



Published in final edited form as:

*Nat Methods*. 2018 November ; 15(11): 955–961. doi:10.1038/s41592-018-0167-z.

## emClarity: Software for High Resolution Cryo-electron Tomography and Sub-tomogram Averaging

Benjamin A. Himes<sup>1,4</sup> and Peijun Zhang<sup>1,2,3,\*</sup>

<sup>1</sup>Department of Structural Biology, University of Pittsburgh School of Medicine, Pittsburgh, PA 15260, USA

<sup>2</sup>Division of Structural Biology, Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, OX3 7BN, UK

<sup>3</sup>Electron Bio-Imaging Centre, Diamond Light Source, Harwell Science and Innovation Campus, Didcot OX11 0DE, UK

<sup>4</sup>Janelia Research Campus, Howard Hughes Medical Institute, Ashburn VA, 20147, USA

### Abstract

Macromolecular complexes are intrinsically flexible and often challenging to purify for structure determination by single particle cryoEM. Such complexes may be studied using cryo-electron tomography combined with sub-tomogram alignment and classification, which in exceptional cases reaches sub-nanometer resolution, yielding insight into structure-function relationships. Extending this approach to specimens that exhibit conformational or compositional heterogeneity, and that may be present at low abundance, remains challenging. To address this challenge, we developed emClarity (<https://github.com/bHimes/emClarity/wiki>), a GPU-accelerated image processing package, which features an iterative tomographic tilt-series refinement algorithm using sub-tomograms as fiducial markers and a 3D-sampling function compensated, multi-scale Principle Component Analysis classification method. We demonstrate substantial improvements in the resolution of maps and in the separation of different functional states of macromolecular complexes, compared to those generated using current state-of-the-art software.

---

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:[http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

\*To whom correspondence may be addressed: [peijun@strubi.ox.ac.uk](mailto:peijun@strubi.ox.ac.uk).

#### Author contributions

This study was conceived and designed by P.Z and B.A.H. B.A.H. developed and tested the code for emClarity. B.A.H. and P.Z. analyzed the results. B.A.H. and P.Z. wrote the paper.

#### Competing financial interests:

The authors declare no competing financial interests.

#### Data Availability

Cryo-EM structural data have been deposited in the EM Data Bank under accession codes EMD-8799 for the yeast 80s ribosome, EMD-8802, EMD-8803, EMD-8804, EMD-8805, and EMD-8806 for rabbit 80s ribosome classes I-V respectively, EMD-8986 for the HIV-1 Gag data.

#### Code availability

The software is freely available from <https://www.github.com/bHimes/emClarity> and as Supplementary Software. A tutorial, documentation and videos are available at <https://www.github.com/bHimes/emClarity/wiki>.

## Introduction

Recent advances in the capabilities of cryo-electron microscopy (cryoEM) have enabled determination of structures approaching atomic resolution. Software tools for the single particle analysis (SPA) of cryoEM data have been developed to probe macromolecular functional dynamics; for example, the maximum-likelihood approach to classification as implemented in *cisTEM*<sup>1</sup> or RELION<sup>2</sup>. For a sample to be suitable for SPA, it must yield tens to hundreds of thousands<sup>3</sup> of particles, purified to high compositional and conformational homogeneity<sup>4</sup> and subsequently imaged in many different orientations. These conditions are often difficult to realize, especially for the large and dynamic assemblies of biological complexes most relevant to cellular activities<sup>5</sup>. An alternative to SPA, cryo-electron tomography (cryoET), can be used to generate three-dimensional (3D) reconstructions of the specimen.

These reconstructions (tomograms) are typically limited to ~3–4 nm resolution due to the low electron dose used to prevent excessive radiation damage to samples<sup>6</sup>. Additionally, the signal-to-noise ratio (SNR) is not distributed evenly in the tomogram, resulting in anisotropic resolution. This is a consequence of both the increased specimen thickness at high tilt angles and the restricted angular sampling ( $\pm 60^\circ$ ) known as the “missing-wedge effect”<sup>7</sup>. However, when many copies of a macromolecule are present in a tomogram, they may be extracted *in silico*, aligned to a common orientation, and averaged using procedures that share many similarities to SPA, which ameliorates the distortions due to the “missing-wedge effect”. One major difference compared to SPA is that the metrics used in the orientation search are biased by the missing-wedge. This bias may be compensated by only considering the regions in Fourier space where both volumes being compared have been measured; the two most common being the constrained cross correlation<sup>8,9</sup> and the constrained Euclidean distance<sup>10</sup>.

Obtaining averages of sub-tomograms at resolutions of 15–20 Å is now relatively routine, and is a prerequisite for further classification of the sample into multiple biological states or functional conformations<sup>11,12</sup>. Compared to SPA, this is arguably the greatest strength of cryoET and sub-tomogram averaging, because each particle exists as a unique, albeit distorted, 3D reconstruction. This allows for a direct analysis of the 3D variance, the value of which has been discussed extensively<sup>13,14</sup>. Despite substantial progress over the last decade, very few structures have been solved at resolutions better than 8 Å using cryoET and sub-tomogram averaging, a critical threshold beyond which flexible fitting approaches to modeling of protein complexes are more reliable<sup>11,15</sup>.

We present here a complete set of GPU-accelerated programs called **emClarity** for **enhanced macromolecular classification and alignment for high-resolution tomography**, with the aim of routinely reaching beyond the critical sub-nanometer resolution for diverse specimens. We have focused our efforts on those areas of image processing that are likely to yield the greatest improvements, as suggested by empirical observation and theoretical calculations<sup>16,17</sup>: accuracy of tilt-series alignment, improved defocus determination and CTF correction, explicit treatment of anisotropic resolution, and more robust classification.

## Results

### emClarity workflow

A typical emClarity workflow is illustrated in figure 1, with new features highlighted in red text. Detailed steps are described below. A comparison of features in available sub-tomogram processing packages, including emClarity, PEET, RELION, Dynamo, Jsubtomo, pyTom and Protomo is presented in Supplementary Table 1. Run times for each major operation are presented in Supplementary Table 2, and run times and speed up factors are compared for several popular Nvidia GPUs in Supplementary Figure 1.

**I) Input data**—The raw data for emClarity are tilt-series, rather than tomograms, and an initial estimate of the tilt-series alignment parameters—tilt-axis angle, in-plane shifts and rotations, magnification and tilt angle—readily obtained using a separate installation of the IMOD software package<sup>18</sup>. For the defocus determination we use periodogram averaging as described previously<sup>19</sup>, refined by fitting a 2D astigmatic CTF and using a low pass filtered power spectrum for background subtraction<sup>20</sup>.

**II) Tomogram WBP, template matching with 3D interactive editing**—By default, emClarity limits the resolution of the initial tomogram used in template matching to 40 Å to prevent model bias. This also means CTF correction is not needed and traditional weighted back-projection is sufficient to reconstruct these initial tomograms.

Sub-tomogram positions and orientations are chosen by the best scoring locally normalized cross-correlation, as described previously<sup>21</sup> with two main exceptions. First, our algorithm is optimized to fit in GPU memory. Second, we do not use distribution fitting to estimate the number of false-positives, rather we provide a maximum intensity projection and interface with IMOD's 3D model editing tools to enable rapid manual cleaning of the results.

**III) 3D-CTF corrected WBP**—To correct for the defocus gradient along the optical axis (sample thickness), emClarity uses a straightforward version of the “Defocus Gradient Corrected Back Projection”<sup>22</sup>. The approach has recently been validated and re-named “3D-CTF correction”<sup>23,24</sup>. To balance accuracy with practical compute time, emClarity determines the acceptable thickness based on the current resolution and defocus (Online Methods). For each slab of this thickness we whiten the power spectrum<sup>25</sup>, multiply by the determined CTF and filter according to cumulative electron dose<sup>26</sup>. For tilted images, strips of a width corresponding to the current accepted defocus are extracted from the inverse Fourier transform of the full image multiplied by the respective CTF.

**IV) Averaging sub-tomograms, CTF amplitude correction, Anisotropic SSNR weighting**—The iterative alignment procedure in emClarity alternates between averaging the sub-tomograms using the current estimate of their orientations, and performing a missing-wedge constrained cross correlation grid search<sup>27</sup>. In addition to the data volumes, an average is also made of the “3D-sampling function” which is similar to the “weighted 3D CTF model” described in figure 1 of Bharat et al.<sup>10</sup> with the additional consideration of the R-weighting applied during tomogram reconstruction (Online Methods eq 1). *To avoid*

*confusion with “3D CTF correction”, we do not adopt the name “3D CTF model”, instead using “3D-sampling function”.*

**V) 3D-sampling function compensated iterative refinement of sub-tomogram alignment**—The iterative refinement procedures commonly used in cryoEM are prone to erroneously fitting noise, known as “over-fitting”<sup>28</sup>. To minimize over-fitting emClarity divides the data from the beginning into two halves which are kept separate during refinement, the so-called “gold standard” approach<sup>29</sup>. Additionally, the references used in the constrained search are carefully filtered as follows. In each cycle the SSNR of the average is estimated by the “gold-standard” FSC. A figure-of-merit weighting<sup>30</sup> derived from this FSC is then combined with CTF amplitude restoration via our adaption of the post-reconstruction volume normalized single-particle Wiener Filter<sup>31</sup> (eq 8 in the original paper). Importantly, this adaptation involves an explicitly accounting of the directional anisotropy in the distribution of signal (Online Methods).

The iterative procedure is a local refinement, improving on the initial global alignment obtained during template matching. Importantly, we rotate the noisy particles back into the microscope reference frame for cross-correlation with the reference volume. This allows symmetry to be applied to the particle improving the SNR in the orientation search. This is not possible in SPA, where the particle is a projection and the reference must be rotated to the particle’s orientation, and to our knowledge, is not implemented in any other sub-tomogram averaging package.

**VI) Iterative refinement of tilt-series alignment**—emClarity implements iterative refinement of the tilt-series alignment by using sub-tomograms as fiducial markers, which we call tomogram constrained-particle refinement (tomoCPR.) It is an approach similar to “particle polishing”<sup>32</sup> implemented for SPA in RELION with two primary differences. First, the reference projections we generate for refining the location of the sub-tomogram fiducial markers in the “raw” tilt-series, includes information from neighboring particles as well as non-particle information that is present in the tomograms (Supplementary Figure 2). Second, tomoCPR constrains spatially proximal particles to behave similarly within a given projection, as in SPA, while also requiring them to vary smoothly as a group from projection to projection through the tilt-series. The set of image transformations (shift, in-plane rotation, tilt-angle, and magnification) are fit for a grid of overlapping patches each containing a fixed number of particles, determined by the total molecular weight. The single set of image transformations that minimize the error for all fiducials in a given patch over all projections are solved using IMOD’s *tiltalign*. As the patches overlap significantly (0.75), the image transformations vary smoothly over neighboring particles.

**VII) 3D-sampling function compensated classification**—Regions of significant variance across a data set may be visualized by overlaying a 3D “variance map” with the average structure. The “missing-wedge” produces significant artifacts that are specific to the orientation of each particle in the sample, but not necessarily its identity or conformation. Left uncorrected these artifacts obscure meaningful differences among particles, reflected in a diffuse variance across the data set (Supplementary Figure 3 a–c). A previously demonstrated technique for estimating the effect of the “missing-wedge” by using a binary

mask, called “wedge masked differences (WMDs)”, was shown to be a good first order correction<sup>33</sup>; however, the accuracy of this model breaks down when higher-resolution features are considered (Supplementary Figure 3 d–e). To allow higher-resolution information in the classification, we replace this binary wedge mask with our 3D-sampling function, resulting in a more accurate estimate of the artifacts introduced by the “missing-wedge” (Supplementary Figure 3 g–i). It is worth noting that this does not “fill in” any missing data, rather it estimates what a given particle should look like by distorting the current sub-tomogram average by that particle’s 3D-sampling function, and clusters based on the difference between this expected value and the observed particle.

**VIII) Multi-scale clustering**—We encode *a priori* biological information through introducing inter-voxel correlations at biologically relevant length scales, such as  $\sim 10$  Å for alpha-helical density, 18–20 Å for RNA helices or small protein domains, and  $\sim 40$  Å for larger protein domains. This is accomplished by selecting features of given length scale via a bandpass filter. A singular value decomposition is run at each length scale using the native MATLAB function SVD, and the singular vectors describing the greatest variance for each length scale are then concatenated into feature vectors for further clustering. While this approach is similar to existing ideas in multi-scale multi-variate statistical analysis applied in other fields<sup>34</sup>, because emClarity *considers each length scale simultaneously*, the approach is capable of providing a richer description of the feature space.

### emClarity improves resolution in sub-tomogram averaging

Given the inherent difficulty in working with extremely low SNR cryoEM data, and the sensitivity of the results to optimal selection of parameters, we elected to demonstrate our software using two publicly available data sets from the Electron Microscopy Pilot Image Archive<sup>35</sup> (EMPIAR). We show these published/deposited maps, juxtaposed with the maps obtained with emClarity in figure 2. Using emClarity, we obtained a structure at 7.8 Å for the yeast 80s ribosome (EMPIAR-10045) as compared to the structure at 12.9 Å solved using RELION version 1.4 (EMD-3228<sup>36</sup>) (Figure 2a). For the mammalian 80s ribosome (EMPIAR-10064), we obtained an improvement from 11.2 Å (solved using pyTOM (EMD-3420<sup>37</sup>)) to 8.6 Å (Figure 2b).

To evaluate the relative impact of each of the individual features implemented in emClarity, we incrementally included them into several reconstructions of the yeast 80s ribosome. To control for errors in alignment and perform a one-to-one comparison with EMD-3228, we used precisely the same particles and orientation parameters from the star files that accompany the raw data (EMPIAR-10045). We compare each map to an external reference map derived from SPA (EMD-2275<sup>38</sup>), via a cross-Fourier Shell Correlation (cross-FSC), starting from the RELION reconstruction as a control (Figure 2c). The results demonstrate the recovery of additional signal from the same data as each subsequent feature is incorporated.

The accuracy of our combined CTF correction approach, including a re-implementation of optimal-exposure filtering and 3D-sampling function based Wiener filtering is reflected in the magenta curve in figure 2c, which shows a substantial improvement over the cross-FSC

of the control (black). The largest single improvement comes from the tomoCPR, shown in green (which included optimal-exposure filtering and Wiener filtering as well.) The ribosome sample had very few gold-fiducial markers and limited overlapping density, perhaps accounting for the high impact of tomoCPR. A more modest improvement is obtained by adding in a per-tilt defocus estimation (cyan cross-FSC curve). When we also explicitly consider anisotropy in the SSNR in our adaptation of the single particle Wiener filter, we see another substantial improvement (red cross-FSC curve). The yeast 80s sample had a preferential orientation which is reflected in the FSC-cones and a plot of the angular distribution in supplementary figure 4. The final and highest resolution curve (blue) represents an alignment carried out in emClarity with all features added, illustrating the total impact on the accuracy of orientation determination.

In addition to improved resolution, we observe a density outside the peptide exit tunnel of the ribosome (figure 2b, white arrow) that is present in the map derived with emClarity, but not in the map derived with pyTOM. Finally, we show the density from a peripheral region with a rigidly docked model of the yeast 80s ribosome (PDB-47VR), clearly displaying alpha helices and the RNA structure in the map from emClarity (Figure 2c, blue box), compared to the map from the original publication (Figure 2c, black box).

To test emClarity on a more challenging case, we analyzed an HIV-1 immature Gag particle dataset<sup>39</sup> which was recently released (EMPIAR-10164). These data yielded the highest resolution (3.4 Å) sub-tomogram average to date (EMD-3782)<sup>23</sup>. Using emClarity, we produced a sub-tomogram average at 3.1 Å resolution (Figure 2d, Supplementary Figure 5). The quality of the density map is clearly manifested in the real-space refinement of the HIV-1 CA structure using Phenix<sup>40</sup> (Figure 2e).

### emClarity improves classification and reveals multiple ribosome functional states

Using multi-scale clustering combined with 3D-sampling function compensated Principal Component Analysis, emClarity helps reveal subtle conformational differences and distinguishes minor populations from noisy and distorted images, as demonstrated with yeast 80s ribosome data from EMPIAR-10045 (Figure 3), and mammalian 80s ribosome data from EMPIAR-10064 (Figure 4).

**Classification of non-translating Yeast 80s ribosomes**—The ribosome is a complex molecular machine composed of RNA and protein which exists in many functional states and interacts with an array of co-factors. The eukaryotic ribosome is composed of two major domains dubbed the large subunit (60s) and small subunit (40s). While the ribosome has a well-conserved catalytic core which mediates the peptidyl transferase reaction<sup>41</sup>, it is increasingly subject to more complex regulation in higher organisms resulting in an expanded set of both RNA and protein components. RNA expansion segments are found primarily at the periphery of the ribosome and are typically *highly dynamic* and difficult to resolve in structural analysis. One good example is es27, an approximately 150 Å RNA helix which predominantly adopts one of two conformations separated by about 90°, shown in orange in figure 3a–e. The first situates the end of the RNA helix just outside the peptide exit tunnel on the 60s subunit (es27<sub>pet</sub>, figure 3a, b, d, e) and the second points toward the

tRNA exit site (es27<sub>L1</sub>, figure 3c). This dynamic domain is generally observed in cryoEM maps as a superposition of these two states, as is the case with the currently published results by ML classification in RELION<sup>36</sup>. A notable exception being ribosomes with accessory complexes bound at the peptide exit tunnel, e.g. Sec61, are known to bias the conformation to the es27<sub>L1</sub><sup>42</sup>.

Another example of a highly dynamic ribosome domain is the L1 stalk – comprised of protein L1, and RNA helices h75, h76 and h79 from the 25s portion of the 60s subunit<sup>43</sup> – which moves through ~ 55Å during a translocation cycle. The motions of L1 are well correlated with several defined functional translational states as observed using single molecule FRET and SPA<sup>44</sup>. Using emClarity, we can discern three of these L1 conformational states isolated from thermal (stochastic) fluctuations of the non-translating yeast 80s ribosome: L1<sub>open</sub>, L1<sub>int</sub>, and L1<sub>closed</sub> shown in green with variable occupancy in the five classes in figure 3. In addition to isolating dynamic states, identifying very sparsely populated classes is a particularly important and challenging task for classification of cryoEM data. We see in figure 3e the dissociated 60s subunit occupying a minor class, only ~4% of the data set or roughly ~140 sub-tomograms. In contrast, the Maximum likelihood approach implemented in RELION found three classes, one designated as a junk class and two relatively indistinguishable classes<sup>36</sup>. This minor population could only be isolated in the case where feature vectors built from the projection on the principal components from at least three length-scales were simultaneously clustered.

**Mammalian 80s ribosome**—In contrast to the non-translating yeast specimen, the mammalian ribosomes imaged in EMPIAR-10064 were prepared from clarified rabbit reticulocyte lysate using a buffer low in Mg<sup>2+</sup> but lacking polyamines, such that cofactors should co-purify excepting perhaps some loss of e-Site tRNA<sup>45</sup>. We extracted 3,090 ribosomes from the four tilt-series deposited as the “mixed-CTEM” data set on EMPIAR, which are collected over a range of defocus values *without* a phase plate. emClarity identified five predominant classes as shown in figure 4a–e. Three of these classes show ribosomes adopting a non-rotated 40s conformation with variable tRNA eeF1A occupancy (class I-III), while two very similar classes adopted a mid-rotated (~5–6°) 40s conformation with eeF2 present (class IV-V).

A rigid body docking of the full 80s mammalian ribosome in the non-rotated POST state from PDB-4UJE<sup>46</sup> showed very clear agreement with the conformation of the 40s subunit, which combined with the co-factors observed suggest classes II and III are POST translational ribosomes differing in retention of E site tRNA while class I is most similar to the “sampling” state. Classes IV and V both have eeF2 bound and differ in rotation of the 40s subunit of 5.9° and 5.0°, respectively (Figure 4f–g). Rigidly docking the 80s yeast structure of eeF2 from PDB-4UJO<sup>47</sup> into classes IV and V show overall good agreement with their eeF2•sordarin•GDP position and our density. There are differences in domain IV of eeF2 which is known to be dynamic and plays a key role in translocation<sup>48,44</sup>. We analyzed these differences qualitatively by comparing the Molecular Dynamics Flexible Fitting (MDFF) model (figure 4h–i) with the docked model solved with Sordarin present (figure 4j–k). The antibiotic Sordarin is highly specific for binding to fungal eeF2 and permits GTP hydrolysis, yet prevents conformational changes that result in subsequent release of eeF2 after

translocation<sup>49</sup>. Sordarin is not present in the sample under study, yet there is a pronounced difference in electron density between domains III-V of eeF2 in class V (figure 4i black arrow) that coincides with the Sordarin binding pocket. This density is not present in class IV which also exhibits a rotation of eeF2-domain IV (figure 4j).

## Discussion

We have created a set of image processing routines incorporated into the program emClarity, which demonstrated improved accuracy in alignment and image restoration compared to current state-of-the-art approaches. We have aimed to make emClarity as easy to use as possible, limiting user specified parameters to the normal microscope and data collection information, as well as an estimate of the particle's radius and mass. The user must also select the angular search range, which is something that may be improved in the future.

We demonstrate a powerful approach for image classification in the presence of the “missing-wedge” effect by combining the correction for wedge differences with multi-scale clustering which helps to encode biologically relevant information for the clustering algorithms. Having isolated classes IV and V of the mammalian ribosome *ex vivo* suggests both that Sordarin binding stabilizes an interaction between eeF2-domain III/V that exists on pathway in functional ribosomes, and that nearby intermediates on the energy landscape not observed in studies using the antibiotic, may be explored and observed by using this approach. In addition to isolating well-resolved class averages with minor populations, and finding nearby minima in the energy landscape, our approach also results in the production of accurate 3D-variance maps.

By highlighting key regions of dynamic behavior, our approach should be useful for direct analysis and the design of complementary biophysical experiments. While these advances in classification are in the pre-processing and dimensionality reduction stage, future work to explore modern approaches in pattern recognition and machine learning, may allow substantial improvement in the technique.

## Online Methods

Please refer to “Life Sciences Reporting Summary” for detailed information on experimental design, reagents and software.

## Datasets

The datasets used in emClarity processing are from Electron Microscopy Public Image Archive (EMPIAR), including the yeast 80s ribosome (EMPIAR-10045), the mammalian 80s ribosome (EMPIAR-10064), and the HIV-1 immature Gag (EMPIAR-10164).

## emClarity programs

emClarity is run from the command line and is easily scripted to run in a manner most suited to a user's particular project. A text parameter file is used to input project specific details, like microscope parameters, mask dimensions, and angular search ranges. We typically make a copy of the parameter file for each cycle of averaging and alignment, which we refer



to as paramC.m, where C refers to the cycle number. The meta-data of each project is tracked in a binary database which named using the “subTomoMeta” parameter. Each tilt-series may have multiple areas slated for reconstruction, tiltN M refers to tilt-series “N” and reconstruction area “M”. A brief description of the major functions (*in italic*) in emClarity is below:

emClarity init paramN.m

Read in the desired dimensions for each sub-region of each tilt-series to reconstruct, initialize the subTomoMeta (metadata binary.)

emClarity ctf estimate tiltN

estimate the ctf for a given tilt-series N

emClarity reScale mapNameIN mapNameOUT AngPixIN AngPixOUT cpu/GPU

resample a map to a new pixel size, particularly for template matching.

emClarity templateSearch tiltN M reference.mrc symmetry gpuIDX

Reconstruct tomogram M from tilt-series N without ctf correction, run template matching on GPU # gpuIDX, randomize results at symmetry related positions.

emClarity ctf 3d paramN.m

Run 3d-CTF corrected weighted-back projection.

emClarity avg paramN.m N RawAlignment

Every cycle begins by creating a sub-tomogram average, calculating the “gold-standard” FSC, and weighting the average accordingly while compensating for amplitude attenuation by the CTF to produce references for the alignment.

emClarity alignRaw paramN.m N

Run the alignment

emClarity removeDuplicates paramN.m N

clean out any subtomograms that have drifted to the same position. Not needed every cycle.

emClarity tomoCPR paramN.m N

Run tomogram constrained particle refinement, this is generally done prior to a step down in binning, e.g. bin4 → bin3.

emClarity ctf update paramN.m N

Only needed after a run of tomoCPR, this updates the tilt-series geometry in the subTomoMeta, and also resamples the raw tilt-series applying rotation, shift and magnification scaling all in Fourier space to reduce interpolation losses of high resolution information.

Since the tilt-series alignments are update, and usually also the binning is reduced, a new round of 3D-CTF reconstructions need to be made with.

If classification is to be run, the cycle starts the same, but with the “flgClassification” parameter enabled.

```
emClarity avg paramN.m N
```

```
emClarity pca paramN.m N previousPCA
```

Run 3D-sampling function compensated PCA at each length scale specified in the “pcaScaleSpace” parameter. The command line argument previousPCA is always zero on the first run. If a random subset (25% or ~3000, whichever is larger) is to be analyzed by setting “Pca\_randSubset” then a subsequent round of pca must be run with “previousPCA” set to one in order to project all of the sub-tomograms along the principal component axes.

```
emClarity cluster paramN.m N
```

Cluster the data based on selected eigenvectors from the pca step.

```
emClarity avg paramN.m N Cluster_cls
```

Notice the last argument (a string) which creates a montage of the class averages selected in the parameter file. Classes with different memberships may be selected.

At the end of processing, the half-sets may be aligned and combined by running:

```
emClarity avg paramN.m N FinalAlignment
```

Align and combine half-sets, optionally creating multiple, differently sharpened maps.

### Image processing

The alignment and classification procedures are generally identical for all the samples, except for the HIV-1 Gag data, which were not classified and had C6 symmetry applied. All parameters unique to each dataset, including the angular search range and iterations used, are shown in Supplementary Tables 3 and 4. emClarity is only tested on Linux operating systems, and all references to command line operations are to be understood in that manner.

**Project set-up and coarse tilt-series alignment:** For each specimen, we make a project directory, which we will refer to generically as “projectDir.” During processing, several sub-directories will be created by emClarity in addition to the user created directories for the raw data, which we recommend calling “rawData” and a folder for the cleaned data that *must be* named “fixedStacks.”

The HIV-1 Gag data consist of dose fractionated frames, which we aligned using the program “unblur” version 1.0 included in the cisTEM package. The aligned frames were summed and saved *without* any exposure-based filtering because this is handled later inside emClarity.

All tilt-series were aligned using the default parameters in IMOD version 4.10.12 using the eTomo interface, with the available gold-fiducial markers. For the ribosome datasets, all gold fiducials (~5–7/tilt) were selected, while for the HIV-1 Gag data ~20–30 closest to the

protein and distributed on both surfaces of the ice were selected. Local alignments with fixed XYZ global coordinates were run for the HIV-1 Gag data only. After generating the final aligned stack, the gold beads were located using *find beads3d*. Only the fiducial model describing the location of the beads is needed, so they were not erased. Note: for EMPIAR 10045 the pixel size in the header must be corrected to 2.17 Å prior to beginning. This may be done with the IMOD program *alterheader* from the command line.

The files describing the projection transformations, any local alignments, and fitted tilt-angles are copied to the fixedStacks directory and renamed.

---

```
>$ mv specimen_name_1_fid.xf      projectDir/fixedStacks/tilt1.xf
>$ mv specimen_name_1_fid.tlt     projectDir/fixedStacks/tilt1.tlt
>$ mv specimen_name_1_local.xf    projectDir/fixedStacks/tilt1.local
>$ mv specimen_name_1_erase.fid   projectDir/fixedStacks/tilt1.erase
```

---

Additionally, a file listing the tilt-angles in the order they were collected must be created for emClarity to apply an appropriate exposure filter.

---

```
projectDir/fixedStacks/tilt1.order
```

---

If outlier pixels are removed in IMOD, this “fixed” stack may be moved to projectDir/fixedStacks/tilt1.fixed, otherwise you may just link the raw data.

```
>$ cd projectDir/fixedStacks
>$ ln -s ../rawData/specimen_name_1.st tilt1.fixed
```

This is repeated for all tilts-series, of which there are 7, 4, and 41 in the yeast, mammalian, and HIV-1 Gag data sets respectively.

**CTF estimation:** The mean defocus at the tilt-axis was then estimated in emClarity for each tilt-series using a  $3.5 \pm 2.5$  µm window covering the range of expected defocus values for all three data sets.

For the HIV-1 Gag data, the per-tilt defocus was determined using “emClarity ctf refine” to produce the power-spectra, which were subsequently fit using *ctffind4* with the–amplitude-spectrum input flag and default parameters.

For the yeast ribosome data which have a thin layer of carbon providing extra signal in the power spectrum, the per-tilt defocus values were refined during tomoCPR. To do so, the height of the cross-correlation peak is maximized by scanning through a small range of defocus values as applied to each reference tile<sup>50</sup>.

**Selecting sub-regions for further analysis**—The selection of sub-regions of each tilt-series for reconstruction are defined by a text file with the minimum and maximum values in x, y, z for each region. The script “recScript2.sh” provided with emClarity was used to first create reconstructions of each tilt-series at a binning of 10 and thickness of 300 covering the full X,Y dimension of the images. Each region is then defined while viewing the reconstruction in IMOD by making an IMOD model with six points per region, xmin, xmax, ymin, ymax, zmin, zmax, in that order.

A second run of “recScript2.sh” creates a projectDir/recon directory and converts these model files into the text files read in by emClarity to be used for the rest of the procedure. These are called tilt1\_recon.coords and list the tilt-series base name, number of regions to reconstruct, and for each region the width, first and last slice in y, thickness, x-origin offset and z-origin offset.

The ribosome data were divided on the x-axis into two regions per tilt-series. The HIV-1 Gag data were divided into quadrants. Additionally, the flag “fscGoldSplitOnTomos=1” is set in the parameter file for the HIV-1 Gag data, so that the even/odd half sets are divided based on tomogram, not randomly on sub-tomograms. This is necessary to avoid mixing neighboring particles which would violate the gold-standard hypothesis.

**Template matching**—References were derived from SPA EMD-3228<sup>39</sup> (yeast 80s ribosome), EMD-5592<sup>51</sup> (human 80s ribosome) and EMD-8403 (HIV-1 Gag)<sup>52</sup> and rescaled to the full pixel size of each data set using “emClarity reScale.” These references were then passed to “emClarity templateSearch” binned to achieve a nominal pixel size ~ 8–12 Å depending on the size of the specimen. All maps and tomograms are automatically low-pass filtered to 40 Å resolution by default in emClarity. Non-CTF corrected tomograms are reconstructed by the templateSearch program as needed for template matching.

The results for the ribosome data set were cleaned manually by comparing the maximum intensity projection maps and the binned tomograms overlaid with an IMOD model showing the x,y,z coordinates of each peak detected.

For the HIV-1 Gag data, emClarity removeNeighbors was used to automatically clean the results based on geometrical restraints. Only peaks that had five neighbors within 100 Å and also oriented within 20° were retained, resulting in 179,168 sub-tomograms to start. (This number dropped to 162,213 in the first round of averaging as particles too close to the edge to allow padding by 1.5 x particleRadius were excluded.

Particles with symmetry pose a special challenge to all missing-wedge compensation approaches as any error in the compensation will result in the particle looking different at its symmetry related orientations. To help with this, we randomize the template matching results around each symmetry related orientation, and subsequently only search an angular range small enough to not reach the neighboring positions.

**Iterative alignment**—Each cycle of alignment is initiated by calculating averages of the two half sets, calculating the gold-standard FSC, and then applying re-weighting each average to generate a FOM weighted reference. These alignments are only for the FSC

calculation and are never considered in the updates to the particle orientations in the iterative alignment, as necessary to maintain independence between the two half-sets.

We alternate searching over just the azimuthal and polar angles, and an in-plane search. For each specimen, we started at a binning of chosen to produce a pixel size of  $\sim 7\text{--}8$  Å. Details may be found in Supplementary Table 3 for the ribosome data sets, and Supplementary Table 4 for the HIV-1 gag data. We then go through three rounds of averaging and alignment, followed by removing any positions that may have drifted to overlap using “emClarity removeduplicates.” We then run a round of tomo-CPR, which requires updating the aligned tilt-series and the 3dCTF corrected tomograms.

```
>$ emClarity ctf update paramX.m
```

```
>$ emClarity ctf 3d paramX.m
```

This reconstruction is generated at a binning one finer than the previous, and the same pattern was repeated until reaching full sampling.

**Classification**—The ribosome data for the yeast 80s were classified in a single pass, using three resolution bands 10,18, and 28Å, 36 of the top eigenvalues were saved, and five from each band were selected (parameter Pca\_coefficients=[7:11;7:11;1:11]) for clustering via kmeans The class averages were then generated by running emClarity avg paramX.m X cluster\_cls

The ribosome data for the mammalian 80s were classified in two passes. First, they were split into groups displaying either a rotated or un-rotated 40s small subunit. To do this, the subTomoMeta file (projectName.mat) was copied to two new files: project\_smallSU.mat and project\_largeSU.mat. The classes are selected for removal by viewing the class average montage in IMOD, and selecting any point in the region of a given class. These models are then used to remove their contributing members in the subTomoMeta.

```
>$ emClarity geometry paramX.m X RemoveClasses [X,0,0] STD.
```

Since both branches of the project access the same raw data, it is convenient to remain in the same project directory, and all subsequent output will be identified by the new subTomoMeta basename.

A subsequent round of classification was run using 12,22,32 Å resolutions. Unlike the yeast 80s which had some Eigen images with clear missing wedge bias, revealed as “streakiness” in the density, the mammalian displayed sufficient true variability to overpower the noise from the missing-wedge bias, and all 36 eigenvectors from each resolution band were used in clustering.

**Analysis**—Models PDB-3J78 for yeast were rigid body docked in using Chimera.

Models PDB-4UJO for mammalian were docked in using Chimera, in combination with the “Segger” plugin.

Models PDB-5193 were docked in using Chimera, refined in real-space using Phenix version 1.13-2998-000, and manually edited in COOT version 0.8.9.

## Algorithmic details

### 3D-Sampling Function—

$$3D \text{ Sampling Function} \equiv SF^{3D} = \sum_{j=1}^S \sum_{i=1}^Z T^{i,j} |CTF_i^{2d}|^2 R^{2d} ExpFilter_i^{2d} \quad \text{eq 1.}$$

The first term in the summation is the combined transformation of projection (i) into the tomogram, and sub-tomogram (j) into the final average. The second term is the normal expression for the CTF limited to third order Seidel aberrations<sup>53</sup>, the third is the R-weighting, and the fourth is the optimal-exposure filter as defined<sup>26</sup>, Z is the number of projections in each tilt-series, and S the number of sub-tomograms.

**Refinement of tilt-series alignment—**Tomo-CPR works by combining the strengths of the 3D-model-based and feature-tracking approaches, while also taking advantage of the robust alignment tools developed for gold-fiducial alignment available in the IMOD package, which must be installed alongside emClarity. Starting with the tomogram (as in the 3D model approach), we additionally replace the density corresponding to our particles of interest at the proper orientation, with a copy of the high SNR sub-tomogram average and then re-project that synthetic tomogram using the IMOD program *tilt* (Supplementary Figure 2a). This re-projection also includes any local alignments previously determined and allows us to create a reference tilt-series along with a model for each sub-tomogram position in the 2D-projection. This model is used to cut tiles out of the data and reference projections for comparison by cross-correlation, while considering the CTF of the data projection at that point, as well as the structural noise from each particle's unique environment (Supplementary Figure 2b,c). These refined positions are then used as input to IMOD's *tiltalign* as if they were derived from gold fiducials, allowing us to take advantage of local refinement and robust fitting, as described previously<sup>54</sup>. The global changes to the projection geometry are applied to the tilt-series, while the local refinements are taken into consideration when the tomograms are reconstructed on the fly by emClarity. In addition to the importance of considering neighboring particles (Supplementary Figure 2d), additional high-contrast features, like the edge of the carbon foil or other particulate matter are pointed out in supplementary figure 2e, where a thin strip of one of the tomo-CPR references, prior to tiling, is shown.

**Statistical optimization of the SNR in the final map—**In addition to alternating phase reversals, the CTF also modulates the amplitudes of the data. Several programs for correcting the CTF phase and amplitude modulations directly in the 2D projections of tilted images are available; the two predominant being CTFPLOTTER and CTFPHASEFLIP<sup>55</sup> included in the IMOD package<sup>18</sup> for measurement and correction respectively, and TOMOCTF<sup>19</sup>. Both may be used to restore the amplitudes in the Fourier transforms of individual projections which inevitably amplifies noise in the process<sup>56</sup>. A more attractive

approach is to correct the phases on the projections, via multiplication by the CTF, and then to address the amplitudes after building the 3D reconstruction (sub-tomogram average).

The amplitude modulations are compensated using a Wiener filter in both SPA and the adaptation of RELION for sub-tomogram averaging<sup>10</sup>. A typical Wiener filter based approach has also been described recently in the structure of the HIV-1 capsid protein<sup>39</sup> but is only implemented through “in-house” software. In addition to the CTF and consideration of increasing sample thickness with tilt angle, our 3D-sampling function also takes into consideration the R-weighting that is applied during tomogram reconstruction and is applied via our adaption of the “volume normalized single-particle Wiener (original eq. 8)<sup>31</sup> We include the original expression, with our nomenclature in eq. s1.

$$F^{SPWi}(\mathbf{q}_{hkl}) = \frac{SF^{3D}}{SF^{3D} + \frac{f_{particle}}{f_{mask}} \left( \frac{1 - FSC_{mask}(\mathbf{q}_{hkl})}{2FSC_{mask}(\mathbf{q}_{hkl})} \right) \left( \frac{1}{n_q} \sum_{q \in \mathbf{q}_{h'k'l'}} SF^{3D} \right)} F^{LSQ}(\mathbf{q}_{hkl}) \text{eq s1}$$

The least squares estimate, which is a Wiener filtered reconstruction with an ad-hoc Wiener constant [CITE] is defined below eq s2.

$$F^{LSQ}(\mathbf{q}_{hkl}) = \frac{\sum_{j=1}^S \sum_{i=1}^Z T^{i,j} |CTF_i^{2d}|^2 R^{2d} ExpFilter_i^{2d} F_i^{2d}}{\sum_{j=1}^S \sum_{i=1}^Z T^{i,j} |CTF_i^{2d}|^2 R^{2d} ExpFilter_i^{2d} F_i^{2d} + 1} \text{eq s2}$$

We have made three major changes to the filter:

1. The 3D-Sampling function is weighted for critically under-sampled regions, where the SSNR estimated by the FSC is less reliable. This is done by:
  - a. choosing a minimum acceptable sampling threshold,  $0.2 * \text{median}(SF^{3D} \neq 0)$
  - b. scaling  $SF^{3D}$  to replace less sampled regions by smoothly transition from this value to some new larger number, chosen by the maximum in the original  $SF^{3D}$ .
2. The FSC, normally calculated over spherical sections is replaced by an anisotropic FSC calculated over conical sections, typically 38, which has been used previously to estimate resolution anisotropy<sup>57</sup>.
3. Finally, the average sampling over spherical shells (final term in the denominator of eq s1) that is used to scale the SSNR estimate to represent the average SSNR in a single sub-tomogram is replaced with a gaussian smoothed version of the 3D-sampling function. Again to account for anisotropy in the sampling.

$$F^{SPW}_i(\mathbf{q}_{hkl}) = \frac{|SF^{3D}|^2}{|SF^{3D}|^2 + \frac{f_{particle}}{f_{mask}} \left( \frac{1 - FSC_{aniso,mask}(\mathbf{q}_{hkl})}{2FSC_{aniso,mask}(\mathbf{q}_{hkl})} \right) (G \otimes |SF^{3D}|^2)} F^{LSQ}(\mathbf{q}_{hkl}) \quad \text{eq}$$

s3

**3D-sampling function compensated classification**—For the special case where subtomograms are all oriented similarly, being adsorbed to a lipid monolayer for example, they may be averaged along the direction of their missing wedge and classified in 2D<sup>58</sup>; however, this is obviously of limited interest for most specimen. Another popular approach involves classifying the

constrained-cross-correlation matrix<sup>9,27</sup>, however, this can be expensive to calculate and also disregards large amounts of information – discussed in detail in the paper on wedge masked differences (WMDs)<sup>33</sup> on which we further expand and enhance.

The WMDs approach seeks to compensate the missing wedge by forming the difference between a given particle and its expected value. The expected value is estimated to be the global average distorted by the particle's missing wedge. These are mean centered, normalized to a variance of one, and arranged into a 2D matrix followed by singular-value decomposition (SVD). The binary wedge used in this approach is only a first approximation, and we replace it with our full 3D-sampling function. This correction allows the classification to include higher-resolution details than previously possible, which is a necessary but not sufficient condition to achieve the classification we report. We find that including the highest variance information from three to four discrete length scales at the same time is required to observe each class.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments:

We thank J. Frank and W. Li for very helpful discussions, D. Bevan for technical assistance with computer clusters, S. Loerch for the help with Phenix real-space refinement and COOT, and T. Brosewitz, F. J. Alvarez, and J. P. Rickgauer for reading the manuscript. We thank F. Schur and J. Briggs for the HIV-1 immature Gag dataset, and X. Fu for testing the emClarity software. This work was supported by the National Institutes of Health (GM085043, GM082251) and the UK Wellcome Trust Investigator Award (206422/Z/17/Z).

## References

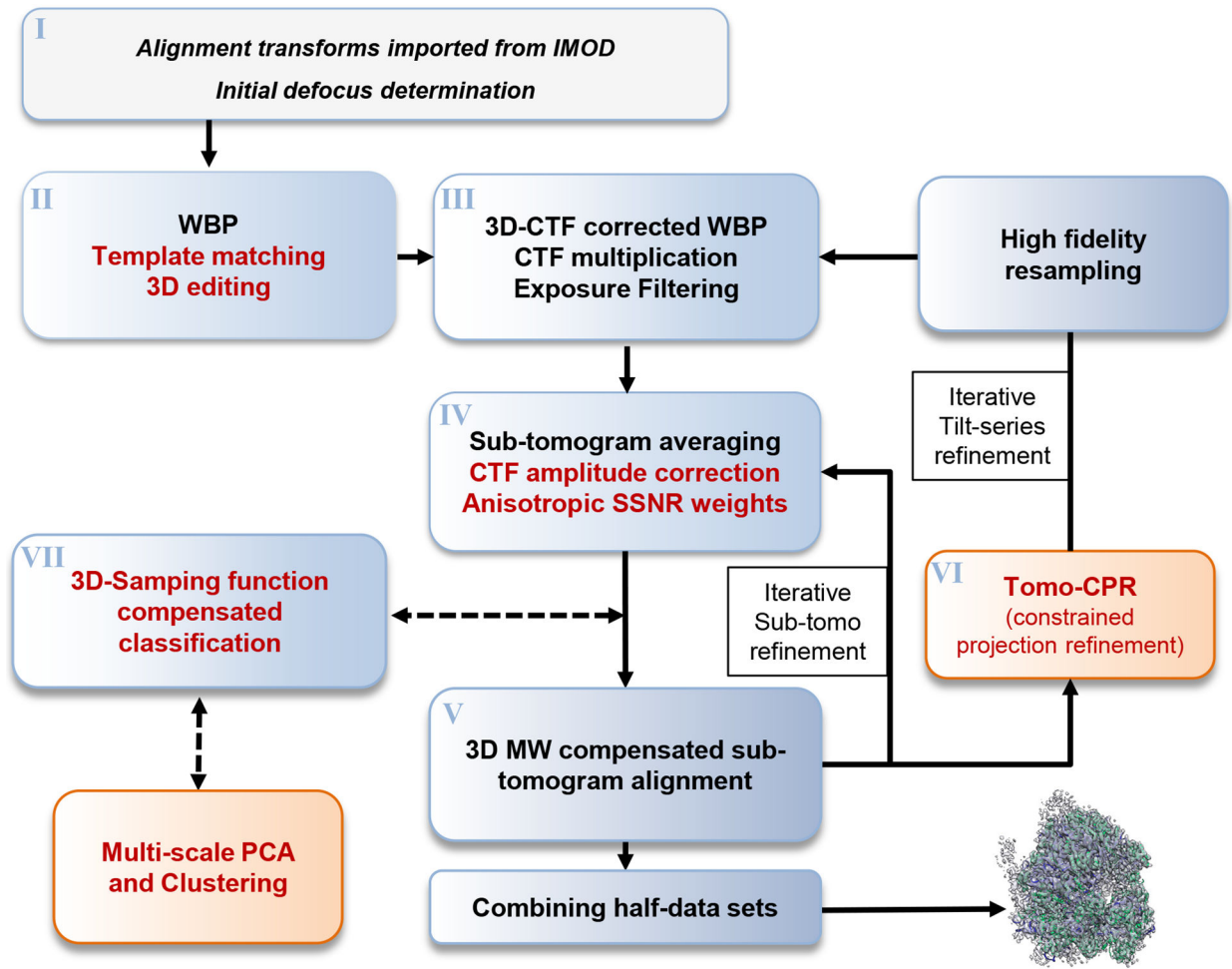
1. Grant T, Rohou A & Grigorieff N cis TEM, user-friendly software for single- particle image processing. 1–24 (2018).
2. Scheres SHW RELION: Implementation of a Bayesian approach to cryo-EM structure determination. *J. Struct. Biol* 180, 519–530 (2012). [PubMed: 23000701]



3. Glaeser RM & Hall RJ Reaching the information limit in cryo-EM of biological macromolecules: experimental aspects. *Biophys. J* 100, 2331–7 (2011). [PubMed: 21575566]
4. Cheng Y, Grigorieff N, Penczek PA & Walz T A primer to single-particle cryo-electron microscopy. *Cell* 161, 439–449 (2015).
5. Oikonomou CM & Jensen GJ Cellular Electron Cryotomography: Toward Structural Biology In Situ. *Annu Rev Biochem* 1–24 (2017). doi:10.1146/annurev-biochem-061516-044741 [PubMed: 28125288]
6. Diebold CA, Koster AJ & Koning RI Pushing the resolution limits in cryo electron tomography of biological structures. *J. Microsc* 248, 1–5 (2012). [PubMed: 22670690]
7. Liu V, Rigort A & Baumeister W Cryo-electron tomography: The challenge of doing structural biology in situ. *J. Cell Biol* 202, 407–419 (2013). [PubMed: 23918936]
8. Frangakis AS et al. Identification of macromolecular complexes in cryoelectron tomograms of phantom cells. *Proc. Natl. Acad. Sci. U. S. A* 99, 14153–14158 (2002). [PubMed: 12391313]
9. Bartesaghi A et al. Classification and 3D averaging with missing wedge correction in biological electron tomography. *J. Struct. Biol* 162, 436–450 (2008). [PubMed: 18440828]
10. Bharat TAM, Russo CJ, Löwe J, Passmore LA & Scheres SHW Advances in Single-Particle Electron Cryomicroscopy Structure Determination applied to Sub-tomogram Averaging. *Structure* 23, 1743–1753 (2015). [PubMed: 26256537]
11. Cassidy CK et al. CryoEM and computer simulations reveal a novel kinase conformational switch in bacterial chemotaxis signaling. *Elife* 4, (2015).
12. Zeev-Ben-Mordehai T et al. Two distinct trimeric conformations of natively membrane-anchored full-length herpes simplex virus 1 glycoprotein B. *Proc. Natl. Acad. Sci. U. S. A* 113, 4176–4181 (2016). [PubMed: 27035968]
13. Penczek PA, Frank J & Spahn CMT A method of focused classification, based on the bootstrap 3D variance analysis, and its application to EF-G-dependent translocation. *J. Struct. Biol* 154, 184–194 (2006). [PubMed: 16520062]
14. Liao HY, Hashem Y & Frank J Efficient Estimation of Three-Dimensional Covariance and its Application in the Analysis of Heterogeneous Samples in Cryo-Electron Microscopy. *Structure* 23, 1129–1137 (2015). [PubMed: 25982529]
15. Trabuco LG, Villa E, Schreiner E, Harrison CB & Schulten K Molecular dynamics flexible fitting: A practical guide to combine cryo-electron microscopy and X-ray crystallography. *Methods* 49, 174–180 (2009). [PubMed: 19398010]
16. Schur FKM, Hagen WJH, de Marco A & Briggs J. a G. Determination of protein structure at 8.5 Å resolution using cryo-electron tomography and sub-tomogram averaging. *J. Struct. Biol* 184, 394–400 (2013). [PubMed: 24184468]
17. Kudryashev M, Castaño-Díez D & Stahlberg H Limiting Factors in Single Particle Cryo Electron Tomography. *Comput. Struct. Biotechnol. J* 1, 1–6 (2012).
18. Kremer JR, Mastronarde DN & McIntosh JR Computer visualization of three-dimensional image data using IMOD. *J. Struct. Biol* 116, 71–6 (1996). [PubMed: 8742726]
19. Fernández JJ, Li S & Crowther R. a. CTF determination and correction in electron cryotomography. *Ultramicroscopy* 106, 587–96 (2006). [PubMed: 16616422]
20. Rohou A & Grigorieff N CTFFIND4: Fast and accurate defocus estimation from electron micrographs. *J. Struct. Biol* 192, 216–221 (2015). [PubMed: 26278980]
21. Hrabe T et al. PyTom: A python-based toolbox for localization of macromolecules in cryo-electron tomograms and subtomogram analysis. *J. Struct. Biol* 178, 177–188 (2012). [PubMed: 22193517]
22. Jensen GJ & Kornberg RD Defocus-gradient corrected back-projection. *Ultramicroscopy* 84, 57–64 (2000). [PubMed: 10896140]
23. Turovová B, Schur FKM, Wan W & Briggs JAG Efficient 3D-CTF correction for cryo-electron tomography using NovaCTF improves subtomogram averaging resolution to 3.4 Å. *J. Struct. Biol* (2017). doi:10.1016/j.jsb.2017.07.007
24. Kunz M & Frangakis AS Three-dimensional CTF correction improves the resolution of electron tomograms. *J. Struct. Biol* 197, 114–122 (2017). [PubMed: 27343995]

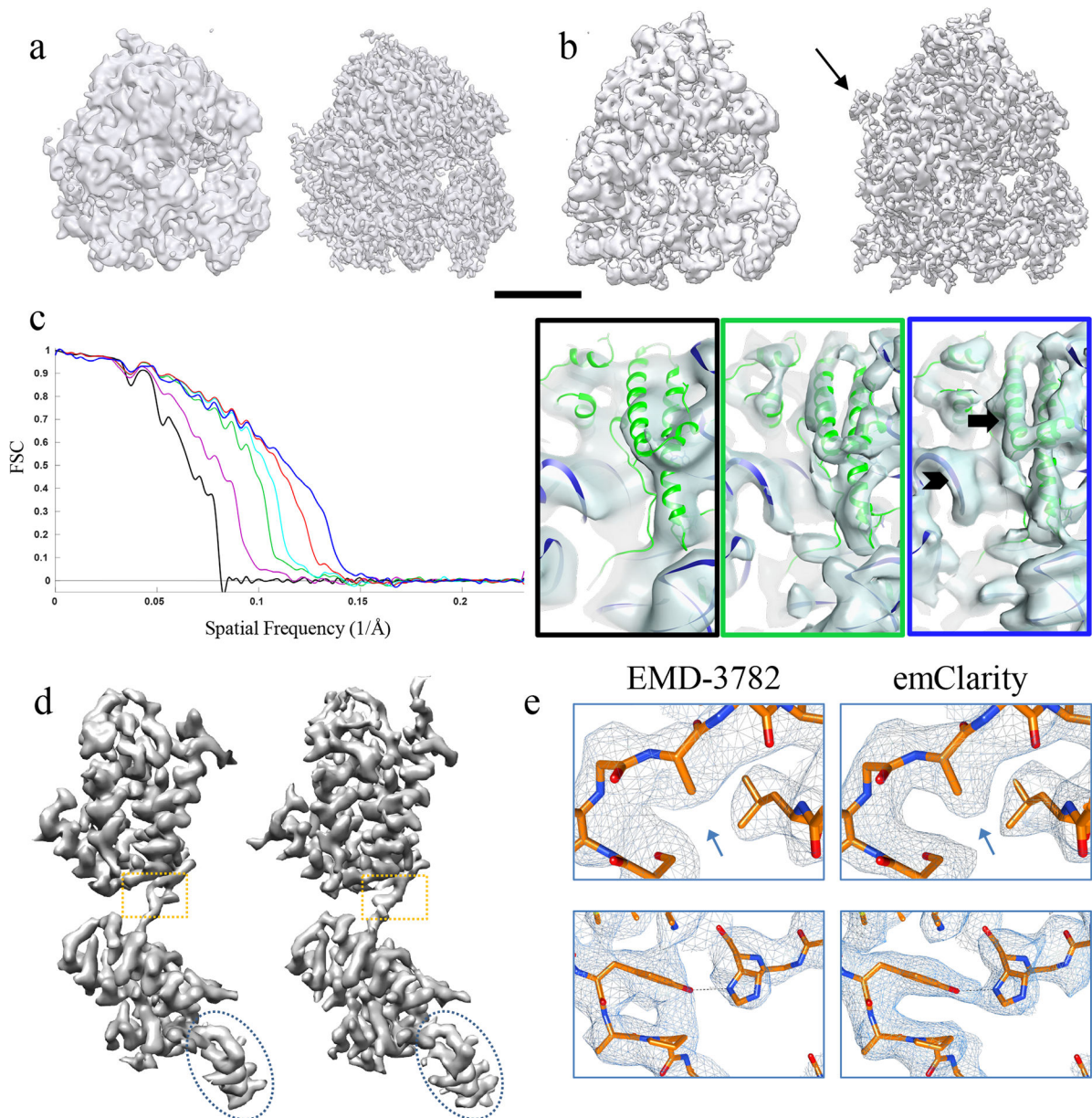
25. Rickgauer JP, Grigorieff N & Denk W Single-protein detection in crowded molecular environments in cryo-EM images. *Elife* 6, 1–22 (2017).
26. Grant T & Grigorieff N Measuring the optimal exposure for single particle cryo-EM using a 2.6 Å reconstruction of rotavirus VP6. *Elife* 4, e06980 (2015). [PubMed: 26023829]
27. Förster F, Pruggnaller S, Seybert A & Frangakis AS Classification of cryo-electron sub-tomograms using constrained correlation. *J. Struct. Biol* 161, 276–286 (2008). [PubMed: 17720536]
28. Stewart A & Grigorieff N Noise bias in the refinement of structures derived from single particles. *Ultramicroscopy* 102, 67–84 (2004). [PubMed: 15556702]
29. Henderson R et al. Outcome of the first electron microscopy validation task force meeting. *Structure* 20, 205–214 (2012). [PubMed: 22325770]
30. Rosenthal PB & Henderson R Optimal Determination of Particle Orientation, Absolute Hand, and Contrast Loss in Single-particle Electron Cryomicroscopy. *J. Mol. Biol* 333, 721–745 (2003). [PubMed: 14568533]
31. Sindelar CV & Grigorieff N Optimal noise reduction in 3D reconstructions of single particles using a volume-normalized filter. *J. Struct. Biol* 180, 26–38 (2012). [PubMed: 22613568]
32. Scheres S. H. w. Beam-induced motion correction for sub-megadalton cryo-EM particles. *Elife* 3, e03665 (2014). [PubMed: 25122622]
33. Heumann JM, Hoenger A & Mastronarde DN Clustering and variance maps for cryo-electron tomography using wedge-masked differences. *J. Struct. Biol* 175, 288–299 (2011). [PubMed: 21616153]
34. Alsberg BK Multiscale cluster analysis. *Anal. Chem* 71, 3092–3100 (1999). [PubMed: 21662901]
35. Marabini R et al. The Electron Microscopy eXchange (EMX) initiative. *J. Struct. Biol* 194, 156–163 (2016). [PubMed: 26873784]
36. Bharat TAM & Scheres SHW Resolving macromolecular structures from electron cryo-tomography data using subtomogram averaging in RELION. *Nat. Protoc* 11, 2054–2065 (2016). [PubMed: 27685097]
37. Khoshouei M, Pfeffer S, Baumeister W, Förster F & Danev R Subtomogram analysis using the Volta phase plate. *J. Struct. Biol* 197, 94–101 (2017). [PubMed: 27235783]
38. Bai XC, Fernandez IS, McMullan G & Scheres SHW Ribosome structures to near-atomic resolution from thirty thousand cryo-EM particles. *Elife* 2013, 2–13 (2013).
39. Schur FKM et al. An atomic model of HIV-1 capsid-SP1 reveals structures regulating assembly and maturation. *Science* (80-. ) 353, 506–508 (2016).
40. Adams PD et al. PHENIX: A comprehensive Python-based system for macromolecular structure solution. *Acta Crystallogr. Sect. D Biol. Crystallogr* 66, 213–221 (2010). [PubMed: 20124702]
41. Gutell RR, Weiser B, Woese CR & Noller HF Comparative anatomy of 16-S-like ribosomal RNA. *Prog. Nucleic Acid Res. Mol. Biol* 32, 155–216 (1985). [PubMed: 3911275]
42. Beckmann R et al. Architecture of the Protein-Conducting Channel Associated with the Translating 80S Ribosome. *Cell* 107, 361–372 (2001). [PubMed: 11701126]
43. Mohan S & Noller HF Recurring RNA structural motifs underlie the mechanics of L1 stalk movement. *Nat. Commun* 8, 14285 (2017). [PubMed: 28176782]
44. Spahn CM et al. Domain movements of elongation factor eEF2 and the eukaryotic 80S ribosome facilitate tRNA translocation. *EMBO J* 23, 1008–1019 (2004). [PubMed: 14976550]
45. Wilson DN & Nierhaus KH The E-site story: the importance of maintaining two tRNAs on the ribosome during protein synthesis. *Cell. Mol. Life Sci* 63, 2725–2737 (2006). [PubMed: 17013564]
46. Budkevich TV et al. Regulation of the Mammalian Elongation Cycle by Subunit Rolling: A Eukaryotic-Specific Ribosome Rearrangement. *Cell* 158, 121–131 (2014). [PubMed: 24995983]
47. Abeyrathne PD, Koh CS, Grant T, Grigorieff N & Korostelev AA Ensemble cryo-EM uncovers inchworm-like translocation of a viral IRES through the ribosome. *Elife* 5, 1–31 (2016).
48. Gomez-Lorenzo MG et al. Three-dimensional cryo-electron microscopy localization of EF2 in the *Saccharomyces cerevisiae* 80S ribosome at 17.5 Å resolution. *EMBO J* 19, 2710–8 (2000). [PubMed: 10835368]

49. Chakraborty B, Mukherjee R & Sengupta J Structural insights into the mechanism of translational inhibition by the fungicide sordarin. *J. Comput. Aided. Mol. Des* 27, 173–184 (2013). [PubMed: 23397219]
50. Meyer RR, Kirkland AI & Saxton WO A new method for the determination of the wave aberration function for high-resolution TEM. 2. Measurement of the antisymmetric aberrations. *Ultramicroscopy* 99, 115–123 (2004). [PubMed: 15093938]
51. Anger AM et al. Structures of the human and *Drosophila* 80S ribosome. *Nature* 497, 80–85 (2013). [PubMed: 23636399]
52. Ning J et al. In vitro protease cleavage and computer simulations reveal the HIV-1 capsid maturation pathway. *Nat. Commun* 7, 13689 (2016). [PubMed: 27958264]
53. Fernando KV & Fuller SD Determination of astigmatism in TEM images. *J. Struct. Biol* 157, 189–200 (2007). [PubMed: 17067820]
54. Mastronarde DN Fiducial Marker and Hybrid Alignment Methods for Single- and Double-axis Tomography in Electron Tomography: Methods for Three-Dimensional Visualization of Structures in the Cell (ed. Frank J) 163–185 (Springer New York, 2006). doi:10.1007/978-0-387-69008-7\_6
55. Xiong Q, Morphew MK, Schwartz CL, Hoenger AH & David N CTF Determination and Correction for Low Dose Tomographic Tilt Series. *J. Struct. Biol* 168, 378–387 (2010).
56. Frank J Electron Microscopy of Macromolecular Assemblies in Three-Dimensional Electron Microscopy of Macromolecular Assemblies 15–69 (Oxford University Press, 2006).
57. Diebolder CA, Faas FGA, Koster AJ & Koning RI Conical fourier shell correlation applied to electron tomograms. *J. Struct. Biol* 190, 215–223 (2015). [PubMed: 25843950]
58. Winkler H et al. Tomographic subvolume alignment and subvolume classification applied to myosin V and SIV envelope spikes. *J. Struct. Biol* 165, 64–77 (2009). [PubMed: 19032983]



**Figure 1. The emClarity workflow for sub-tomogram averaging and classification.**

The most important improvements introduced in emClarity are highlighted in red text, while novel additions to the process [AU: please clarify if you mean the general workflow for cryo-ET image processing here?] are indicated in orange boxes. Dashed lines indicate optional branches which may be included any number of times during the iterative alignment.



**Figure 2. Improvement in resolution of sub-tomogram averaging using emClarity.**

(a-b) Comparison of the sub-tomogram average of yeast 80s ribosome (a) and of rabbit 80s ribosome (b) by RELION (EMD-3228) (a, left) or by pyTOM (EMD-3420) (b, left) and by emClarity (right). Arrow points to an additional feature only revealed in emClarity. Scale bar, 100 Å. (c) Cross-FSC between the sub-tomogram averages by emClarity and the SPA cryoEM map (EMD-2275) of yeast 80s ribosome, each accumulating the previous improvement: original orientation parameters from Relion (black), CTF estimation and correction with the optimal exposure filter (magenta), one round of tomo-CPR (green), per-tilt defocus estimation (cyan), consideration of resolution anisotropy (red), and with all features plus alignment in emClarity (blue). Right, representative views of sub-tomogram averages, with the frame color matching the plot colors, and a rigid body docking of yeast

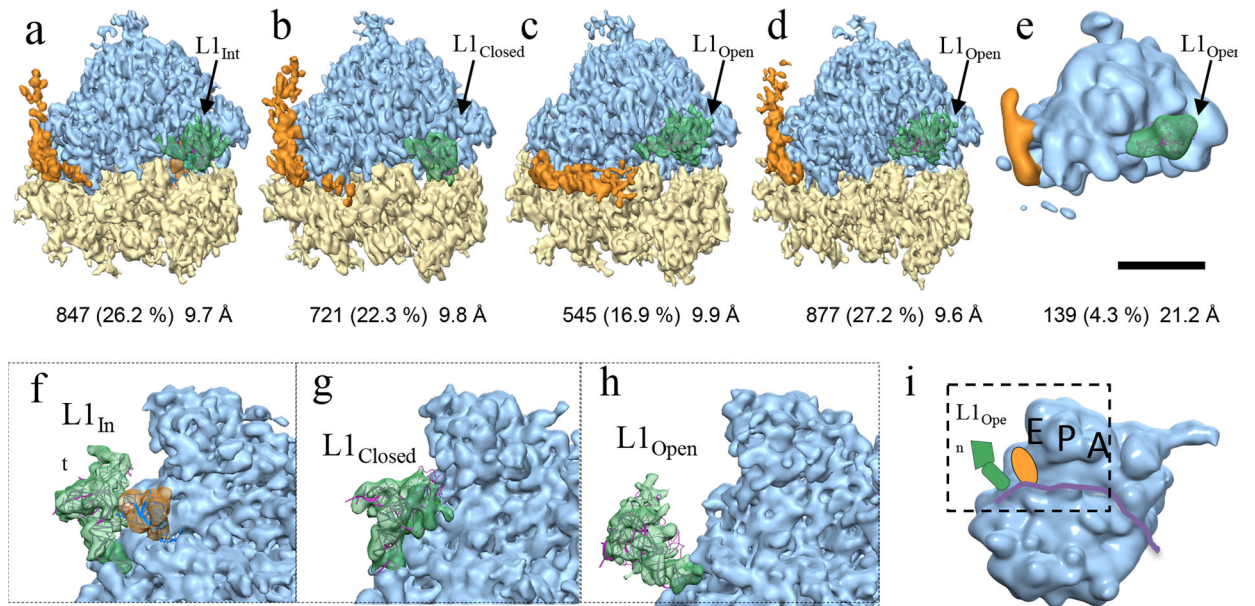
80s atomic model (PDB-47VR). The arrow and chevron highlight the resolved alpha helices and RNA structures, respectively. These experiments were repeated at least three times with nominally identical results. (d) Comparison of the sub-tomogram average of HIV-1 immature CA-SP1 monomer (EMD-3782) (left) and by emClarity (right). (e) Enlarged views of boxed area in d overlaid with a real-space refined model.

Author Manuscript

Author Manuscript

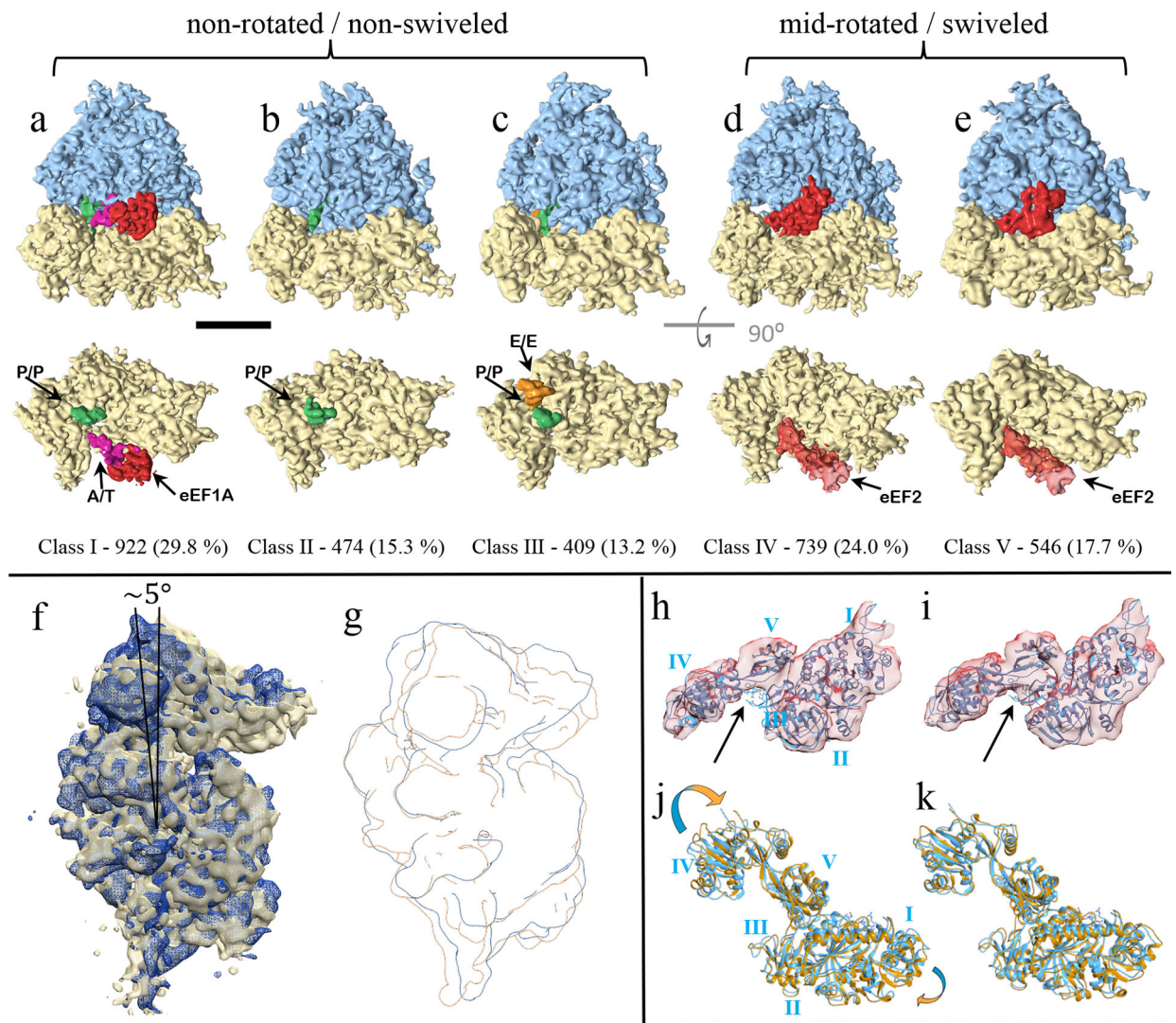
Author Manuscript

Author Manuscript



**Figure 3. Classification of yeast 80s ribosome (EMPIAR-10045) with 3D-CTF compensated missing-wedge and multi-scale PCA.**

(a-e) Four major classes and a minor class contributing 96.9% of sub-tomograms are shown with number and percentage of contributing units and resolution indicated below. The highly dynamic L1 protuberance and RNA expansion-segment 27 are shown in green and orange, respectively. Scale bar, 100 Å. (f-h) Enlarged views, with the small-subunit removed for clarity, of the L1 protuberance (green) in an intermediate position (f, class a), fully closed (g, class b), and fully open (h, classes c-e) respectively. The riboprotein and selected RNA helix components of the L1 protuberance (rpL1,h76,h79 from PDB 3J78) shown in magenta ribbon after rigid body docking into the respective density. (i) A cartoon illustrates the region displayed in (f-h) from the inter-subunit space with the E,P,A sites labelled and L1 (green) E-site tRNA (orange) and mRNA channel (purple ribbon.)



**Figure 4. Classification of translating mammalian 80s ribosome with 3D-CTF compensated missing-wedge and multi-scale PCA.**

(a-c) Classes I-III represent a post translocational state with the co-factors shown in the lower row from the inter-subunit surface with the 60s subunit removed for clarity. (d-e) Classes IV-V have a mid-rotated 40s and a sweveled head corresponding to a pre-translocational intermediate. Scale bar represents 100 Å. (f-g) Class IV in dark blue overlaid with class III in gold, showing the mid-rotated 40s state. (h-i) MDFF of eEF2 (orange) with the density from class-IV (h) and class-V (i) starting from PDB-4ujo (cyan ribbon). (j-k) The rigid body docking of PDB-4ujo into the density from class IV/V, respectively. Arrows point to this density which is occupied by the antibiotic Sordarin in PDB-4ujo, but is not present in the sample used in this study.