cosine is 0.0 at the Fourier space origin,
106 and 1.0 at a radius corresponding to $1/4 \text{ \AA}^{-1}$, and serves to emphasize high-resolution Thon rings,
107 which are less susceptible to artefacts caused by imperfect background subtraction. For all
108 outputs from the program (diagnostic image of the amplitude spectrum, 1D plots, etc.), the
109 background-subtracted, but non-masked, version of the amplitude spectrum is used. (5) Users
110 receive a warning if the box size of the amplitude spectrum and the estimated defocus parameters
111 suggest that significant CTF aliasing occurred (Penczek et al., 2014).

112

113 *Particle picking*

114 Putative particles are found by matching to a soft-edged disk template, which is related to a
115 convolution with Gaussians (Voss et al., 2009) but uses additional statistics based on an
116 algorithm originally described by (Sigworth, 2004). The use of a soft-edged disk template as
117 opposed to structured templates has two main advantages. It greatly speeds up calculation,
118 enabling picking in ‘real time’, and alleviates the problem of templates biasing the result of all
119 subsequent processing towards those templates (Henderson, 2013; Subramaniam, 2013; van
120 Heel, 2013). Any bias that is introduced will be towards a featureless “blob” and will likely be
121 obvious if present.

122 Rather than fully describing the original algorithm by (Sigworth, 2004), we will emphasize here
123 where we deviated from it. The user must specify three parameters: the radius of the template
124 disk, the maximum radius of the particle, which sets the minimum distance between picks, and
125 the detection threshold value, given as a number of standard deviations of the (Gaussian)
126 distribution of scores expected if no particles were present in the input micrograph. Values of 1.0
127 to 6.0 for this threshold generally give acceptable results. All other parameters mentioned below
128 can usually remain set to their default values.

129 Prior to matched filtering, micrographs are resampled by Fourier cropping to a pixel size of 15 Å
130 (the user can override this by changing the “Highest resolution used in picking” value from its
131 default 30 Å), and then filtered with a high-pass cosine-edged aperture to remove very low-
132 frequency density ramps caused by variations in ice thickness or uneven illumination.

133 The background noise spectrum of the micrograph is estimated by computing the average
134 rotational power spectrum of 50 areas devoid of particles, and is then used to “whiten” the
135 background (shot + solvent) noise of the micrograph. Normalization, including CTF effects, and
136 matched filtering are then performed as described (Sigworth, 2004), except using a single
137 reference image and no principal components’ decomposition. When particles are very densely
138 packed on micrographs, this approach can significantly over-estimate the background noise
139 power so that users may find they have to use lower thresholds for picking. It might also be
140 expected that under those circumstances, micrographs with much lower particle density will
141 suffer from a higher rate of false-positive picks.

142 One difficulty in estimating the background noise spectrum of the micrograph is to locate areas
143 devoid of particles without a priori knowledge of their locations. Our algorithm first computes a
144 map of the local variance and local mean in the micrograph (computed over the area defined by
145 the maximum radius given by the user (Roseman, 2004; van Heel, 1982)) and the distribution of
146 values of these mean and variance maps. The average radial power spectrum of the 50 areas of
147 the micrograph with the lowest local variance is then used as an estimate of the background noise
148 spectrum. Optionally, the user can set a different number of areas to be used for this estimate (for
149 example if the density of particles is very high or very low) or use areas with local variances
150 closest to the mode of the distribution of variances, which may also be expected to be devoid of
151 particles.

152 Matched-filter methods are susceptible to picking high-contrast features such as contaminating
153 ice crystals or carbon films. (Sigworth, 2004) suggests subtracting matched references from the
154 extracted boxes and examining the remainder in order to discriminate between real particles and
155 false positives. In the interest of performance, we decided instead to pick using a single artificial

156 reference (disk) and to forgo such subtraction approaches. To avoid picking these kinds of
157 artifacts, the user can choose to ignore areas with abnormal local variance or local mean. We find
158 that ignoring high-variance areas often helps avoid edges of problematic objects, e.g. ice crystals
159 or carbon foils, and that avoiding high- and low-mean areas helps avoid picking from areas
160 within them, e.g. the carbon foil itself or within an ice crystal (Figure 3). The thresholds used are
161 set to $Mo + 2 FWHM$ for the variance and $Mo \pm 2 FWHM$ for the mean, where Mo is the mode
162 (i.e. the most-commonly-occurring value) and $FWHM$ the full width at half-maximum of the
163 distribution of the relevant statistic. For micrographs with additional phase plate phase shifts
164 between 0.1 and 0.9π , where much higher contrast is expected, the variance threshold is
165 increased to $Mo + 8 FWHM$. We have found that in favorable cases many erroneous picks can
166 be avoided. Remaining false-positive picks are removed later during 2D classification.

167 Because of our emphasis on performance, our algorithm can be run nearly instantaneously on a
168 typical $\sim 4K$ image, using a single processor. In the Action panel, the user is presented with an
169 “Auto preview” mode to enable interactive adjustment of the picking parameters (Figure 3). In
170 this mode, the micrograph is displayed with optional and adjustable low-pass and high-pass
171 filters, and the results of picking using the currently selected parameters are overlaid on top.
172 Changing one or more of the parameters leads to a fast re-picking of the displayed micrograph,
173 so that the parameters can be optimized in real-time. Once the three main parameters have been
174 adjusted appropriately, the full complement of input micrographs can be picked, usually in a few
175 seconds or minutes.

176 A possible disadvantage of using a single disk template exists when the particles to be picked are
177 non-uniform in size or shape (e.g. in the case of an elongated particle). In this case, it may be
178 expected that a single template would have difficulty in picking all the different types and views

179 of particles present, and that in this case using a number of different templates would lead to a
180 more accurate picking. In practice, we found that with careful optimization of the parameters,
181 elongated particles and particles with size variation (Figure 3) were picked adequately.

182 The underlying implementation of the algorithm supports multiple references as well as
183 reference rotation. These features may be exposed to the graphical user interface in future
184 versions, for example enabling the use of 2D class averages as picking templates (Scheres,
185 2015).

186

187 *2D classification*

188 2D classification is a relatively quick and robust way to assess the quality of a single-particle
189 dataset. *cisTEM* implements a maximum likelihood algorithm (Scheres et al., 2005) and
190 generates fully CTF-corrected class averages that typically display clear high-resolution detail,
191 such as secondary structure. Integration of the likelihood function is done by evaluating the
192 function at defined angular steps $d\alpha$ that are calculated according to

$$193 \quad d\alpha = R/D \quad (1)$$

194 where R is the resolution limit of the data and D is the diameter of the particle (twice the mask
195 radius that is applied to the iteratively-refined class averages). *cisTEM* runs a user-defined
196 number of iterations n defaulting to 20. To speed up convergence, the resolution limit is adjusted
197 as a function of iteration cycle l ($0 \leq l < n$):

$$198 \quad R = R_{start} + l(R_{finish} - R_{start})/(n - 1) \quad (2)$$

199 where R_{start} and R_{finish} are user-defined resolution limits at the first and last iteration,
 200 defaulting to 40 Å and 8 Å, respectively. The user also sets K , the number of classes to calculate.
 201 Depending on this number and the number of particles N in the dataset, only a percentage p of
 202 the particles are included in the calculation. These particles are randomly reselected for each
 203 iteration and p is typically small, for example 0.1, in the first 10 iterations (p_{0-9}), then increases
 204 to 0.3 for iteration 10 to 14 (p_{10-14}) and finishes with five iterations including all data (p_{15-19}):

$$\begin{aligned}
 205 \quad p_{0-9} &= \begin{cases} 300K/N, & 300K/N < 1 \\ 1, & 300K/N \geq 1 \end{cases} \\
 206 \quad p_{10-14} &= \begin{cases} 0.3, & p_{0-9} < 0.3 \\ p_{0-9}, & p_{0-9} \geq 0.3 \end{cases} \quad (3) \\
 207 \quad p_{15-19} &= 1 .
 \end{aligned}$$

208 For example, for a dataset containing $N = 100,000$ particles, $p_{0-9} = 0.15$, i.e. 15% of the data
 209 will be used to obtain $K = 50$ classes. Apart from speeding up the calculation, the stepwise
 210 increase of the resolution limit and the random selection of subsets of the data also reduce the
 211 chance of overfitting (see also the calculation of ab-initio 3D reconstructions and 3D refinement
 212 below) and, therefore, increase the convergence radius of the 2D classification algorithm.

213 For the calculation of the likelihood function, the particle images \mathbf{X}_i are noise-whitened by
 214 dividing their Fourier transforms $\mathcal{F}\{\mathbf{X}_i\}$ by the square root of the radially average noise power
 215 spectrum, NPS :

$$216 \quad \mathcal{F}\{\tilde{\mathbf{X}}_i\}(\mathbf{g}) = \mathcal{F}\{\mathbf{X}_i\}(\mathbf{g}) / \sqrt{NPS(g)} \quad (4)$$

217 where \mathbf{g} is the 2D reciprocal space coordinate and $g = |\mathbf{g}|$ its magnitude. The noise power
 218 spectrum is calculated from the boxed particle images using the area outside the circular mask

219 set by the user according to the expected particle size. To increase accuracy, it is further
220 averaged across 2000 randomly selected particles. The background (density outside the mask) is
221 further normalized by adding a constant to each particle that yields a background average of
222 zero.

223 Finally, at the beginning of each iteration, noise features in the class averages \mathbf{A}_i are suppressed
224 by resetting negative values below a threshold t_i to the threshold:

$$225 \quad t_i = -0.3 \max_j A_{i,j} \quad (5)$$

226 where j runs over all pixels in average \mathbf{A}_i .

227

228 *3D refinement (FrealignX)*

229 The refinement of 3D reconstructions in *cisTEM* uses a version of Frealign (Lyumkis et al.,
230 2013) that was specifically designed to work with *cisTEM*. Most of Frealign’s control
231 parameters are exposed to the user in the “Manual Refine” Action panel (Figure 4). The “Auto
232 Refine” and “Ab-Initio” panels also use Frealign but manage many of the parameters
233 automatically (see below). Frealign’s algorithm was described previously (Grigorieff, 2007;
234 Lyumkis et al., 2013) and this section will mostly cover important differences, including a new
235 objective function used in the refinement, different particle weighting used in reconstructions,
236 optional likelihood-based blurring, as well as new masking options.

237 **Matched filter** To make Frealign compatible with *cisTEM*’s GUI, the code was completely
238 rewritten in C++, and it will be referred to here as Frealign v10, or FrealignX. The new version
239 makes use of a matched filter (McDonough and Whalen, 1995) to maximize the signal in cross

240 correlation maps calculated between particle images and reference projections. This requires
 241 whitening of the noise present in the images and resolution-dependent scaling of the reference
 242 projections to match the signal in the noise-whitened images. Both can be achieved if the spectral
 243 signal-to-noise ratio (SSNR) of the data is known. As part of a 3D reconstruction, Frealign
 244 calculates the resolution-dependent *PSSNR*, the radially averaged SSNR present in the particle
 245 images before they are affected by the CTF (Sindelar and Grigorieff, 2012). Using *PSSNR* and
 246 the CTF determined for a particle, the SSNR in the particle image can be calculated as

$$247 \quad SNR(\mathbf{g}) = PSSNR(g) \times CTF^2(\mathbf{g}) \quad (6)$$

248 (as before, \mathbf{g} is the 2D reciprocal space coordinate and $g = |\mathbf{g}|$). Here, *SNR* is defined as the
 249 ratio of the variance of the signal and the noise. The Fourier transform $\mathcal{F}\{\tilde{\mathbf{X}}_i\}$ of the noise-
 250 whitened particle image $\tilde{\mathbf{X}}_i$ can then be calculated as

$$251 \quad \mathcal{F}\{\tilde{\mathbf{X}}_i\}(\mathbf{g}) = \frac{\mathcal{F}\{\mathbf{X}_i\}(\mathbf{g})}{\sqrt{|\mathcal{F}\{\mathbf{X}_i\}|_r^2(g)}} \sqrt{1 + SNR(\mathbf{g})} \quad (7)$$

252 where $\mathcal{F}\{\mathbf{X}_i\}$ is the Fourier transform of the original image \mathbf{X}_i , $|\cdot|$ is the absolute value, and
 253 $|\mathcal{F}\{\mathbf{X}_i\}|_r^2$ is the radially averaged spectrum of the squared 2D Fourier transform amplitudes of
 254 image \mathbf{X}_i . To implement Eq. (7), a particle image is first divided by its amplitude spectrum,
 255 which includes power from both signal and noise, and then multiplied by a term that amplifies
 256 the image amplitudes according to the signal strength in the image. The reference projection \mathbf{A}_i
 257 can be matched by calculating

$$258 \quad \mathcal{F}\{\tilde{\mathbf{A}}_i\}(\mathbf{g}) = \frac{\mathcal{F}\{\mathbf{A}_i\}(\mathbf{g})}{\sqrt{|\mathcal{F}\{\mathbf{A}_i\}|_r^2(g)}} \sqrt{SNR(\mathbf{g})}. \quad (8)$$

259 Eq. (8) scales the variance of the signal in the reference to be proportional to the measured
 260 signal-to-noise ratio in the noise-whitened images. The main term in the objective function $O(\phi)$
 261 maximized in FrealignX is therefore given by the cross-correlation function

$$262 \quad CC(\phi) = \frac{Re(\mathcal{F}_{R1,R3}\{\tilde{\mathbf{A}}_i(\phi)\}^* \mathcal{F}_{R1,R3}\{\tilde{\mathbf{X}}_i\})}{\|\mathcal{F}_{R1,R3}\{\tilde{\mathbf{A}}_i(\phi)\}\| \|\mathcal{F}_{R1,R3}\{\tilde{\mathbf{X}}_i\}\|} \quad (9a)$$

263 where ϕ is a set of parameters describing the particle view, x,y position, magnification and
 264 defocus, $Re(\cdot)$ is the real part of a complex number, $\|\cdot\|$ is the Euclidean norm, i.e. the square
 265 root of the sum of the squared pixel values, and $\mathcal{F}_{R1,R3}\{\cdot\}^*$ is the conjugate complex value of the
 266 Fourier transform $\mathcal{F}_{R1,R3}\{\cdot\}$. The subscripts R1 and R3 specify the low- and high-resolution
 267 limits of the Fourier transforms included in the calculation of Eq. (9a), as specified by the user.
 268 To reduce noise overfitting, the user has the option to specify also a resolution range in which the
 269 absolute value of the cross terms in the numerator of Eq. (9a) are used (Grigorieff, 2000; Stewart
 270 and Grigorieff, 2004), instead of the signed values (option ‘‘Signed CC Resolution Limit’’ under
 271 ‘‘Expert Options’’ in the ‘‘Manual Refine’’ Action panel). In this case

$$272 \quad CC(\phi) = \frac{Re(\mathcal{F}_{R1,R2}\{\tilde{\mathbf{A}}_i(\phi)\}^* \mathcal{F}_{R1,R2}\{\tilde{\mathbf{X}}_i\}) + |Re(\mathcal{F}_{R2,R3}\{\tilde{\mathbf{A}}_i(\phi)\}^* \mathcal{F}_{R2,R3}\{\tilde{\mathbf{X}}_i\})|}{\|\mathcal{F}_{R1,R3}\{\tilde{\mathbf{A}}_i(\phi)\}\| \|\mathcal{F}_{R1,R3}\{\tilde{\mathbf{X}}_i\}\|} \quad (9b)$$

273 where R2 is specified by the ‘‘Signed CC Resolution Limit.’’ The objective function also includes
 274 a term $R(\phi|\Theta)$ to restrain alignment parameters (Chen et al., 2009; Lyumkis et al., 2013;
 275 Sigworth, 2004), which currently only includes the x,y positions:

$$276 \quad R(\phi|\Theta) = -\frac{\sigma^2}{M} \left(\frac{(x-\bar{x})^2}{2\sigma_x^2} + \frac{(y-\bar{y})^2}{2\sigma_y^2} \right) \quad (10)$$

277 where σ is the standard deviation of the noise in the particle image and Θ represents a set of
 278 model parameters including the average particle positions in a dataset \bar{x} and \bar{y} , and the standard

279 deviations of the x,y positions from the average values, σ_x and σ_y , and M is the number of pixels
280 in the mask applied to the particle before alignment. The complete objective function is therefore

$$281 \quad O(\phi) = CC(\phi) + R(\phi|\Theta) . \quad (11)$$

282 The maximized values determined in a refinement are converted to particle scores by
283 multiplication with 100.

284 **CTF refinement** FrealignX can refine the defocus assigned to each particle. Given typical
285 imaging conditions with current instrumentation (300 kV, direct electron detector), this may be
286 useful when particles have a size of about 400 kDa or larger. Depending on the quality of the
287 sample and images, these particles may generate sufficient signal to yield per-particle defocus
288 values that are more accurate than the average defocus values determined for whole micrographs
289 by CTFFIND4 (see above). Refinement is achieved by a simple one-dimensional grid search of a
290 defocus offset applied to both defocus values determined in the 2D CTF fit obtained by
291 CTFFIND4. FrealignX applies this offset to the starting values in a refinement, typically
292 determined by CTFFIND4, and evaluates the objective function, Eq. (11), for each offset. The
293 offset yielding the maximum is then used to assign refined defocus values. In a typical
294 refinement, the defocus offset is searched in steps of 50 Å, in a range of ± 500 Å. In the case of
295 β -galactosidase (see below), a single round of defocus refinement changed the defocus on
296 average by 60 Å; the RMS change was 80 Å. For this refinement, the resolution for the signed
297 cross terms equaled the overall refinement resolution limit (3.1 Å), i.e. no unsigned cross terms
298 were used. The refinement produced a marginal improvement of 0.05 Å in the Fourier Shell
299 Correlation (FSC) threshold of 0.143, suggesting that the defocus values determined by
300 CTFFIND4 were already close to optimal. In a different dataset of rotavirus double-layer
301 particles, a single round of defocus refinement changed the defocus on average by 160 Å; the

302 RMS change was 220 Å. In this case, the refinement increased the resolution from ~3.0 Å to
303 ~2.8 Å.

304 **Masking** FREALIGN has a 3D masking function to help in the refinement of structures that
305 contain significant disordered regions, such as micelles in detergent-solubilized membrane
306 proteins. To apply a 3D mask, the user supplies a 3D volume that contains positive and negative
307 values. *cis*TEM will binarize this volume by zeroing all voxels with values less than or equal to
308 zero, and setting all other voxels to 1, indicating the region of the volume that is inside the mask.
309 A soft cosine-shaped falloff of specified width (e.g. 10 Å) is then applied to soften the edge of
310 the masked region and avoid sharp edges when the mask is applied to a 3D reconstruction. The
311 region of the reconstruction outside the mask can be set to zero (simple multiplication of the
312 mask volume), or to a low-pass filtered version of the original density, optionally downweighted
313 by multiplication by a scaling factor set by the user. At the edge of the mask, the low-pass
314 filtered density is blended with the unfiltered density inside the mask to produce a smooth
315 transition. Figure 5 shows the result of masking the reconstruction of an ABC transporter
316 associated with antigen processing (TAP, (Oldham et al., 2016)). The mask was designed to
317 contain only density corresponding to protein and the outside density was low-pass filtered at 30
318 Å resolution and kept with a weight of 100% in the final masked reconstruction. The
319 combination of masking and low-pass filtering in this case keeps a low-pass filtered version of
320 the density outside the mask in the reconstruction, including the detergent micelle. Detergent
321 micelles can be a source of noise in the particle images because the density represents disordered
322 material. However, at low, 20 to 30 Å resolution, micelles generate features in the images that
323 can help in the alignment of the particles. In the case of TAP, this masking prevented noise

324 overfitting in the detergent micelle and helped obtain a reconstruction at 4 Å resolution (Oldham
 325 et al., 2016).

326 **3D reconstruction** In Frealign, a 3D reconstruction \mathbf{V}_k of class average k and containing N
 327 images is calculated as (Lyumkis et al., 2013; Sindelar and Grigorieff, 2012)

$$328 \quad \mathbf{V}_k = \mathcal{F}^{-1} \left\{ \frac{\sum_{i=1}^N \frac{q_{ik}}{\sigma_i^2} \mathcal{R}(\phi_i, w_{ik} \cdot CTF_i \cdot \mathcal{F}\{\hat{\mathbf{X}}_i\})}{\sum_{i=1}^N \frac{q_{ik}}{\sigma_i^2} \mathcal{R}(\phi_i, w_{ik} \cdot CTF_i^2) + 1/PSSNR_k} \right\} \quad (12)$$

329 where q_{ik} is the probability of particle i belonging to class k , σ_i is the standard deviation of the
 330 noise in particle image i , ϕ_i are its alignment parameters, w_{ik} the score-based weights (Eq. (14),
 331 see below), CTF_i the CTF of the particle image, $\mathcal{R}(\phi_i, \cdot)$ the reconstruction operator merging
 332 data into a 3D volume according to alignment parameters ϕ_i , $PSSNR$ the radially averaged
 333 particle SSNR derived from the FSC between half-maps (Sindelar and Grigorieff, 2012), $\hat{\mathbf{X}}_i$
 334 noise-whitened image i , and $\mathcal{F}^{-1}\{\cdot\}$ the inverse Fourier transform. For the calculation of the 3D
 335 reconstructions, as well as 3D classification (see below) the particle images are not whitened
 336 according to Eq. (7). Instead, they are whitened using the radially- and particle-averaged power
 337 spectrum of the background around the particles:

$$338 \quad \mathcal{F}\{\hat{\mathbf{X}}_i\}(\mathbf{g}) = \frac{\mathcal{F}\{\mathbf{X}_i\}(\mathbf{g})}{\sqrt{|\mathcal{F}\{B(\mathbf{X}_i)\}|_r^2(g)}} \quad (13)$$

339 where $B(\mathbf{X}_i)$ is a masked version of image \mathbf{X}_i with the area inside a circular mask centered on the
 340 particle replaced with the average values at the edge of the mask, and scaled variance to produce
 341 an average pixel variance of 1 in the whitened image $\hat{\mathbf{X}}_i$. Using the procedure in Eq. (13) has the
 342 advantage that whitening does not depend on the knowledge of the SSNR of the data, and
 343 reconstructions can therefore be calculated even when the SSNR is not known.

344 **Score-based weighting** In previous versions of Frealign, resolution-dependent weighting was
 345 applied to the particle images during reconstruction (the Frealign parameter was called “PBC”,
 346 (Grigorieff, 2007)). The weighting function took the form of a B-factor dependent exponential
 347 that attenuates the image data at higher resolution. FrealignX still uses B-factor weighting but the
 348 weighting function is now derived from the particle scores (see above) as

$$349 \quad w(\text{score}, \mathbf{g}) = e^{-\frac{BSC}{4}(\text{score} - \overline{\text{score}})g^2} . \quad (14)$$

350 *BSC* converts the difference between a particle score and $\overline{\text{score}}$, the score average, into a B-
 351 factor. Setting *BSC* to zero will turn off score-based particle weighting. Scores typically vary by
 352 about 10, and values for *BSC* that produce reasonable discrimination between high-scoring and
 353 low-scoring particles are between 2 and 10 Å², resulting in B-factor differences between particles
 354 of 20 to 100 Å².

355 **3D Classification** FrealignX uses a maximum-likelihood approach for 3D classification
 356 (Lyumkis et al., 2013). Assuming that all images were noise-whitened according to Eq. (13),
 357 which scales the variance of each image such that the average standard deviation of the noise in a
 358 pixel is 1, the probability density function (PDF) of observing image \mathbf{X}_i , given alignment
 359 parameters ϕ_i and reconstruction \mathbf{V}_k , is calculated as (Lyumkis et al., 2013)

$$360 \quad \Gamma(\mathbf{X}_i | \phi_{ik}, \mathbf{V}_k) = \left(\frac{1}{2\pi}\right)^{\tilde{M}} \exp\left[-\frac{\|\hat{\mathbf{X}}_i - CTF_i \cdot \wp(\mathbf{V}_k, \phi_{ik})\|_{\tilde{M}}^2}{2}\right] \gamma(\phi_{ik} | \Theta_k) . \quad (15)$$

361 As before, ϕ_{ik} are the alignment parameters (usually just Euler angles and x,y shifts) determined
 362 for image i with respect to class average k , \wp is the projection operator producing an aligned 2D
 363 projection of reconstruction \mathbf{V}_k according to parameters ϕ_{ik} , $\|\hat{\mathbf{X}}_i - CTF_i \cdot \wp(\mathbf{V}_k, \phi_{ik})\|_{\tilde{M}}^2$ is the
 364 sum of the squared pixel value differences between whitened image $\hat{\mathbf{X}}_i$ and the reference

365 projection inside a circular mask defining the area of the particle with user-defined diameter, \tilde{M}
366 is the number of pixels inside this mask, and $\gamma(\phi_{ik}|\Theta_k)$ is a hierarchical prior describing the
367 probability of observing alignment parameters ϕ_{ik} given model parameters Θ_k (see Eq. (10)).
368 Eq. (15) does not include marginalization over alignment parameters. Marginalization could be
369 added to improve classification when particle alignments suffer from significant errors.
370 However, this is currently not implemented in *cis*TEM. Given the joint probability, Eq. (15),
371 determined in a refinement, the probability q_{ik} of particle i belonging to class k can be updated
372 as (Lyumkis et al., 2013)

$$373 \quad q_{ik} = \frac{\Gamma(\mathbf{X}_i|\Theta_{ik}, \mathbf{V}_k)\pi_k}{\sum_{k=1}^K \Gamma(\mathbf{X}_i|\Theta_{ik}, \mathbf{V}_k)\pi_k} \quad (16)$$

374 where the summation in the denominator is taken over all classes and the average probabilities
375 π_k for a particle to belong to class k are given by the average values of q_{ik} determined in a prior
376 iteration, calculated for the entire dataset of N particles:

$$377 \quad \pi_k = \frac{1}{N} \sum_{i=1}^N q_{ik} . \quad (17)$$

378 An example of 3D classification is shown in Figure 6 for F_1F_0 -ATPase, revealing different
379 conformational states of the γ subunit (Zhou et al., 2015).

380 **Focused classification** 3D classification can be improved by focusing on conformationally- or
381 compositionally-variable regions of the map. To achieve this, a mask is applied to the particle
382 images and reference projections, the area of which is defined as the projection of a sphere with
383 user-specified center (within the 3D reconstruction) and radius. This 2D mask is therefore
384 defined independently for each particle, as a function of its orientation. When using focused
385 classification, \tilde{M} in Eq. (15) is adjusted to the number of pixels inside the projected mask and the

386 sum of the squared pixel value differences in Eq. (15) is limited to the area of the 2D mask. By
 387 applying the same mask to image and reference, only variability inside the masked region is used
 388 for 3D classification. Other regions of the map are ignored, leading to a “focusing” on the region
 389 of interest. The focused mask also excludes noise contained in the particle images outside the
 390 mask and therefore improves classification results that often depend on detecting small
 391 differences between particles and references. A typical application of a focused mask is in the
 392 classification of ribosome complexes that may exhibit localized conformational and/or
 393 compositional variability, for example the variable conformations of an IRES (Abeyrathne et al.,
 394 2016) or different states of tRNA accommodation (Loveland et al., 2017).

395 **Likelihood-based blurring** In some cases, the convergence radius of refinement can be
 396 improved by blurring the reconstruction according to a likelihood function. This procedure is
 397 similar to the maximization step in a maximum likelihood approach (Scheres, 2012). The
 398 likelihood-blurred reconstruction is given by

$$399 \quad \mathbf{V}_k^n = \frac{\sum_{i=1}^N \frac{1}{\sigma_i^2} \int \phi_{\alpha xy} \Gamma(\mathbf{X}_i | \phi_i, \mathbf{V}_k^{n-1}) \mathcal{R}(\phi_i, w_i \cdot CTF_i; \mathbf{X}_i) d\phi_{\alpha xy}}{\sum_{i=1}^N \frac{q_i k_i}{\sigma_i^2} \mathcal{R}(\phi_i, w_i \cdot CTF_i^2) + 1/PSSNR_k} \quad (18)$$

400 where, in the case of FREALIGNX, $\phi_{\alpha xy}$ only includes the x,y particle positions and in-plane
 401 rotation angle α , which are a subset of the alignment parameters ϕ_i , and \mathbf{V}_k^{n-1} is the
 402 reconstruction from an earlier refinement iteration. As before, $\Gamma(\mathbf{X}_i | \phi_i, \mathbf{V}_k^{n-1})$ is the probability
 403 of observing image i , given alignment parameters ϕ_i and reconstruction \mathbf{V}_k^{n-1} . Integration over
 404 these three parameters can be efficiently implemented and, therefore, does not produce a
 405 significant additional computational burden.

406 **Resolution assessment** The resolution of reconstructions generated by FrealignX is assessed
407 using the FSC criterion (Harauz and van Heel, 1986) using the 0.143 threshold (Rosenthal and
408 Henderson, 2003). FSC curves in *cis*TEM are calculated using two reconstructions (“half-maps”)
409 calculated either from the even-numbered and odd-numbered particles, or by dividing the dataset
410 into 100 equal subsets and using the even- and odd-numbered subsets to calculate the two
411 reconstructions (in the *cis*TEM GUI, the latter is always used). The latter method has the
412 advantage that accidental duplication of particles in a stack is less likely to affect the FSC
413 calculation. All particles are refined against a single reference and, therefore, the calculated FSC
414 values may be biased towards higher values (Grigorieff, 2000; Stewart and Grigorieff, 2004).
415 This bias extends slightly beyond the resolution limit imposed during refinement, by
416 approximately $2/D_{mask}$, where D_{mask} is the mask radius used to mask the reconstructions (see
417 above). During auto-refinement (see below), the resolution limit imposed during refinement is
418 carefully adjusted to stay well below the estimated resolution of the reconstruction and the
419 resolution estimate is therefore unbiased (Scheres and Chen, 2012). However, users have full
420 control over all parameters during manual refinement and will have to make sure that they do not
421 bias the resolution estimate by choosing a resolution limit that is close to, or higher than, the
422 estimated resolution of the final reconstruction. Calculated FSC curves are smoothed using a
423 Savitzky–Golay cubic polynomial that reduces the noise often affecting FSC curves at the high-
424 resolution end.

425 The FSC calculated between two density maps is dependent on the amount of solvent included
426 inside the mask applied to the maps. A larger mask that includes more solvent background will
427 yield lower FSC values than a tighter mask. To obtain an accurate resolution estimate in the
428 region of the particle density, one possibility is to apply a tight mask that closely follows the

429 boundary of the particle. This approach bears the risk of generating artifacts because the particle
 430 boundary is not always well defined, especially when the particle includes disordered domains
 431 that generate weak density in the reconstruction. The approach in Frealign avoids tight masking
 432 and instead calculates an FSC curve using generously masked density maps, corrected for the
 433 solvent content inside the mask (Sindelar and Grigorieff, 2012). The corrected FSC curve is
 434 referred to as *Part_FSC* and is calculated from the uncorrected FSC_{uncor} as (Oldham et al.,
 435 2016)

$$436 \quad Part_FSC_{half-maps} = \frac{fFSC_{uncor}}{1+(f-1)FSC_{uncor}}, \quad (19)$$

437 where f is the ratio of mask volume to estimated particle volume. The particle volume can be
 438 estimated from its molecular mass M_w as $\frac{\text{\AA}^3}{0.81\text{Da}} M_w$ (Matthews, 1968). FSC curves obtained with
 439 the generous masking and subsequent solvent correction yield resolution estimates that are very
 440 close to those obtained with tight masking (Fig. 7C). Eq. (19) assumes that both maps have
 441 similar SSNR values, as is normally the case for the two reconstructions calculated from two
 442 halves of the dataset, indicated by the subscript *half – maps*. If one of the maps does not
 443 contain noise from solvent background, for example when calculating the FSC between a
 444 reconstruction and a map derived from an atomic model, the solvent-corrected FSC is given as

$$445 \quad Part_FSC_{model-map} = \sqrt{\frac{fFSC_{uncor}^2}{1+(f-1)FSC_{uncor}^2}}. \quad (20)$$

446 **Speed optimization** FrealignX has been optimized for execution on multiple CPU cores. Apart
 447 from using optimized library functions for FFT calculation and vector multiplication (Intel Math
 448 Kernel Library), the processing speed is also increased by on-the-fly cropping in real and
 449 reciprocal space of particle images and 3D reference maps. Real-space cropping reduces the

450 interpolation accuracy in reciprocal space and is therefore limited to global parameter searches
451 that do not require the highest accuracy in the calculation of search projections. Reciprocal-space
452 cropping is used whenever a resolution limit is specified by the user or in an automated
453 refinement (ab-initio 3D reconstruction and auto-refinement). For the calculation of in-plane
454 rotated references, reciprocal-space padding is used to increase the image size four-fold,
455 allowing fast nearest-neighbor resampling in real space with sufficient accuracy to produce
456 rotated images with high fidelity.

457

458 *Ab-initio 3D reconstruction*

459 Ab-initio reconstruction offers a convenient way to proceed from single particle images to a 3D
460 structure when a suitable reference is not available to initialize 3D reconstruction and refinement.
461 Different ab-initio methods have been described (Hohn et al., 2007; Punjani et al., 2017; Reboul
462 et al., 2018) and *cisTEM*'s implementation follows a strategy published originally by (Grigorieff,
463 2016). It is based on the premise that iterative refinement of a reconstruction initialized with
464 random angular parameters is likely to converge on the correct structure if overfitting is avoided
465 and the refinement proceeds in small steps to reduce the chance of premature convergence onto
466 an incorrect structure. The procedure is implemented as part of *cisTEM*'s GUI and uses
467 *FrealignX* to perform the refinements and reconstructions.

468 After initialization with random angles, *cisTEM* performs a user-specified number of global
469 alignment parameter searches, recalculating the reconstruction after each search and applying an
470 automatic masking procedure to it before the next global search. Similar to 2D classification (see
471 above), only a randomly selected subset of the data is used in each iteration and the resolution

472 limit applied during the search is increased with every iteration. The number of iterations n
473 defaults to 40, the starting and final resolution limits R_{start} and R_{finish} default to 20 Å and 8 Å,
474 respectively, and the starting and final percentage of included particles in the reconstruction,
475 p_{start} and p_{finish} default to $2500K/N$ and $10,000K/N$, respectively (results larger than 1 are
476 reset to 1), with K the number of 3D classes to be calculated as specified by the user, and N the
477 number of particles in the dataset. If symmetry is applied, N is replaced by NO_{sym} where O_{sym} is
478 the number of asymmetric units present in one particle. The resolution limit is then updated in
479 each iteration l as in Eq. (2), and the percentage is updated as

$$480 \quad p = p_{start} + l(p_{finish} - p_{start})/(n - 1), \quad (21)$$

481 again resetting results larger than 1 to 1. *cis*TEM actually performs a global search for a
482 percentage $3p$ of the particle stack, i.e. three times as many particles as are included in the
483 reconstructions for each iteration. The particles included in the reconstructions are then chosen to
484 be those with the highest scores as calculated by FrealignX.

485 The global alignment parameters are performed using the “general” FrealignX procedure with
486 the following changes. Firstly, the *PSSNR* is not directly estimated from the FSC calculated at
487 each round. Instead, for the first 3 iterations, a default *PSSNR* is calculated based on the
488 molecular weight. From the 4th iteration onwards, the *PSSNR* is calculated from the FSC,
489 however if the calculated *PSSNR* is higher than the default *PSSNR*, the default *PSSNR* is taken
490 instead. This is done in order to avoid some of the overfitting that will occur during refinement.
491 Secondly, during a normal global search the top h (where h defaults to 20) results of the grid
492 search are locally refined, and the best locally refined result is taken. In the ab-initio procedure,
493 however, the result of the global search for a given particle image is taken randomly from all

494 results that have a score which lies in the top 15% of the difference between the worst score and
495 the best score.

496 During the reconstruction steps, the calculated σ for each particle is reset to 1 prior to 3D
497 reconstruction and score weighting is disabled. This is done because the σ and score values are
498 not meaningful until an approximately correct solution is obtained.

499 The reconstructions are automatically masked before each new refinement iteration to suppress
500 noise features that could otherwise be amplified in subsequent iterations. The same masking
501 procedure is also applied during auto-refinement (see below). It starts by calculating the density
502 average $\bar{\rho}$ of the reconstruction and resetting all voxel values below $\bar{\rho}$ to $\bar{\rho}$. This thresholded
503 reconstruction is then low-pass filtered at 50 Å resolution and turned into a binary mask by
504 setting densities equal or below a given threshold t to zero and all others to 1. The threshold is
505 calculated as

$$506 \quad t = \bar{\rho}_{filtered} + 0.03(\bar{\rho}_{max_500} - \bar{\rho}_{filtered}) \quad (22)$$

507 where $\bar{\rho}_{filtered}$ is the density average of the low-pass filtered map and $\bar{\rho}_{max_500}$ is the average of
508 the 500 highest values in the filtered map. The largest contiguous volume in this binarized map is
509 then identified and used as a mask for the original thresholded reconstruction, such that all
510 voxels outside of this mask will be set to $\bar{\rho}$. Finally, a spherical mask, centered in the
511 reconstruction box, is applied by resetting all densities outside the mask to zero.

512 The user has the option to repeat the ab-initio procedure multiple times using the result from the
513 previous run as the starting map in each new run, to increase the convergence radius if necessary.

514 In the case of symmetric particles, the default behavior is to perform the first 2/3^{rds} of the
515 iterations without applying symmetry. The non-symmetrized map is then aligned to the expected

516 symmetry axes and the final 1/3rd of the iterations are carried out with the symmetry applied.
517 This default behavior can be changed by the user such that symmetry is always applied, or is
518 never applied.

519 Alignment of the model to the symmetry axes is achieved using the following process. A brute
520 force grid search over rotations around the x, y and z axes is set up. At each position on the grid
521 the 3D map is rotated using the current x, y and z parameters, and then its projection along Euler
522 angle (0, 0, 0) is calculated. All of the symmetry-related projections are then also calculated, and
523 for each one a cross-correlation map is calculated using the original projection as a reference,
524 and the peak within this map is found. The sum of all peaks from all symmetry related directions
525 is taken and the x,y,z rotation that most closely aligns the original 3D map along the symmetry
526 axes should provide the highest peak sum. To improve robustness, this process is repeated for
527 two additional angles (-45, -45, -45 and 15, 70, -15) that were chosen with the aim of including
528 different-looking areas when the map to be aligned is unusual in some way. The x,y,z rotation
529 that results in the largest sum of all peaks, over all three angles, is taken as the final rotation
530 result. Shifts for this rotation are then calculated based on the found 2D x,y shifts between the
531 initial and symmetry related projections, with the z shift being set to 0 for C symmetries. The
532 symmetry alignment is also included as a command-line program, which can be used to align a
533 volume to the symmetry axis when the ab-initio is carried out in C1 only, or when using a
534 reference obtained by some other means.

535

536 *Automatic refinement*

537 Like ab-initio 3D reconstruction, auto-refinement makes use of randomly selected subsets of the
538 data and of an increasing resolution limit as refinement proceeds. However, unlike the ab-initio
539 procedure, the percentage of particles p_l and the resolution limit R_l used in iteration l depend on
540 the resolution of the reconstructions estimated on iteration $l - 1$. When the estimated resolution
541 improved in the previous cycle,

$$542 \quad p_l = \max[p_R, p_{l-1}] \quad (23)$$

543 with

$$544 \quad p_R = 8000Ke^{75/R_{l-1}^2}/N \quad (24)$$

545 where K is the number of 3D classes to be calculated and N the number of particles in the
546 dataset. As before, if the particle has symmetry, N is replaced by NO_{sym} where O_{sym} is the
547 number of asymmetric units present in one particle. If the calculated p_l exceeds 1, it is reset to 1.
548 The resolution limit is estimated as

$$549 \quad R = FSC_{0.5} - 2/D_{mask} \quad (25)$$

550 where $FSC_{0.5}$ is the point at which the FSC, unadjusted for the solvent within the mask (see
551 above) crosses the 0.5 threshold and D_{mask} is the user-specified diameter of the spherical mask
552 applied to the 3D reference at the beginning of each iteration, and to the half-maps used to
553 calculate the FSC. The term $2/D_{mask}$ accounts for correlations between the two half-maps due
554 to the masking (see above). When the resolution did not improve in the previous iteration,

$$555 \quad p_l = 1.5p_{l-1} \quad (26)$$

556 (reset to 1 if resulting in a value larger than 1). At least five refinement iterations are run and
557 refinement stops when p_l reaches 1 (all particles are included) and there was no improvement in
558 the estimated resolution for the last three iterations.

559 If multiple classes are refined, the resolution limit in Eq. (25) is set independently for each class,
560 however the highest resolution used for classification is fixed at 8 Å. At least nine iterations are
561 run and refinement stops when p_l reaches 1, the average change in the particle occupancy in the
562 last cycle was 1% or less, and there was no improvement in the estimated resolution for the last
563 three iterations.

564 In a similar manner to the ab-initio procedure, σ values for each particle are set to 1 and score
565 weighting is turned off. This is done until the refinement resolution is better than 7 Å, at which
566 point it is assumed the model is of a reasonable quality.

567

568 *Map sharpening*

569 Most single-particle reconstructions require some degree of sharpening that is usually achieved
570 by applying a negative B-factor to the map. *cisTEM* includes a map sharpening tool that allows
571 the application of an arbitrary B-factor. Additionally, maps can be sharpened by whitening the
572 power spectrum of the reconstruction beyond a user-specified resolution (the default is 8 Å). The
573 whitening amplifies terms at higher resolution similar to a negative B-factor but avoids the over-
574 amplification at the high-resolution end of the spectrum that sometimes occurs with the B-factor
575 method due to its exponential behavior. Whitening is applied after masking of the map, either
576 with a hollow spherical mask of defined inner and outer radius, or with a user-defined mask
577 supplied as a separate 3D volume. The masking removes background noise and makes the

578 whitening of the particle density more accurate. Both methods can be combined in *cis*TEM,
579 together with a resolution limit imposed on the final reconstruction. The whitened and B-factor-
580 sharpened map can optionally be filtered with a figure-of-merit curve calculated using the FSC
581 determined for the reconstruction (Rosenthal and Henderson, 2003; Sindelar and Grigorieff,
582 2012).

583

584 *GUI design and workflow*

585 *cis*TEM's GUI required extensive development because it is an integral part of the processing
586 pipeline. GUIs have become more commonplace in cryo-EM software tools to make them more
587 accessible to users (Conesa Mingo et al., 2018; Desfosses et al., 2014; Moriya et al., 2017;
588 Punjani et al., 2017; Scheres, 2012; Tang et al., 2007). Many of the interfaces are designed as so-
589 called wrappers of command-line driven tools, i.e. they take user input and translate it into a
590 command line that launches the tool. Feedback to the user takes place by checking output files,
591 which are also the main repository of processing results, such as movie frame alignments, image
592 defocus measurements and particle alignment parameters. As processing strategies become more
593 complex and the number of users new to cryo-EM grows, the demands on the GUI increase in
594 the quest for obtaining the best possible results. Useful GUI functions include guided user input
595 (so-called wizards) that adjust to specific situations, graphical presentation of relevant results,
596 user interaction with running processes to allow early intervention and make adjustments, tools
597 to manipulate data (e.g. masking), implementation of higher-level procedures that combine more
598 primitive processing steps to achieve specific goals, and a global searchable database that keeps
599 track of all processing steps and result. While some of these functions can be or have been
600 implemented in wrapper GUIs, the lack of control of these GUIs over the data and processes

601 makes a reliable implementation more difficult. For example, keeping track of results from
602 multiple processing steps, some of them perhaps repeated with different parameters or run many
603 times during an iterative refinement, can become challenging if each step produces a separate
604 results file. Communicating with running processes via files can be slow and is sometimes
605 unreliable due to file system caching. Communication via files may complicate the
606 implementation of higher-level procedures, which rely on the parsing of results from the more
607 primitive processing steps.

608 The *cisTEM* GUI is more than a wrapper as it implements some of the new algorithms in the
609 processing pipeline directly, adjusting the input of running jobs as the refinement proceeds. It
610 enables more complex data processing strategies by tracking all results in a single searchable
611 database. All processing steps are run and controlled by the GUI, which communicates with
612 master and slave processes through TCP/IP. *cisTEM* uses an SQL database, similar to Appion
613 (Lander et al., 2009), to store all results (except image files), offers input functions that guide the
614 user or set appropriate defaults, and implements more complex procedures to automate
615 processing where possible. *cisTEM*'s design is flexible to allow execution in many different
616 environments, including single workstations, multiple networked workstations and large
617 computer clusters.

618 User input and the display of results is organized into different panels that make up *cisTEM*'s
619 GUI, each panel dedicated to specific processing steps (for examples, see Figure 1, 3, 4). This
620 design guides users through a standard workflow that most single particle projects follow: movie
621 alignment, CTF determination, particle picking, 2D classification, 3D reconstruction, refinement
622 and classification, and sharpening of the final reconstructions. Three types of panels exist,
623 dealing with Assets, Actions and Results. Assets are mostly data that can be used in processing

624 steps called Actions. They include Movies, Images, Particle Positions and Volumes. One type of
625 Asset, a Refinement Package, defines the data and parameters necessary to carry out refinement
626 of a 3D structure (or a set of structures if 3D classification is done), it contains a particle stack, as
627 well as information about the sample (e.g. particle size and molecular weight) along with
628 parameters for each particle (e.g. orientations and defocus values). Actions comprise the above
629 mentioned workflow steps, with additional options for ab-initio 3D reconstruction, as well as
630 automatic and manual 3D refinement to enable users to obtain the best possible results from their
631 data. The results of most of these Actions are stored in the database and can be viewed in the
632 related Results panels, which display relevant data necessary to evaluate the success of each
633 processing step. The option to sort and select results by a number of different metrics is available
634 in the movie alignment and CTF estimation Results panels. For example images can be sorted /
635 selected based on the CTF fit resolution (Rohou and Grigorieff, 2015). Movie alignment, 3D
636 refinement and reconstruction also produce new Image and Volume Assets, respectively.

637 Importing or generating new Assets is accomplished by wizards that solicit the necessary user
638 input and perform checks were possible to avoid nonsensical input. In the more complex case of
639 creating a new Refinement Package Asset, the wizard allows the specification of input data, for
640 example based on particle picking results or the selection of 2D and 3D classes. Once an Action
641 has been launched, results are displayed as soon as they become available, together with an
642 overall progress bar, giving users an estimate of how long a processing step will take and of
643 whether the results are as expected. If desired, an Action can be aborted and restarted with a
644 different set of parameters, or the Action can be run again after regular termination to test
645 different parameters. In the latter case, all prior results remain accessible and users can specify
646 those to be used for the next step in the workflow. This provides users with the flexibility to pick

647 and choose the best results in cases where different parts of a dataset require different settings to
648 yield optimal results.

649

650 *Parallelization*

651 *cisTEM* uses a home-grown scheme to accelerate processing in parallel environments. Image
652 processing of single-particle data is an embarrassingly parallel problem, i.e. the parallelization of
653 most tasks can be achieved simply by dividing the data to be processed into smaller chunks that
654 are each processed by a separate program thread, without the need for inter-process
655 communication. Only certain steps require merging of data, such as the calculation of a 3D
656 reconstruction from the entire dataset. *cisTEM* parallelizes processing steps by running multiple
657 instances of the same program, each dealing with a subset of the data, then directly
658 communicating with the launched processes over TCP/IP sockets. This enables the *cisTEM* GUI
659 to distribute jobs and receive results in real time. Communication is directed through a
660 “manager” process, which enables jobs to be run on a cluster, while the GUI itself can run on a
661 local workstation.

662 How each copy of the program is run is specified by the user by setting up a “Run Profile”. This
663 profile is a user defined command, or script that will be run to launch the job, and is designed to
664 be flexible to enable the user to set up parallelization in many different environments. For
665 example, users can design profiles to run on multiple machines via SSH, or to submit to a cluster
666 (e.g. using qsub) etc., or even merge the two in a single profile. One disadvantage of this system
667 is that it may be difficult to create profiles for clusters that require many jobs to be submitted
668 using one command.

669 Another advantage of using a home-grown scheme over existing schemes (e.g. MPI) occurs
670 when jobs are run on a multi-node computing cluster. In this case, jobs will complete even if the
671 full number of requested processors is not available. For example, if a user requests 300 CPUs
672 for a processing step but only 100 are available, *cisTEM* launches 300 jobs of which 200 will
673 remain in the job scheduler's queue. Data processing starts immediately with the 100 jobs that
674 are allowed to run and will complete even if the remaining jobs never run. In such a scenario, an
675 MPI-based job could only run when 300 CPUs become available, potentially delaying execution.
676 In the few cases where a step requires merging of an entire dataset, for example in a 3D
677 reconstruction, parallelization is achieved by calculating multiple intermediate 3D
678 reconstructions for subsets of the data, dumping the intermediate reconstructions to disk and
679 merging them after all reconstruction jobs have finished. It can therefore help to designate a fast
680 disk as a scratch disk to allow rapid dumping and reading of the relatively large files (200 MB –
681 5 GB).

682

683 *Benchmarking with β -galactosidase*

684 A high-resolution dataset of β -galactosidase (Bartesaghi et al., 2015) was used to benchmark
685 Relion 2 (Kimanius et al., 2016) and is also used here to illustrate the workflow of *cisTEM* and
686 assess the time for the completion of different processing steps. The entire dataset was
687 downloaded from the EMPIAR database (Iudin et al., 2016) and consists of 1539 movies
688 containing 38 frames, recorded at super-resolution on a K2 Summit camera (Gatan, Inc.,
689 Pleasanton, CA) and stored locally as tif files using LZW compression (the conversion to tiff and
690 compression was performed by *mrc2tif* (Mastronarde and Held, 2017)). The pixel size of the
691 super-resolution frames was 0.3185 Å, and images were binned to a pixel size of 0.75 Å after

692 movie processing. For 2D classification and ab-initio 3D reconstruction, particles were boxed
693 using 384 x 384 pixel boxes. For auto- and manual refinement, the particles were re-boxed into
694 648 x 648 pixel boxes (boxing is part of the creation of Refinement Packages, see above). For all
695 processing steps, a Dell Precision T7910 workstation containing two E5-2699 v4 Xeon
696 processors with a total of 44 cores was used. Processing parameters were left on default settings
697 except for CTF determination, which was performed at 3.5 Å resolution using the movie
698 averages instead of the frames, and particle picking, which used optimized parameters based on
699 previewing a few selected images (Figure 3). The resolution limit during refinement (auto,
700 manual and CTF) never exceeded 3.1 Å. The data were stored on a local SSD Raid 0 disk for fast
701 access. Table 1 lists the timings of the different processing steps using all 44 CPU cores. Results
702 obtained at different points in the workflow are shown in Figure 7.

703

Processing step	Details	Time (hours)
Movie processing	1539 movies, 38 frames, super-resolution	1.1
CTF determination	Using aligned movie average as input	0.1
Particle picking	131,298 particles	0.1
2D classification	50 classes, 28 selected with 119,523 particles	0.8
Ab initio 3D reconstruction	40 iterations	0.8
Auto refinement	8 iterations, final resolution 2.2 Å	1.4
Manual refinement	1 iteration (incl. defocus), final resolution 2.2 Å	0.4
Total		4.7

704

705 **Table 1** Benchmarking of *cis*TEM using a high-resolution dataset of β -galactosidase (Bartesaghi
706 et al., 2015).

707

708 **Discussion**

709 The implementation of a complete image processing workflow in *cis*TEM offers users a uniform
710 experience and guarantees smooth transitions between processing steps. It also helps developers
711 maintain the software as all the tools and algorithms are developed in-house.

712 The main focus of *cis*TEM is on the processing of single-particle cryo-EM data and high-
713 resolution 3D reconstruction. Future releases of *cis*TEM may include on-the-fly processing of
714 data as it is collected, particle-based movie alignment, support for helical particles, improved 3D
715 masking tools, more reliable resolution and quality indicators, as well as miscellaneous tools
716 such as the determination of the detective quantum efficiency of electron detectors.

717 Since *cis*TEM does not rely on third-party libraries, such as Python, MPI or CUDA, that usually
718 have to be installed and compiled separately on the target system, ready-to-run binaries can be
719 made available for download that are optimized for different architectures. Using the wxWidgets
720 library also means that *cis*TEM can be compiled for different operating systems, including
721 Linux, Windows and OSX. Using a configure script, different options for the fast Fourier
722 transforms (FFTs) can be specified, including the FFTW (<http://www.fftw.org>) and Intel MKL
723 (<http://software.intel.com/en-us/mkl>) libraries. The downloadable binaries are statically linked
724 against the MKL library as this exhibits superior speeds compared to the FFTW library on Intel-
725 based CPUs.

726 While the lack of support for GPUs simplifies the installation and execution of *cis*TEM, it can
727 also be a limitation on workstations that are optimized for GPU-accelerated code. These
728 workstations often do not have many CPU cores and, therefore, *cis*TEM will run significantly
729 more slowly than code that can take advantage of the GPU hardware. Users who would like to
730 run both CPU and GPU-optimized software may therefore have to invest in both types of
731 hardware.

732

733 **Materials and Methods**

734 *Development of cis*TEM

735 The entire *cis*TEM image processing package was written in C++ using the wxWidgets toolkit
736 (<http://wxwidgets.org>) to implement the GUI, as well as the libtiff library (<http://www.libtiff.org>)
737 to support the tiff image format, the SQLite library (<https://sqlite.org>) to implement the SQL
738 database, and Intel's MKL library (<http://software.intel.com/en-us/mkl>) for the calculation of
739 Fourier transforms and vector products. Optionally, *cis*TEM can also be linked against the
740 FFTW library (<http://www.fftw.org>) to replace the MKL library. The code was written and
741 edited using Eclipse (<http://www.eclipse.org>) and GitHub (<http://github.com>) was used for
742 version control.

743

744 *Benchmark dataset*

745 The performance of *cis*TEM was benchmarked using a cryo-EM dataset of β -galactosidase
746 (Bartesaghi et al., 2015), entry EMPIAR-10061 in the EMPIAR database (Iudin et al., 2016).

747

748 *Image and data formats*

749 *cis*TEM stores all image data using the MRC format (Crowther et al., 1996). Additionally,
750 particle parameters can be imported from, and exported to Frealign (Grigorieff, 2016) and Relion
751 (Scheres, 2012).

752

753 **Acknowledgements**

754 The authors are grateful for feedback from early testers of *cis*TEM, including Ruben Diaz-
755 Avalos, Sarah Loerch, Priyanka Abeyrathne, Peter Rickgauer, Ben Himes, Andrei Korostelev,
756 Anna Loveland, Gabriel Demo, Jue Chen, Dmitry Lyumkis, Hiro Furukawa, Wei Lu and Juan
757 Du.

758

759 **Competing Interests**

760 The authors declare no competing interest.

761

762

763 **Figure Captions**

764 **Figure 1** Movie alignment panel of the *cis*TEM GUI. All Action panels provide background
765 information on the operation they control, as well as a section with detailed explanations of all
766 user-accessible parameters. All Action panels also have an Expert Options section that exposes
767 additional parameters.

768 **Figure 2** Thon ring pattern calculated for micrograph “0000” of the high-resolution dataset of β -
769 galactosidase (Bartesaghi et al., 2015) used to benchmark *cis*TEM. The left pattern was
770 calculated from the average of aligned frames while the right pattern was calculated using the
771 original movie with 3-frame sub-averages. The pattern calculated using the movie shows
772 significantly stronger rings compared to the other pattern.

773 **Figure 3** Particle picking panel of the *cis*TEM GUI. The panel shows the preview mode, which
774 allows interactive tuning of the picking parameters for optimal picking. The red circles
775 overlaying the image of the sample indicate candidate particles. The picking algorithm avoids
776 areas of high variance, such as the ice contamination visible in the image.

777 **Figure 4** Manual refinement panel with Expert Options exposed. Most of the parameters needed
778 to run FrealignX can be accessed on this panel. The panel also allows application of a 3D mask,
779 which can be imported as a Volume Asset.

780 **Figure 5** 3D masking with low-pass filtering outside the mask. A) Orthogonal sections through
781 the 3D reconstruction of the transporter associated with antigen processing (TAP), an ABC
782 transporter (Oldham et al., 2016). Density corresponding to the protein, as well as the detergent
783 micelle (n-Dodecyl b-D-maltoside; highlighted with arrows), is visible. B) Orthogonal sections
784 through a 3D mask corresponding to the sections shown in A). The sharp edges of this mask are

785 smoothed before the mask is applied to the map. C) Orthogonal sections through the masked 3D
786 reconstruction. The regions outside the mask are low-pass filtered at 30 Å resolution to remove
787 high-resolution noise from the disordered detergent micelle, but keeping its low-resolution signal
788 to help particle alignment.

789 **Figure 6** 3D classification of a dataset of F₁F₀-ATPase, revealing different conformational states
790 (reproduced from Figure 6A and B in Zhou et al., 2015). Sections through the F₁ domain
791 showing the γ subunit (arrows) in three different states related by 120° rotations are shown on the
792 left. A surface rendering of the map corresponding to State 1a is shown on the right. Scale bars,
793 25 Å.

794 **Figure 7** Processing results of the β -galactosidase dataset (Bartesaghi et al., 2015) used to
795 benchmark *cis*TEM. A) Different stages of the ab-initio reconstruction procedure, starting from a
796 reconstruction from randomly assigned Euler angles. The process takes less than an hour to
797 complete on a high-end CPU-based workstation. B) High-resolution detail of the refined β -
798 galactosidase reconstruction with an average resolution of 2.2 Å, showing sidechain details for
799 most amino acids. C) FSC plots for the refined β -galactosidase reconstruction. The black curve
800 was calculated using a tight mask applied to the half maps (Masked_FSC). An correction for
801 potential masking artifacts (Chen et al., 2013) did not lead to adjustments of this curve. The red
802 curved was calculated with a more generous spherical mask and adjusted for the solvent
803 background within that mask (Part_FSC, Eq. (19)). The resolution limit of 3.1 Å, which was not
804 exceeded during refinement, as well as the FSC = 0.143 threshold are indicated by lines.

805 **Figure 7-source data 1** Source data for the curves shown in Fig. 7C.

806 **References**

- 807 Abeyrathne, P.D., Koh, C.S., Grant, T., Grigorieff, N., and Korostelev, A.A. (2016). Ensemble
808 cryo-EM uncovers inchworm-like translocation of a viral IRES through the ribosome. *Elife* 5.
- 809 Bartesaghi, A., Matthies, D., Banerjee, S., Merk, A., and Subramaniam, S. (2014). Structure of
810 beta-galactosidase at 3.2-Å resolution obtained by cryo-electron microscopy. *Proc Natl Acad Sci*
811 *U S A* 111, 11709-11714.
- 812 Bartesaghi, A., Merk, A., Banerjee, S., Matthies, D., Wu, X., Milne, J.L., and Subramaniam, S.
813 (2015). 2.2 Å resolution cryo-EM structure of beta-galactosidase in complex with a cell-
814 permeant inhibitor. *Science* 348, 1147-1151.
- 815 Chen, J.Z., Settembre, E.C., Aoki, S.T., Zhang, X., Bellamy, A.R., Dormitzer, P.R., Harrison,
816 S.C., and Grigorieff, N. (2009). Molecular interactions in rotavirus assembly and uncoating seen
817 by high-resolution cryo-EM. *Proc Natl Acad Sci U S A* 106, 10644-10648.
- 818 Chen, S., McMullan, G., Faruqi, A.R., Murshudov, G.N., Short, J.M., Scheres, S.H., and
819 Henderson, R. (2013). High-resolution noise substitution to measure overfitting and validate
820 resolution in 3D structure determination by single particle electron cryomicroscopy.
821 *Ultramicroscopy* 135, 24-35.
- 822 Cheng, Y., Grigorieff, N., Penczek, P.A., and Walz, T. (2015). A primer to single-particle cryo-
823 electron microscopy. *Cell* 161, 438-449.
- 824 Conesa Mingo, P., Gutierrez, J., Quintana, A., de la Rosa Trevin, J.M., Zaldivar-Peraza, A.,
825 Cuenca Alba, J., Kazemi, M., Vargas, J., Del Cano, L., Segura, J., *et al.* (2018). Scipion web
826 tools: Easy to use cryo-EM image processing over the web. *Protein Sci* 27, 269-275.

827 Crowther, R.A., Henderson, R., and Smith, J.M. (1996). MRC image processing programs. *J*
828 *Struct Biol* *116*, 9-16.

829 Desfosses, A., Ciuffa, R., Gutsche, I., and Sachse, C. (2014). SPRING - an image processing
830 package for single-particle based helical reconstruction from electron cryomicrographs. *J Struct*
831 *Biol* *185*, 15-26.

832 Frank, J., Radermacher, M., Penczek, P., Zhu, J., Li, Y., Ladjadj, M., and Leith, A. (1996).
833 SPIDER and WEB: processing and visualization of images in 3D electron microscopy and
834 related fields. *J Struct Biol* *116*, 190-199.

835 Grant, T., and Grigorieff, N. (2015a). Automatic estimation and correction of anisotropic
836 magnification distortion in electron microscopes. *J Struct Biol* *192*, 204-208.

837 Grant, T., and Grigorieff, N. (2015b). Measuring the optimal exposure for single particle cryo-
838 EM using a 2.6 Å reconstruction of rotavirus VP6. *Elife* *4*, e06980.

839 Grigorieff, N. (2000). Resolution measurement in structures derived from single particles. *Acta*
840 *Crystallogr D Biol Crystallogr* *56*, 1270-1277.

841 Grigorieff, N. (2007). FREALIGN: high-resolution refinement of single particle structures. *J*
842 *Struct Biol* *157*, 117-125.

843 Grigorieff, N. (2016). Frealign: An Exploratory Tool for Single-Particle Cryo-EM. *Methods*
844 *Enzymol* *579*, 191-226.

845 Harauz, G., and van Heel, M. (1986). Exact Filters for General Geometry 3-Dimensional
846 Reconstruction. *Optik* *73*, 146-156.

847 Henderson, R. (2013). Avoiding the pitfalls of single particle cryo-electron microscopy: Einstein
848 from noise. *Proc Natl Acad Sci U S A* *110*, 18037-18041.

849 Hohn, M., Tang, G., Goodyear, G., Baldwin, P.R., Huang, Z., Penczek, P.A., Yang, C., Glaeser,
850 R.M., Adams, P.D., and Ludtke, S.J. (2007). SPARX, a new environment for Cryo-EM image
851 processing. *J Struct Biol* *157*, 47-55.

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

854 Kimanius, D., Forsberg, B.O., Scheres, S.H., and Lindahl, E. (2016). Accelerated cryo-EM
855 structure determination with parallelisation using GPUs in RELION-2. *Elife* *5*.

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

859 Loveland, A.B., Demo, G., Grigorieff, N., and Korostelev, A.A. (2017). Ensemble cryo-EM
860 elucidates the mechanism of translation fidelity. *Nature* *546*, 113-117.

861 Lyumkis, D., Brilot, A.F., Theobald, D.L., and Grigorieff, N. (2013). Likelihood-based
862 classification of cryo-EM images using FREALIGN. *J Struct Biol* *183*, 377-388.

863 Mastronarde, D.N., and Held, S.R. (2017). Automated tilt series alignment and tomographic
864 reconstruction in IMOD. *J Struct Biol* *197*, 102-113.

865 Matthews, B.W. (1968). Solvent content of protein crystals. *J Mol Biol* *33*, 491-497.

866 McDonough, R.N., and Whalen, A.D. (1995). *Detection of Signals in Noise*, 2nd edn (San
867 Diego: Academic Press).

868 McMullan, G., Faruqi, A.R., and Henderson, R. (2016). Direct Electron Detectors. *Methods*
869 *Enzymol* 579, 1-17.

870 McMullan, G., Vinothkumar, K.R., and Henderson, R. (2015). Thon rings from amorphous ice
871 and implications of beam-induced Brownian motion in single particle electron cryo-microscopy.
872 *Ultramicroscopy* 158, 26-32.

873 Mindell, J.A., and Grigorieff, N. (2003). Accurate determination of local defocus and specimen
874 tilt in electron microscopy. *J Struct Biol* 142, 334-347.

875 Moriya, T., Saur, M., Stabrin, M., Merino, F., Voicu, H., Huang, Z., Penczek, P.A., Raunser, S.,
876 and Gatsogiannis, C. (2017). High-resolution Single Particle Analysis from Electron Cryo-
877 microscopy Images Using SPHIRE. *J Vis Exp*.

878 Oldham, M.L., Grigorieff, N., and Chen, J. (2016). Structure of the transporter associated with
879 antigen processing trapped by herpes simplex virus. *Elife* 5.

880 Penczek, P.A., Fang, J., Li, X., Cheng, Y., Loerke, J., and Spahn, C.M. (2014). CTER-rapid
881 estimation of CTF parameters with error assessment. *Ultramicroscopy* 140, 9-19.

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

884 Reboul, C.F., Eager, M., Elmlund, D., and Elmlund, H. (2018). Single-particle cryo-EM-
885 Improved ab initio 3D reconstruction with SIMPLE/PRIME. *Protein Sci* 27, 51-61.

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

888 Roseman, A.M. (2004). FindEM--a fast, efficient program for automatic selection of particles
889 from electron micrographs. *J Struct Biol* 145, 91-99.

890 Rosenthal, P.B., and Henderson, R. (2003). Optimal determination of particle orientation,
891 absolute hand, and contrast loss in single-particle electron cryomicroscopy. *J Mol Biol* 333, 721-
892 745.

893 Scheres, S.H. (2012). RELION: implementation of a Bayesian approach to cryo-EM structure
894 determination. *J Struct Biol* 180, 519-530.

895 Scheres, S.H. (2015). Semi-automated selection of cryo-EM particles in RELION-1.3. *J Struct*
896 *Biol* 189, 114-122.

897 Scheres, S.H., and Chen, S. (2012). Prevention of overfitting in cryo-EM structure determination.
898 *Nat Methods* 9, 853-854.

899 Scheres, S.H., Valle, M., Nunez, R., Sorzano, C.O., Marabini, R., Herman, G.T., and Carazo,
900 J.M. (2005). Maximum-likelihood multi-reference refinement for electron microscopy images. *J*
901 *Mol Biol* 348, 139-149.

902 Sheth, L.K., Piotrowski, A.L., and Voss, N.R. (2015). Visualization and quality assessment of
903 the contrast transfer function estimation. *J Struct Biol* 192, 222-234.

904 Sigworth, F.J. (2004). Classical detection theory and the cryo-EM particle selection problem. *J*
905 *Struct Biol* 145, 111-122.

906 Sindelar, C.V., and Grigorieff, N. (2012). Optimal noise reduction in 3D reconstructions of
907 single particles using a volume-normalized filter. *J Struct Biol* 180, 26-38.

908 Stewart, A., and Grigorieff, N. (2004). Noise bias in the refinement of structures derived from
909 single particles. *Ultramicroscopy* *102*, 67-84.

910 Subramaniam, S. (2013). Structure of trimeric HIV-1 envelope glycoproteins. *Proc Natl Acad*
911 *Sci U S A* *110*, E4172-4174.

912 Subramaniam, S., Earl, L.A., Falconieri, V., Milne, J.L., and Egelman, E.H. (2016). Resolution
913 advances in cryo-EM enable application to drug discovery. *Curr Opin Struct Biol* *41*, 194-202.

914 Tang, G., Peng, L., Baldwin, P.R., Mann, D.S., Jiang, W., Rees, I., and Ludtke, S.J. (2007).
915 EMAN2: an extensible image processing suite for electron microscopy. *J Struct Biol* *157*, 38-46.

916 Thon, F. (1966). Zur Defokussierungsabhängigkeit des Phasenkontrastes bei der
917 elektronenmikroskopischen Abbildung. *Z Naturforschung* *21a*, 476-478.

918 van Heel, M. (1982). Detection of Objects in Quantum-Noise-Limited Images. *Ultramicroscopy*
919 *7*, 331-341.

920 van Heel, M. (2013). Finding trimeric HIV-1 envelope glycoproteins in random noise. *Proc Natl*
921 *Acad Sci U S A* *110*, E4175-4177.

922 van Heel, M., Harauz, G., Orlova, E.V., Schmidt, R., and Schatz, M. (1996). A new generation
923 of the IMAGIC image processing system. *J Struct Biol* *116*, 17-24.

924 Voss, N.R., Yoshioka, C.K., Radermacher, M., Potter, C.S., and Carragher, B. (2009). DoG
925 Picker and TiltPicker: software tools to facilitate particle selection in single particle electron
926 microscopy. *J Struct Biol* *166*, 205-213.

927 Zhang, K. (2016). Gctf: Real-time CTF determination and correction. *J Struct Biol* *193*, 1-12.

928 Zhou, A., Rohou, A., Schep, D.G., Bason, J.V., Montgomery, M.G., Walker, J.E., Grigorieff, N.,
929 and Rubinstein, J.L. (2015). Structure and conformational states of the bovine mitochondrial
930 ATP synthase by cryo-EM. *Elife* 4, e10180.

931



Overview



Assets



Actions



Results



Settings



Align Movies



Find CTF



Find Particles



2D Classify



Ab-Initio 3D



Auto Refine



Manual Refine



Generate 3D



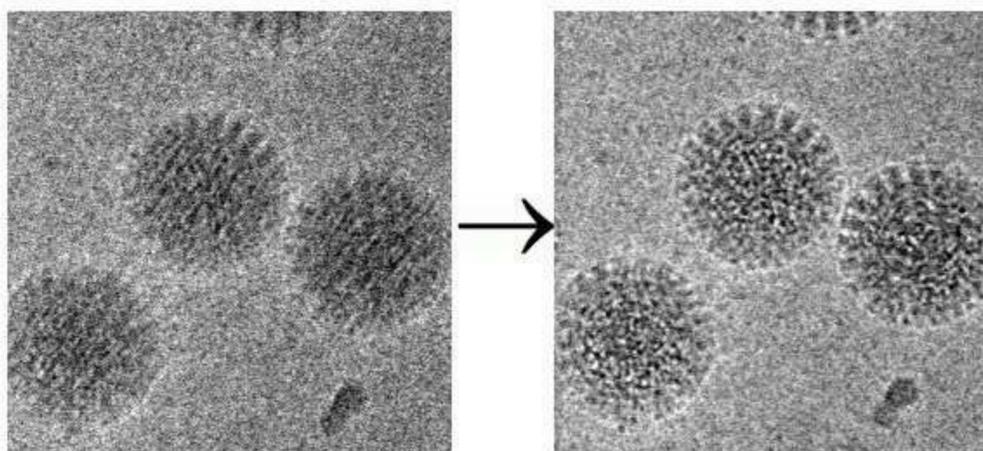
Sharpen 3D

Input Group : All Movies (20)

Show Expert Options

Movie Alignment

Physical drift and beam induced motion (Brilot et al., 2012; Campbell et al., 2012; Li et al., 2013; Scheres, 2014) of the specimen leads to a degradation of information within images, and will ultimately limit the resolution of any reconstruction. Aligning a movie prior to calculating the sum will prevent a large amount of this degradation and lead to better data. This panel therefore attempts to align movies based on the Unblur algorithm described in (Grant and Grigorieff, 2015).



Additionally an exposure weighted sum can be calculated which attempts to maximize the signal-to-noise ratio in the final sums by taking into account the radiation damage the sample has suffered as the movie progresses. This exposure weighting is described in (Grant and Grigorieff, 2015a). During this step, any anisotropic magnification distortions present in the images can also be corrected (Grant and Grigorieff, 2015b), using the parameters that were entered during movie import.

Program Options

Input Group : The group of movie assets that will be aligned. For each movie within the group, the output will be an image representing the sum of the aligned movie. This movie will be automatically added to the image assets list.

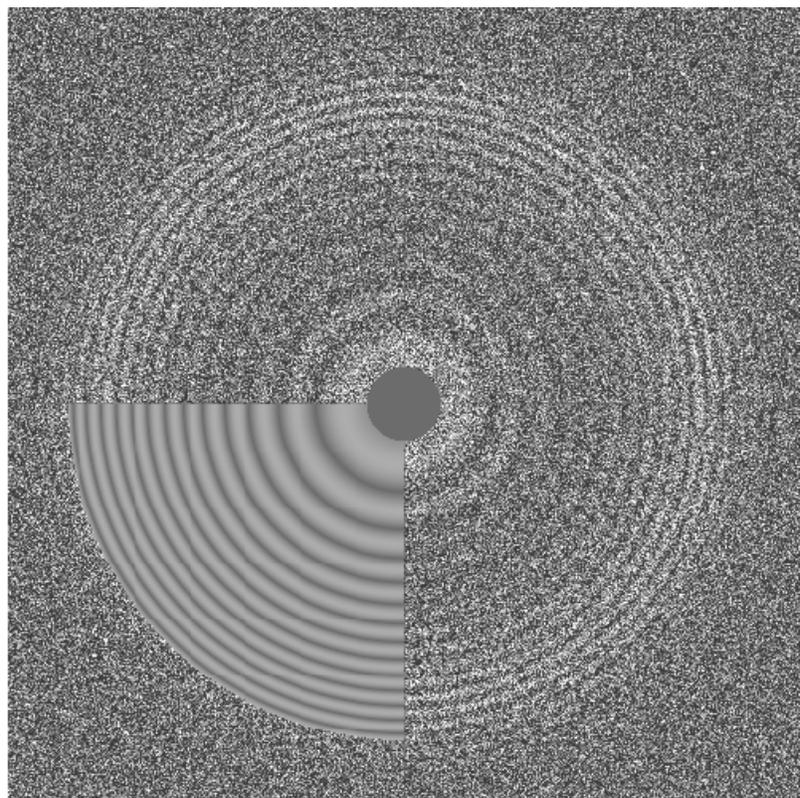
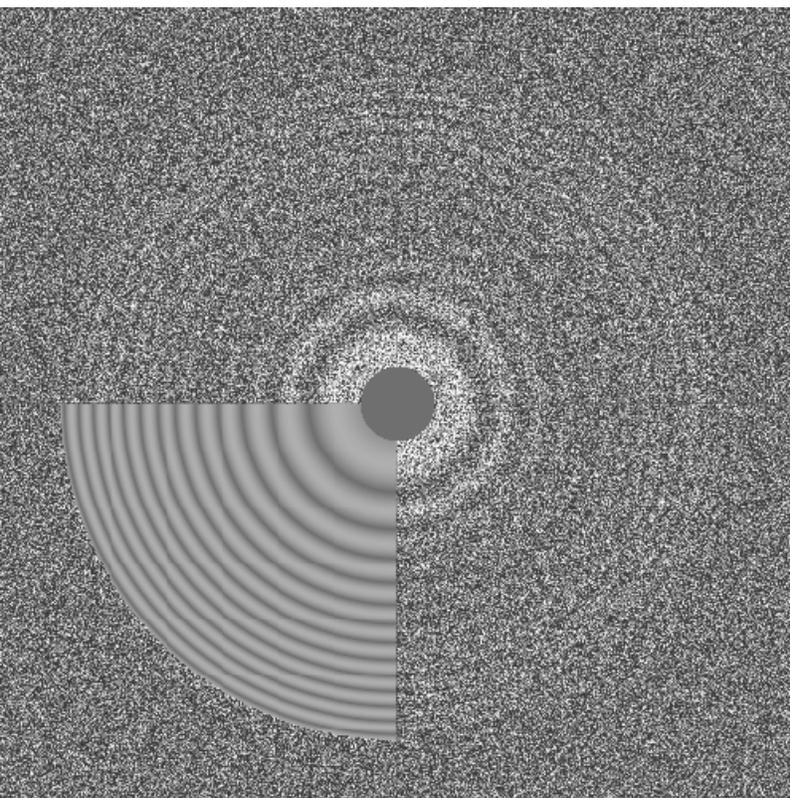
Run Profile : The selected run profile will be used to run the job. The run profile describes how the job should be run (e.g. how many processors should be used, and on which different computers). Run profiles are set in the Run Profile panel, located under settings.

Expert Options

Minimum Shift : This is the minimum shift that can be applied during the initial refinement stage. Its purpose is to prevent images aligning to detector artifacts that may be reinforced in the initial sum which is used as the first reference. It is applied only during the first alignment round, and is ignored after that.

Run Profile : Local 88 (89)

Start Alignment





Overview



Assets



Actions



Results



Settings



Align Movies



Find CTF



Find Particles



2D Classify



Ab-Initio 3D



Auto Refine



Manual Refine



Generate 3D



Sharpen 3D

Input Group : All Images (1539)



Picking algorithm : Ab initio

Show Expert Options

Maximum particle radius (Å) : 60.00

 Circles around particles Scale bar High-pass filter Low-pass filter

20.00

Å

Undo

Redo

Characteristic particle radius (Å) : 50.00

Threshold peak height : 2.8

Select Preview Image :-

EMD-2984_0004_frames_4_0.mrc

Preview

 Auto preview**Expert Options**

Highest resolution used in picking (Å) : 30.00

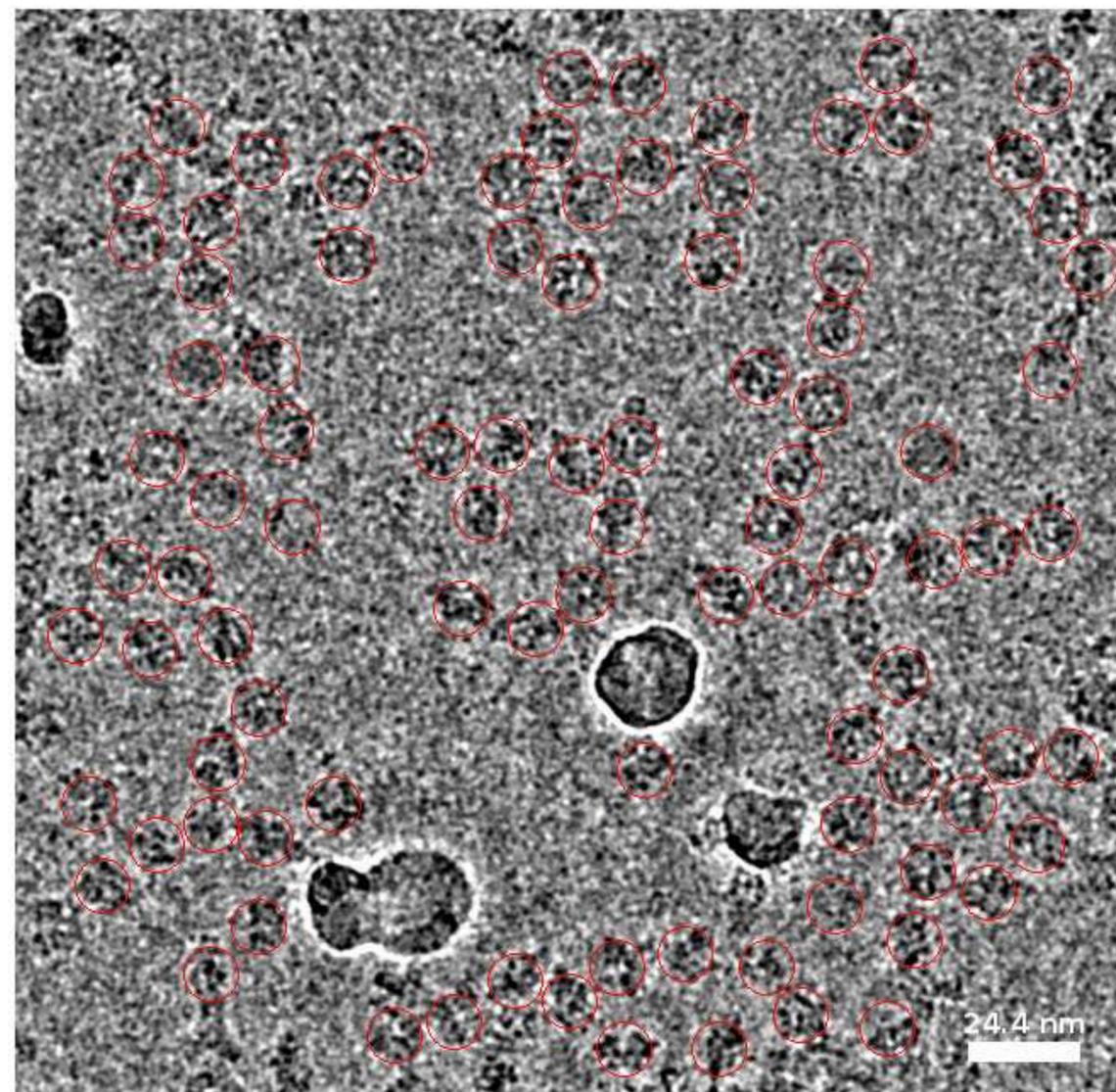
 Set minimum distance from edges (pix.) : 128 Use radial averages of templates Rotate each template this many times : 72 Avoid high variance areas Avoid areas with abnormal local mean Show estimated background spectrum Show positions of background boxes

Number of background boxes : 50

Algorithm to find background areas : Lowest variance

Run Profile : Local (40)

Start Picking





Overview



Assets



Results



Settings



Settings

Align Movies Find CTF Find Particles 2D Classify Ab-Initio 3D Auto Refine **Manual Refine** Generate 3D Sharpen 3D

Input Refinement Package : Refinement Package #3 648 px cleaned

Input Parameters : Imported Parameters

Use a Mask? : mask

Local Refinement Global Search

No. of Cycles to Run : 1

Hi-Res Limit (Å) : 3.10

Show Expert Options

Class No.	Active Reference Volume
1	Generate from params.

Parameters To Refine Reset All Defaults

- Psi
- Theta
- Phi
- X Shift
- Y Shift

General Refinement

Low-Resolution Limit (Å) : 285.00

Outer Mask Radius (Å) : 123.50

Inner Mask Radius (Å) : 0.00

Signed CC Resolution Limit (Å) : 0.00

Percent Used (%) : 100.00

Global Search

Global Mask Radius (Å) : 152.00

Number of Results to Refine : 20

Also Refine Input Parameters? Yes No

Angular Search Step (°) : 2.88

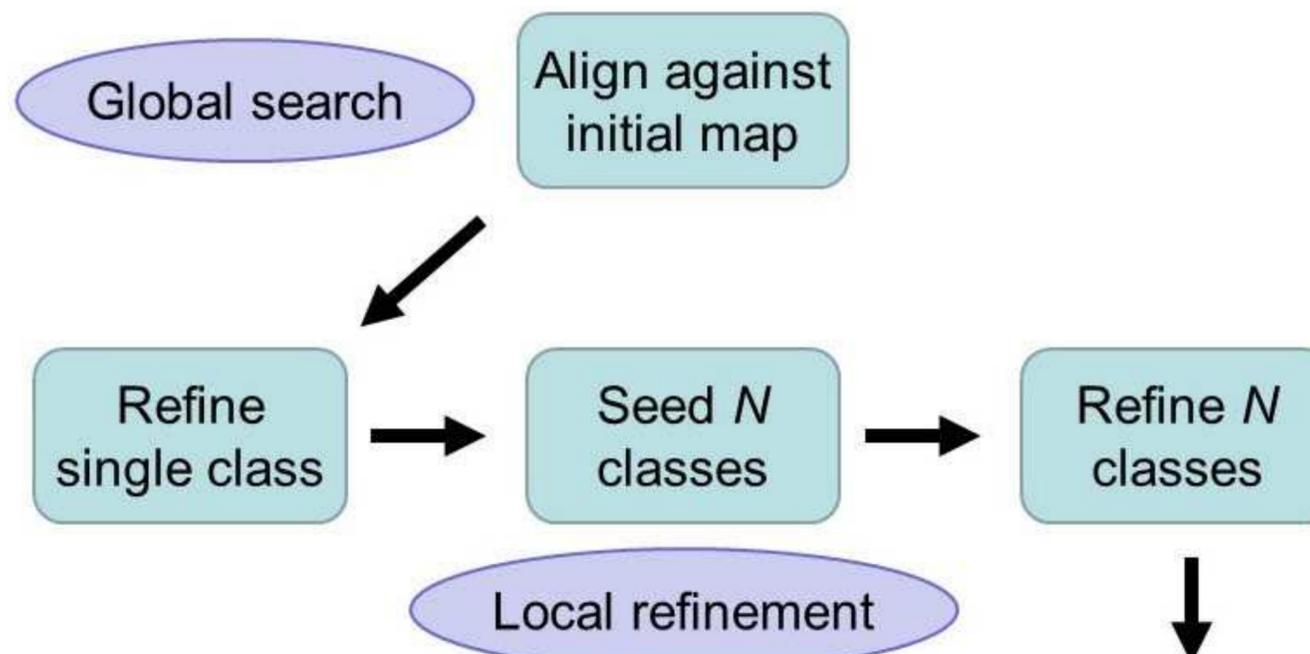
Search Range in X (Å) : 28.50

Refinement Run Profile : Local 88 (89)

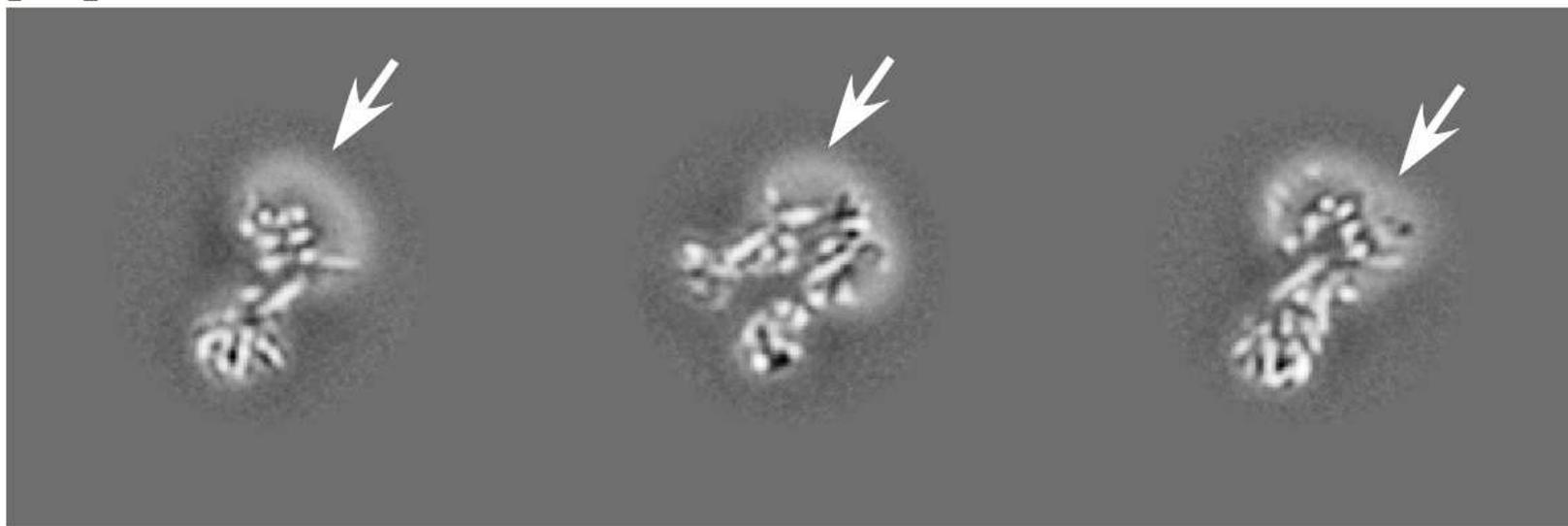
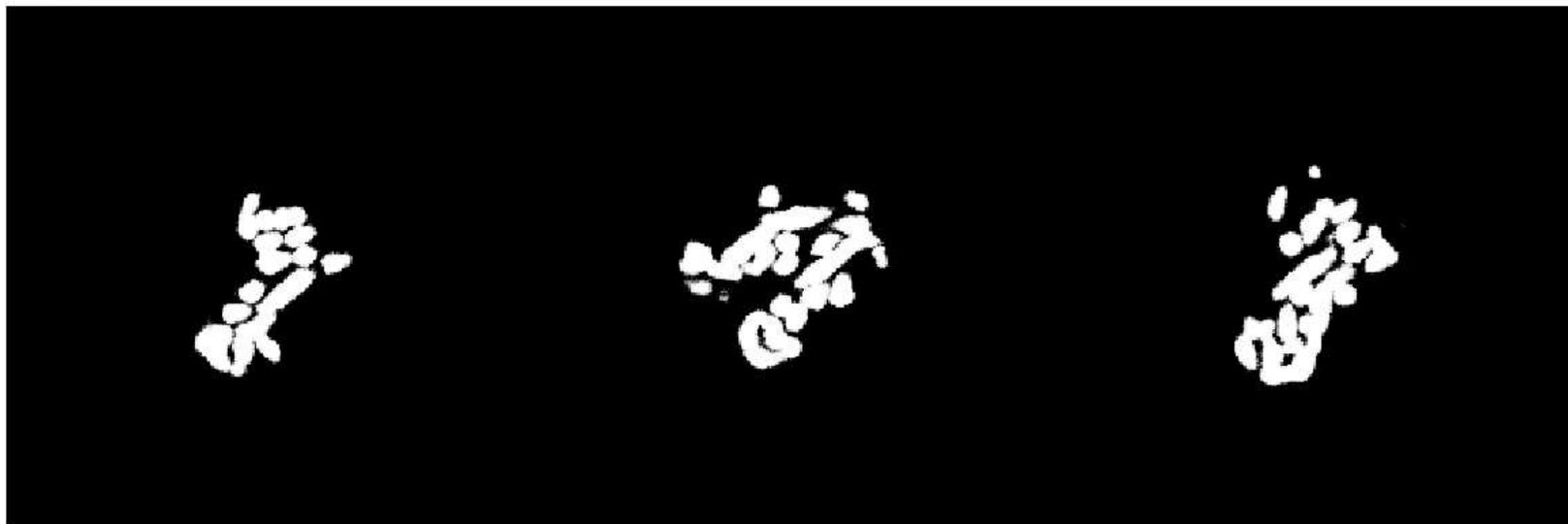
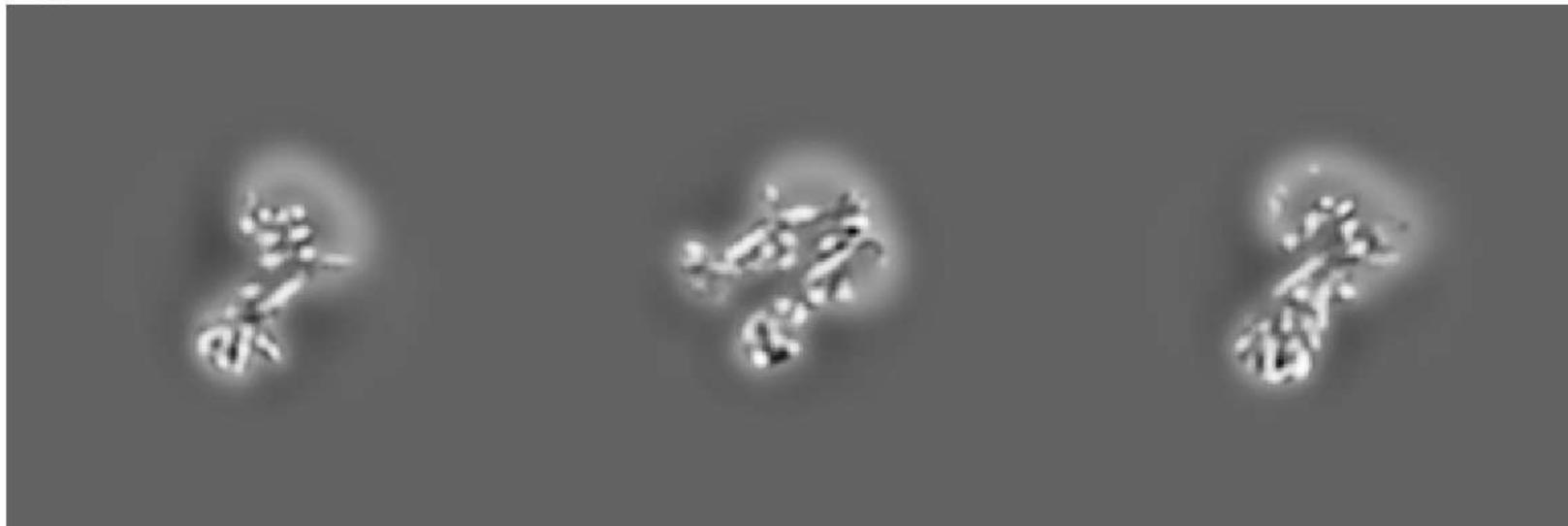
Reconstruction Run Profile : Local 44 (45)

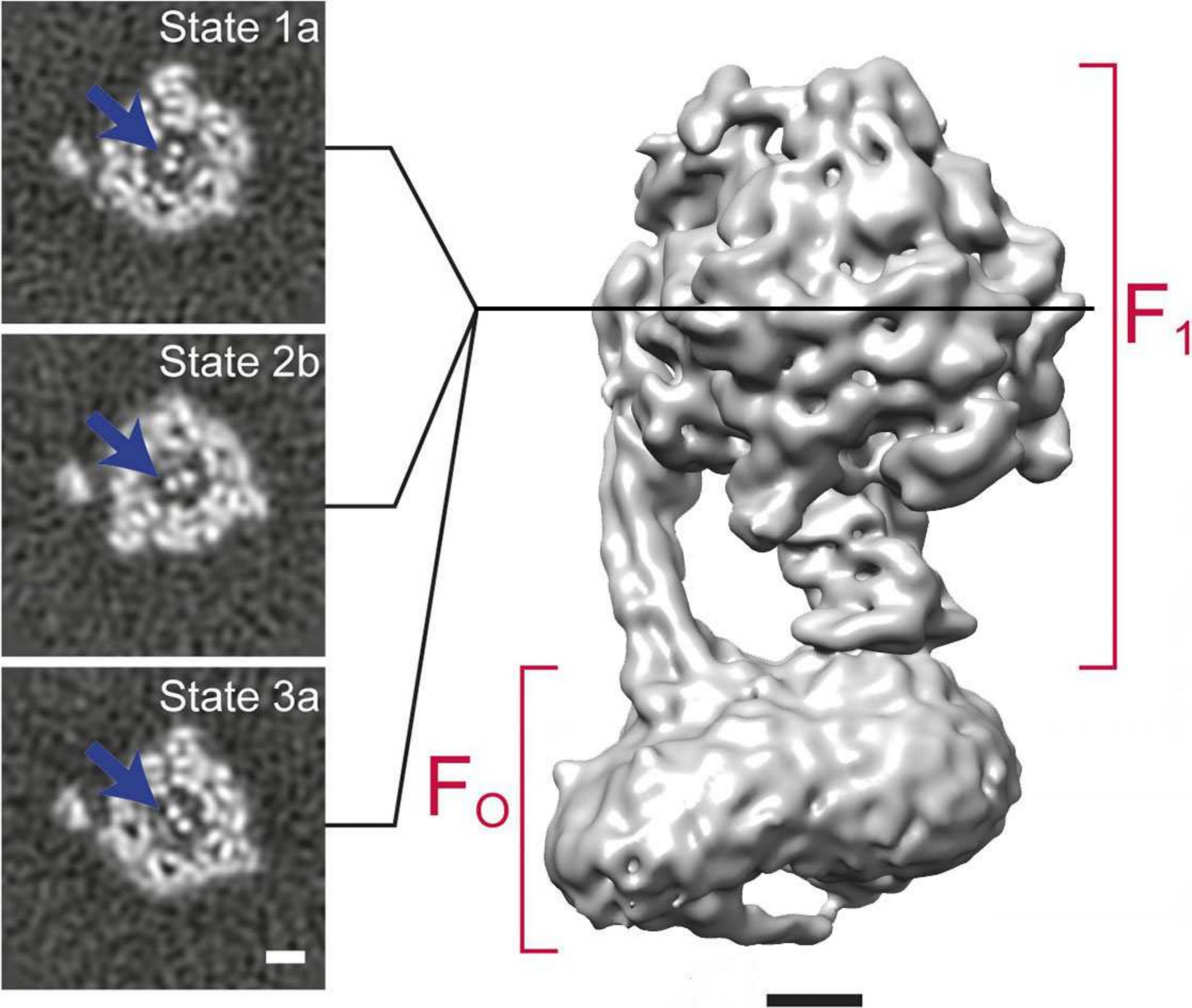
3D Refinement & Reconstruction (FrealignX)

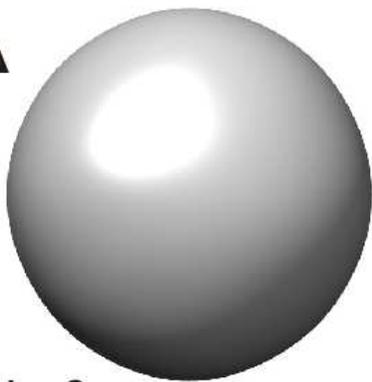
The goal of refinement and reconstruction is to obtain 3D maps of the imaged particle at the highest possible resolution. Refinement typically starts with a preexisting structure that serves as a reference to determine initial particle alignment parameters using a global parameter search. In subsequent iterations, these parameters are refined and (optionally) the dataset can be classified into several classes with distinct structural features. This panel allows the user to define a refinement job that includes a set number of iterations (refinement cycles) and number of desired classes to be generated (Lyumkis et al. 2013). The general refinement strategies and options are similar to those available with Frealign and are described in Grigorieff, 2016:



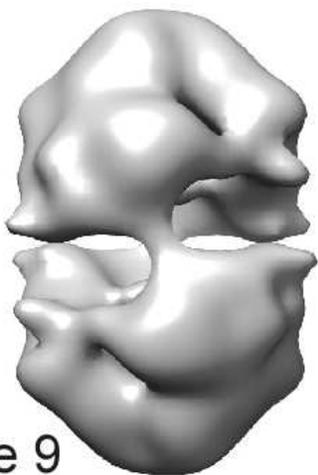
Start Refinement

A**B****C**

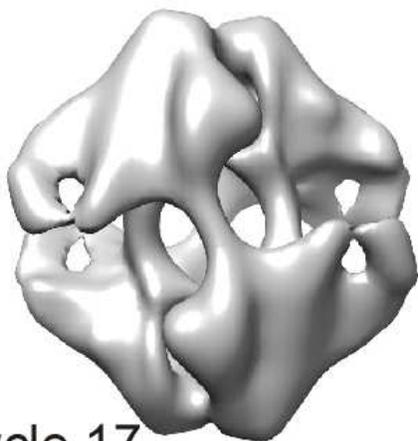


A

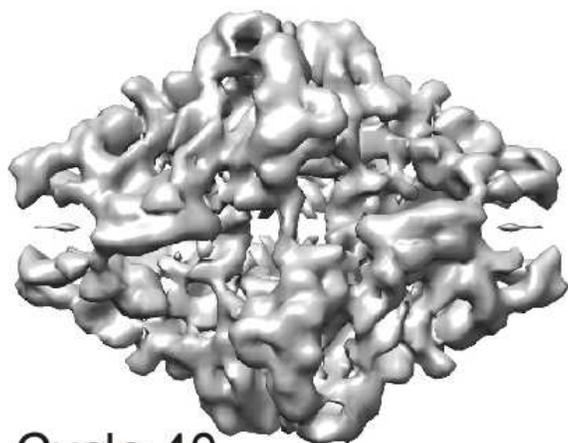
Cycle 0



Cycle 9



Cycle 17



Cycle 40

B

Asn 414

Ala 413

Glu 412

Asp 411

Val 410

Val 409

Tyr 408

C