

# In situ structure determination using single particle cryo-electron microscopy images

**Jing Cheng**

Institute of Biophysics, Chinese Academy of Sciences

**Bufan Li**

Institute of Biophysics, Chinese Academy of Sciences

**Long Si**

Institute of Biophysics, Chinese Academy of Sciences

**Xinzheng Zhang** (✉ [xzzhang@ibp.ac.cn](mailto:xzzhang@ibp.ac.cn))

Institute of Biophysics, Chinese Academy of Sciences

---

## Article

**Keywords:** cryo-EM, isSPA, tomography

**Posted Date:** November 2nd, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-94639/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# *In situ* structure determination using single particle cryo-electron microscopy images

Jing Cheng<sup>1,2</sup>, Bufan Li<sup>1,2</sup>, Long Si<sup>1,2</sup>, Xinzheng Zhang<sup>1,2,3,\*</sup>

<sup>1</sup> National Laboratory of Biomacromolecules, CAS Center for Excellence in Biomacromolecules, Institute of Biophysics, Chinese Academy of Sciences, 100101 Beijing, People's Republic of China.

<sup>2</sup> University of Chinese Academy of Sciences, 100049 Beijing, People's Republic of China.

<sup>3</sup> Center for Biological Imaging, CAS Center for Excellence in Biomacromolecules, Institute of Biophysics, Chinese Academy of Sciences, 100101 Beijing, China.

\* For correspondence: xzzhang@ibp.ac.cn

## Abstract

Cryo-electron microscopy (cryo-EM) tomography is a powerful tool for *in situ* structure determination. However, this method requires the acquisition of tilt series, and its time consuming throughput of acquiring tilt series severely slows determination of *in situ* structures. By treating the electron densities of non-target protein as non-Gaussian distributed noise, we developed a new target function that greatly improves the efficiency of the recognition of the target protein in a single cryo-EM image without acquiring tilt series. Moreover, we developed a sorting function that effectively eliminates the false positive detection, which not only improves the resolution during the subsequent structure refinement procedure but also allows using homolog proteins as models to recognize the target protein. Together, we developed an *in situ* single particle analysis (isSPA) method. Our isSPA method was successfully applied to solve structures of glycoproteins on the surface of a non-icosahedral virus and Rubisco inside the carboxysome. The cryo-EM data from both samples were collected within 24 hours, thus allowing fast and simple structural determination *in situ*.

## 1. Introduction

43 Determining the *in situ* structure of working protein machineries in their native  
44 context allows for more physiological structural information, together with identifying  
45 the interactions with other proteins nearby. One of the best technologies to determine  
46 *in situ* structures is cryo-electron tomography (Beck & Baumeister, 2016; Lučić, Leis,  
47 & Baumeister, 2008) When combined with sub-tomogram averaging technique (M.  
48 Chen et al., 2019; Leigh et al., 2019) that increases the signal noise ratio (SNR) of target  
49 protein complexes with multiple copies in the tomogram by aligning and averaging the  
50 three dimensional (3D) volume of the protein complex, *in situ* protein structures have  
51 previously been determined at sub-nanometer resolution on non-cryo-sectioning  
52 sample (Dodonova et al., 2017; Himes & Zhang, 2018; Mattei et al., 2018; Pfeffer et  
53 al., 2015; Schur, 2016; F. K. M. Schur et al., 2015; Turoňová, Schur, Wan, & Briggs,  
54 2017; Wan et al., 2017) or nanometer resolution on cryo-sectioning sample (Bäuerlein  
55 et al., 2017; Bykov et al., 2017; Freeman Rosenzweig et al., 2017; Q. Guo et al., 2018;  
56 Mosalaganti et al., 2018). However, tomography requires the acquisition of a tilt series  
57 of the target protein complex. A tilt series typically contains more than 30 images taken  
58 at a range of tilt angles typically acquired within a time of 30 min, resulting in a  
59 slowdown of data collection throughput. The recent development of a stable sample  
60 stage in electron microscopy allows faster data collection by decreasing the waiting  
61 time of the stage (Chreifi, Chen, Metskas, Kaplan, & Jensen, 2019). However, the  
62 throughput is still dozens of times slower than the collection of single particle data.

63 In a single particle image, when the target protein complex is located in an *in situ*  
64 environment, the density of the target protein complex in the image is overlapped by  
65 other densities from surrounding proteins or biological molecules. The overlapping  
66 densities can be considered as noise that decreases the SNR of the image, especially at  
67 the range of low frequencies, where the shot noise is much lower than signals. The low-  
68 frequency signals exhibiting high SNR are essential for determining the initial position  
69 and orientation of the protein complex which is required by applying a conventional  
70 iterative single particle algorithm. A previous study showed that by using a high-  
71 resolution model of the target protein complex, the initial position and orientation of  
72 this protein complex can be determined from the protein background by incorporating  
73 the high-frequency signals of the target protein into the search (Rickgauer, Grigorieff,  
74 & Denk, 2017). However, the usage of the high resolution structure of the target protein  
75 renders the method less practical, since this structure is usually not available yet. A  
76 simple whitening filter was applied to both the reference and the image before  
77 calculating the correlation coefficient (Rickgauer et al., 2017), which did not take the  
78 overlapping density and SNR oscillation into account. This oscillation stems from the  
79 signal oscillation induced by the contrast transfer function (CTF). Furthermore, the shot  
80 noise in the image follows a Gaussian distribution, resulting in a smooth background in  
81 Fourier space. Therefore, CTF-like weighting has been widely used in score function  
82 (F. Guo & Jiang, 2014; Jasenko Zivanov et al., 2018). However, the overlapping  
83 densities can be regarded as part of the noise that fails to follow Gaussian distribution.  
84 How to optimize the score function when considering the noise distributed in a non-  
85 Gaussian manner remains to be investigated.

86 Our and other's previous studies showed that in cases where the density of target

87 protein complexes is overlapped by other densities, after extracting information of the  
 88 initial center and orientation of the protein complex from single particle results (Zhu et  
 89 al., 2018) or from sub tomogram averaging (Song et al., 2019), the structure of a protein  
 90 complex can be effectively refined using traditional local refinement procedures  
 91 without having to subtract overlapping densities. Both methods provide accurate initial  
 92 center and orientation parameters. However, the determination of the initial center and  
 93 orientation of a protein complex by incorporating a high-frequency signal with low  
 94 SNR introduces false positive solutions which limits the resolution of the reconstruction.  
 95 An effective method to reduce this problem and improve the resolution is urgently  
 96 required.

97 Here, we developed a single particle-like work flow to determine the *in situ*  
 98 structure of protein complexes by combining an optimized picking function to provide  
 99 the initial orientation and location of a target protein and a sorting algorithm to  
 100 effectively distinguish the correct solution from the false positive solutions to improve  
 101 the resolution.

## 102 103 **2. Results**

### 104 **2.1. Theoretical background**

105 To localize the target protein in a crowded environment, a cross-correlation  
 106 coefficient (cc) between a template of the target protein and raw cryo-EM image, had  
 107 previously been applied as a picking function (Rickgauer et al., 2017). However, cryo-  
 108 EM data are characterized by a frequency-dependent SNR. The template from a cryo-  
 109 EM reconstruction also contains frequency-dependent noise. The SNR of the template  
 110 can be calculated according to the Fourier Shell Correlation (FSC). The noise in both  
 111 the cryo-EM image and the template contribute errors to the cc. Therefore, an  
 112 appropriate frequency-dependent weighting function should be applied to the picking  
 113 function to minimize the error in cc. In this work, the densities of the target protein in  
 114 a cryo-EM image are overlapped with the densities from other surrounding proteins *in*  
 115 *situ*. We treat the densities from surrounding proteins as noise, together with the shot  
 116 noises. However, different from the shot noise, the densities from the surrounding  
 117 proteins are determined by the structure factor of proteins and modified by the contrast  
 118 transfer function.

119 Our picking function is based on cc with a frequency-dependent weight applied,  
 120 which can be expressed in Fourier space as

$$121 \quad cc = \sum_k W(k) \cdot X(k) \cdot M^*(k), \quad (1)$$

122 Here,  $k$  is the spatial frequency,  $W(k)$  is a weighting function,  $M^*(k)$  is the conjugate  
 123 complex value of the Fourier transform of a projection of a 3D template and  $X(k)$  is a  
 124 raw image in Fourier form. Signal  $S(k)$  is modulated by CTF in raw image, and  
 125 considering shot noise  $N(k)$  and protein noise which representing overlapping protein  
 126 densities  $PN(k)$ ,  $X(k)$  and  $M^*(k)$  are written as below

$$127 \quad X(k) = CTF(k) \cdot S(k) + N(k) + CTF(k) \cdot PN(k), \quad (2)$$

$$128 \quad M^*(k) = S^*(k) + N'^*(k), \quad (3)$$

129 thus we can rewrite  $cc$

$$130 \quad cc = \sum_k W(k) \cdot (CTF(k) \cdot S(k) + N(k) + CTF(k) \cdot PN(k)) \cdot (S^*(k) + N'^*(k)). \quad (4)$$

131 Weighting function is calculated to maximize the SNR of  $cc$

$$132 \quad SNR_{cc} = \frac{\left\{ \sum_k W(k) \cdot CTF(k) \cdot S(k) \cdot S^*(k) \right\}^2}{\left\{ \sum_k W(k) \cdot N(k) \cdot S^*(k) + W(k) \cdot N'^*(k) \cdot CTF(k) \cdot S(k) + W(k) \cdot N'^*(k) \cdot N(k) + \sum_k W(k) \cdot PN(k) \cdot CTF(k) \cdot S^*(k) + W(k) \cdot PN(k) \cdot CTF(k) \cdot N'^*(k) \right\}^2}, \quad (5)$$

$$133 \quad \frac{d(SNR_{cc})}{dW} = 0. \quad (6)$$

134 The numerator of Eq. (5) represents the square of signal of  $cc$ , which contributes to  
 135 particle detection, and the denominator represents the variance of noise in  $cc$  by using  
 136 the summation notation to calculate noise in  $cc$  and treating the average noise value of  
 137  $cc$  as zero since both  $N$  (Rosenthal & Henderson, 2003) and  $PN$  (Scheres, 2012) are  
 138 suggested as zero-mean variations. We simplified this calculation by including in only  
 139 two frequency terms ( $k_1$  and  $k_2$ ), and the corresponding variables are simplified to  $W_1$ ,  
 140  $N_1$ ,  $N'^*_1$ ,  $PN_1$ ,  $CTF_1$ , and  $S_1$  for  $W(k_1)$ ,  $N(k_1)$ ,  $N'(k_1)$ ,  $PN(k_1)$ ,  $CTF(k_1)$  and  $S(k_1)$ , and  $W_2$ ,  
 141  $N_2$ ,  $N'^*_2$ ,  $PN_2$ ,  $CTF_2$ , and  $S_2$  for  $W(k_2)$ ,  $N(k_2)$ ,  $N'(k_2)$ ,  $PN(k_2)$ ,  $CTF(k_2)$   
 142

$$143 \quad SNR_{cc} = \left[ \frac{W_1 CTF_1 |S_1|^2 + W_2 CTF_2 |S_2|^2}{W_1 N_1 S_1^* + W_1 N_1'^* CTF_1 S_1 + W_1 N_1'^* N_1 + W_1 PN_1 CTF_1 S_1^* + W_1 PN_1 CTF_1 N_1'^* + W_2 N_2 S_2^* + W_2 N_2'^* CTF_2 S_2 + W_2 N_2'^* N_2 + W_2 PN_2 CTF_2 S_2^* + W_2 PN_2 CTF_2 N_2'^*} \right]^2. \quad (7)$$

144 and  $S(k_2)$  denotes the absolute value, ignoring the cross terms (noise variables are  
 145 treated as random noise), and assuming that noise  $N'$  present in 3D template is much  
 146 smaller than noise  $N$  present in 2D raw image, the result is

$$147 \quad W(k) = \frac{CTF(k) \cdot FSC(k)}{\left| N(k) \right|^2 + \left| PN(k) \right|^2 \cdot CTF^2(k)}. \quad (8)$$

148 Here  $FSC(k)$  describes the Fourier Shell Correlation (FSC) between the perfect model  
 149 and the 3D template, which is introduced by (Rosenthal & Henderson, 2003) as

$$150 \quad FSC(k) = \frac{SNR(k)}{1 + SNR(k)}. \quad (9)$$

151 In this situation,  $SNR(k)$  equals the ratio of  $|S(k)|^2$  to  $|N'(k)|^2$ .

152 To remove the impacts caused by structure factor and the B-factor damping in  
 153 images, phase-flipped images and projections were signal-whitened.

154 
$$M(k) = \frac{M(k)}{\sqrt{|S(k)|_r^2}}, \quad (10)$$

155 
$$X(k) = \frac{X(k)}{\sqrt{|S(k)|_r^2}}. \quad (11)$$

156 The term  $|S(k)|_r^2$  is the radially averaged signal intensity, so the weighting function  
 157 applied to signal-whitened images was

158 
$$W(k) = \frac{CTF(k) \cdot FSC(k)}{\frac{1}{SSNR(k)} + \frac{|PN(k)|^2 \cdot CTF^2(k)}{|S(k)|_r^2}} \quad (12)$$

159 Assuming that the structural factors of different proteins are similar, the ratio of  
 160 protein noise intensity to projection intensity could be approximated as a constant  $n$ ,  
 161 which describes how proteins are overlapped. Our optimized weighting function is  
 162 simplified to

163 
$$W(k) = \frac{CTF(k) \cdot FSC(k)}{\frac{1}{SSNR(k)} + n \cdot CTF^2(k)} \quad (13)$$

164  
 165 In the absence of protein noise  $PN(k)$ , weighting function is calculated as

166 
$$W(k) = CTF(k) \cdot FSC(k) \cdot SSNR(k) \quad (14)$$

167 The SSNR used in the weighting function was determined from the power spectrum  
 168 of the image by measures the ratio of CTF affected contrast to non-CTF affected  
 169 contrast. First we fitted the non-CTF affected contrast presented in the power spectrum  
 170 by fitting zero points of CTF with an exponential function. As in situ sample is usually  
 171 thick ( $>100\text{nm}$ ), a significant amount of none-CTF modulated noise is introduced by  
 172 inelastic scattering electrons. Since inelastic electrons can be removed by energy filter  
 173 which is recommended for data collection on thick sample, we did not take this into  
 174 account. We calculated the differences between crest values and trough values induced  
 175 by CTF oscillations on power spectrum to calculate the CTF- affected contrast. The  
 176 CTF oscillation at high frequency range can be easily weakened by defocus variation  
 177 in the image caused by either the distribution of proteins along the incident electron  
 178 beam in thick sample or a slight unintentional tilt of the grid. Thus,, we used values  
 179 below  $1/8 \text{ \AA}^{-1}$  for fitting the other exponential function. The two sets of exponential  
 180 function were used to calculate SSNR in the later processing.

181  
 182 Weighting functions from a typical cryo-EM micrograph with different values of  $n$

183 applied were calculated and plotted in Fig. 1. Here, SSNR was estimated from the  
184 power spectrum of the image (see Materials and Methods) and  $FSC$  was set to 1. An  
185 increase of  $n$  indicates more noise from overlapping protein densities, which decreases  
186 the SNR. However, the SNR at low frequency range decreases faster than that at high  
187 frequency range. Therefore, the weight of score at high frequency range increases along  
188 with an increase of  $n$ . Moreover, the shape of the peak of the oscillation of the weighting  
189 function expanded in  $x$  direction with an increasing  $n$ . When  $n$  is much larger than  
190  $1/SSNR(k)$  (Fig. 1,  $n = 999$ ) (i.e. ignoring the shot noise), our weighting function is  
191 similar to a whitening filter (Rickgauer et al., 2017). Under different conditions when  
192  $n$  is zero (e.g. ignoring the overlapping densities), our weighting function is a dot  
193 product of  $CTF(k)$  (absolute value when applied to our signal-whitened images),  
194  $SSNR(k)$  and  $FSC(k)$  (see Materials and Methods). Without considering the overlapping  
195 densities, similar weighting functions have been used in CisTEM (Grant, Rohou, &  
196 Grigorieff, 2018) and for CTF refinement in RELION (J. Zivanov et al., 2018). In our  
197 picking function,  $n$  is normally estimated to 3 or 4 in real cryo-EM data (see below).

## 198 **2.2. Particle detection efficiency**

199 We tested our picking function by finding different protein complexes of different  
200 sizes in cryo-EM images of different icosahedral viruses (HSV, Alphavirus, Reovirus).  
201 In these viruses, protein complexes are overlapped with densities from other proteins  
202 and genome of the virus, which mimics an *in situ* environment. High resolution single  
203 particle analysis (Yuan et al., 2018) (L. Chen et al., 2018) on these viruses ensured an  
204 accurate determination of rotational and translational parameters of each 2D image. We  
205 extracted the center and the orientation of target protein complexes on the virus and  
206 used these known parameters as positive controls. Thirty virus particle images from  
207 each of the HSV-2 and Reovirus datasets were selected for testing. Before applying a  
208 weighting function, 2D images were signal-whitened. Projections of perfect 3D model  
209 were whitened first, then the square root of each weighting function was applied to both  
210 2D images and projections. The projections of the initial model were generated by  
211 *EMAN* (Ludtke, Baldwin, & Chiu, 1999) with the incremental rotation angle set to 5  
212 degrees. Images and projections were binned by 2 to reduce the CPU hours. For both  
213 HSV-2 and Reovirus datasets, any results with the translational error larger than 5 pixels  
214 or the rotational error larger than 6 degrees were considered as false positive results.

215 In our weighting function,  $n$  regulates the ratio of protein density noise to shot noise.  
216 To pick these protein complexes, we tested different picking functions, with  $n$  ranging  
217 from 0 to 50. The results are displayed as precise-recall curves (Figs. 2B and C), where  
218 precision is the ratio of true detections to false positives and recall is the ratio of detected  
219 particles to total particles presented. When  $n$  increases from 0 to 3, the ratio of correct  
220 results increases in the datasets of HSV-2 and Reovirus. The ratio becomes steady when  
221  $n$  is between 3 and 6. Thus, we used 3 as default value of  $n$  in the picking function. The  
222 ratio decreases when  $n$  increases further. The precise-recall curve of the whitening filter  
223 where  $n$  approaches infinity shows a much lower ratio of correct results compared with  
224 the curve of  $n = 3$  (Figs. 2B and C).

225 Our picking function was tested on recognizing protein complexes of different  
226 molecular weights on Alphavirus data. As shown in Fig. 2D, the ratio of correct results

227 versus all results decreases quickly as the size of the protein complexes decreases. This  
228 finding suggested that our picking function also generates a significant amount of  
229 incorrect particles. These false positive results are very similar to the model for picking  
230 according to the high score of our picking function, which will induce model biased  
231 feature in the reconstruction in the frequency that used for picking and noise in the  
232 structure in higher frequency in the refinement procedure.

233 When we are detecting a specific projection with global orientations and locations,  
234 all of the other densities (including targets at other orientations and other kinds of  
235 proteins) are contributed as noise. Therefore the SNR of low frequency signal is ruined,  
236 which makes the experiment we did on viruses different from regular single particle  
237 analysis. We think the high abundance of targets is the main superiority when compared  
238 with cellular environment and that may have an effect on detection efficiency. To  
239 investigate this, we did simulations on Alphavirus data set.

240 Assuming that different ratios of target protein are randomly removed from the  
241 micrographs to change the abundance of the protein, although the recall of the protein  
242 (percentage of the target proteins that are picked) remains unchanged using a fixed  
243 cutoff threshold of CC, the number of picked target protein decreases along with the  
244 decreasing of abundance. Since the background noise remains almost the same, we  
245 assume that the number of false positive result remains the same too. In such a way, we  
246 simulated the precision-recall curves of a 960kD protein with different abundances as  
247 shown in Fig. 2E. Here, when the precision drops to 0.1, the recall of the 960kD protein  
248 with 200 times decreasing of abundance is similar to that of a 480kD protein. Therefore,  
249 the size of the protein is a more important factor than the abundance.

### 250 **2.3.Application on test data**

251 To investigate the ability of our method to determine the protein structure in a  
252 crowded environment, we selected a part of the HSV-2 capsid as a target protein. This  
253 part is approx. 900 kD in molecular weight and consists of VP5 and VP26 trimetric  
254 capsid protein as well as the surrounding triplex (two copies of VP23 and one copy of  
255 VP19C) (Yuan et al., 2018).

256 We set  $n$  to 4 and used the frequencies ranging from  $1/100 \text{ \AA}^{-1}$  to  $1/8 \text{ \AA}^{-1}$  for particle  
257 detection. Possible locations and orientations of the target protein complex were  
258 calculated by our picking function and sorted according to the score of our picking  
259 function. After merging the results with similar orientations (within 7 degrees) in  
260 neighbor locations (within 10 pixels) into a single result, the first 500 putative target  
261 protein complexes with highest score from each virus particle were selected for further  
262 data processing, in which around 88% were false positive results according to the  
263 criteria we set in the Methods section. We performed 3D classification in RELION  
264 skipping alignment using the centers and orientations provided by our picking function.  
265 As shown in Fig. 3A, 2 out of 10 classes containing the lowest percentages of false  
266 positive result were selected, among which ~50% were false positive results. Further  
267 classification failed to improve the ratio of correct result.

268 Next, we performed auto refinement using only local search, which resulted in a  
269 reconstruction at  $4.3 \text{ \AA}$ -resolution. The FSC curve as shown in Fig.3B drops quickly at  
270 the frequency of  $\sim 1/8 \text{ \AA}^{-1}$  and exhibits a shoulder around the frequency of  $1/5 \text{ \AA}^{-1}$ ,

271 indicating a reference bias problem (below  $1/8 \text{ \AA}^{-1}$ ) and noise problem (above  $1/8 \text{ \AA}^{-1}$ )  
272 in the reconstruction. We also tried to use fewer putative target protein complexes at  
273 the top of the sorting list to increase the ratio of correct results. However, this approach  
274 also decreased the number of correct results. When we reduced the number of putative  
275 target protein complexes in the refinement procedure, the resolution of the  
276 reconstruction improved (Fig. S2) presumably due to the increase of the ratio of correct  
277 result before decreasing because of the lack of particles.

278 To further reduce the ratio of false positive results, we calculated the score between  
279 reference and the raw image according to refined parameters using phase residual  
280 (Methods and Materials). For this calculation, we only used the frequencies ranging  
281 from  $1/8 \text{ \AA}^{-1}$  to  $1/5 \text{ \AA}^{-1}$ . The particles were sorted according to their score. As shown in  
282 Fig. 3E, the sorting efficiently separated the correct results from false positive results  
283 by two Gaussian-like peaks. When the sorting was based on the score from our picking  
284 function using the frequencies ranging from  $1/20 \text{ \AA}^{-1}$  to  $1/8 \text{ \AA}^{-1}$ , the power of separation  
285 decreased markedly (Fig. 3D). This range of frequency was involved in particle picking,  
286 which performed a global search of location and orientation of the target protein on a  
287 whole virus. For instance, in a combination of translational (step size of  $2.76 \text{ \AA}$ ) and  
288 rotational parameters (step size of 5 degrees),  $7.5 \times 10^{10}$  possible locations of the  
289 protein complex were searched, from which the top 500 possible locations were  
290 selected by the program. Thus, each false positive result was selected from  $7.5 \times 10^{10}$   
291 possible locations. However, further refinement using only local search improved the  
292 resolution to  $\sim 4.3 \text{ \AA}$ . The local search strictly limited the possible locations. Thus, in  
293 the range of frequencies from  $1/8 \text{ \AA}^{-1}$  to  $1/4.3 \text{ \AA}^{-1}$ , the false positive result exhibits  
294 differences from the model. Therefore, excluding the range of frequencies that involves  
295 in picking function for sorting exhibits less reference bias. In addition, the refinement  
296 only performing local search was based on maximum likelihood score function in  
297 RELION. We tested “MaxValueProbDistribution” generated in RELION as score to  
298 sort the particles; however, the correct results were barely differentiated from the false  
299 positive results (Fig. 3C). The sorting according to parameter of  
300 “NrOfSignificantSamples” exhibited similar result to that of  
301 “MaxValueProbDistribution”. Thus, it is possible that using scores different from the  
302 one used in the refinement for sort also helps to reduce the false positive detections.

303 After sorting by the score calculated using only the frequencies from  $1/8 \text{ \AA}^{-1}$  to  $1/5$   
304  $\text{ \AA}^{-1}$ , 40,000 particles (the top 40% particles contained  $\sim 93\%$  correct results) were  
305 selected. Further refinement of this dataset led to a resolution of  $4.0 \text{ \AA}$  (Fig. 3B). By  
306 adding in a further 6,000 viral particles, 180,000 particles were selected after non-  
307 alignment 3D classification and sorting. The gold standard resolution was  $3.7 \text{ \AA}$ . By  
308 combining with CTF refinement in RELION using our optimized weighting function,  
309 resolution was improved to  $3.5 \text{ \AA}$ . The original CTF refinement procedure implanted  
310 in RELION produced a map of  $3.6 \text{ \AA}$  resolution.

#### 311 **2.4.Determining protein structure using a homologous structure as a picking** 312 **model**

313 Since the expected structure from *in situ* structure determination is usually unknown,  
314 we explored the possibility of using homologous structure as a picking model. The

315 differences between homologous structure and expected structure were treated as noise  
316 ( $N'$ ), resulting in a different term named  $FSC_m$  (between the cryo-EM map of  
317 homologous structure and the expected structure) in the weighting function to replace  
318 the  $FSC$ .

$$319 \quad W(k) = \frac{CTF(k) \cdot FSC_m(k)}{\frac{1}{SSNR(k)} + n \cdot CTF^2(k)} \quad (16)$$

320 To search for a similar protein complex on HSV-2 capsid core, we used a  
321 homologous protein complex present in the HSV-1 capsid core as a model. First, we  
322 extracted the 3D model of the homologous protein complex from a 4.2 Å map of HSV-  
323 1 (Dai & Zhou, 2018). The FSC curve between protein complexes from HSV-1 and  
324 HSV-2 showed similarity in the structures with the FSC value decreasing to 0.7 at the  
325 frequency of  $1/8 \text{ \AA}^{-1}$ . Three million potential particles of protein complex on HSV-2  
326 capsid were selected on the basis of the score produced by our picking function. After  
327 local 3D classification and further selection by sorting, 60,000 particles were finally  
328 selected. As shown in Fig. 4B, the local refinement resulted in a 4.0 Å resolution map.  
329 We calculated the FSC curves between this map and the corresponding 3.1 Å map from  
330 the single particle result of HSV-2 and between this map and the corresponding 4.2 Å  
331 map of HSV-1 (Dai & Zhou, 2018). As shown in Fig.4C and D, our refined structure is  
332 closer to the corresponding structure in HSV-2 than that in HSV-1. Assuming that the  
333 3.1 Å map represents a perfect map, the FSC between 3.1 Å map and our map reported  
334 a resolution of 4.0 Å using a threshold of 0.5. This result is in agreement with the  
335 resolution (4.0 Å) reported by FSC between two half maps using a threshold of 0.143.  
336 Together, these results show that the procedure we used in the *in situ* reconstruction  
337 avoids reference bias induced by homologous models.

338 To evaluate the ability of finding protein complexes by homologous models of  
339 different similarity between their structures and the structure of the target protein  
340 complex, we simulated homologous models by applying different scale factors to the  
341 structure of the target protein. As shown in Fig. S3, the similarities between the re-  
342 scaled model and the original map are indicated by FSC curves. The precision-recall  
343 curves show that the ability of finding a protein complex decreases along with the  
344 reduction of the similarity between homologous model and the protein complex (Fig.  
345 4E and F). We downloaded PDB files of different homologous proteins and calculated  
346 the potential density maps, and then calculated the FSCs between pairs of homologous  
347 maps. As shown in Fig.S4, the similarity between proteins varies greatly. These  
348 homologous proteins with high similarities can be used as picking models in our method.  
349 **2.5.Determining *in situ* structures at high resolution using isSPA**

350 To evaluate our *in situ* structure determination program, we first processed a dataset  
351 of a Bunyavirus by 2D and 3D classification. As shown in Fig. 5B, only approx. 28%  
352 of the viral particles (3,140 particles) exhibited an icosahedral symmetry. Further  
353 refinement on the icosahedral viral particles resulted in an 11.8 Å map due to the  
354 flexibility. To compensate for this limitation in flexibility, one block (one pentamer and  
355 5 surrounding hexamers) centered on a pentamer was extracted and refined (Shaikh,

356 Hegerl, & Frank, 2003), which led to a map of 8.7 Å resolution. A centered sub-block  
357 (pentamer) and a sub-block adjacent to the pentamer that centered on a hexamer  
358 segmented from the 8.7 Å block were used as 3D models to pick the protein complex  
359 on the non-icosahedral virus particles previously excluded from data processing.  
360 Through global detection by isSPA, twenty potential solutions were selected from each  
361 virus and ~27,000 particles in total were selected according to 3D non-alignment  
362 classification resulting in a 7.7 Å resolution map. In the search for hexamers, one  
363 hexamer from the 5 surrounding a pentamer was segmented and set as the model, and  
364 200 potential solutions were selected from each virus. After 3D non-alignment  
365 classification, ~87,000 particles were selected and refined to 8.3 Å resolution.

366 In the second sample, carboxysomes were purified from *Halothiobacillus*  
367 *neapolitanus* and the cryo-EM data were collected on this sample. The size of the  
368 carboxysomes ranged from 100 nm to 150 nm. We also purified Rubisco from fractured  
369 carboxysomes by adding extra freeze-thaw cycles before the centrifugation step for  
370 carboxysome purification. We collected the cryo-EM data of purified Rubisco, and  
371 obtained a 2.7 Å map of this complex (around 500 kD). This complex was then used as  
372 a model to pick the Rubisco inside the carboxysome in the cryo-EM images. The  
373 frequencies between 1/100 Å<sup>-1</sup> to 1/8 Å<sup>-1</sup> were used for picking. From the collected  
374 micrographs, ~5,200 carboxysome images were extracted, and each image was cross-  
375 correlated with the 2.7 Å model at a frequencies range from 1/100 Å<sup>-1</sup> to 1/8 Å<sup>-1</sup>. This  
376 procedure yielded a large group of locations and orientations, from which on average  
377 150 solutions were picked per carboxysome for further processing. The non-alignment  
378 3D classification was performed using RELION. Approximately 150,000 particles  
379 were selected, and a further auto local refinement with angular sampling of 0.9 degrees  
380 and translational sampling of 1.04 Å in RELION reported a 4.3 Å resolution map. To  
381 remove the false positive detections, the refined particles were sorted with phase  
382 residual using frequencies from 1/8 Å<sup>-1</sup> to 1/5 Å<sup>-1</sup>. We tested three cutting thresholds  
383 (0.07, 0.08 and 0.09) and selected ~82,000, ~46,000 and ~21,000 particles to the  
384 second-round refinement individually, resulting in 4.0 Å, 3.9 Å and 3.9 Å resolution.  
385 Further CTF refinement subsequently improved the resolutions to 3.9 Å, 3.9 Å and 3.7  
386 Å. According to the fitting of two Gaussian distributions (Fig. 5H), the portion of true  
387 solutions at threshold 0.09 was estimated to ~90%.

388 The success on Rubisco data set may owe to the high symmetry (D4) and its high  
389 abundance. We have shown above that our work flow can solve protein picked with  
390 precision below 0.1. According to the description of section 3.1, the recall of the 960kD  
391 protein with 200 times decreasing of abundance keeps the similar as a 480kD protein  
392 at precision 0.1. Thus, it is possible that a ~1MD protein in Carboxysome thickness  
393 (~120 nm) with 200 times lower abundance than that of Rubisco can be solved at a  
394 pseudo atomic resolution.

395

### 396 **3. Discussion**

397 In this work, we showed that the image-processing routines of isSPA method can  
398 resolve proteins (larger than 400 kD) at high resolution *in situ* with non-cryo-sectioning  
399 data when the thickness of the sample is around 120nm. The ability of finding the initial

400 rotational and translational parameters of the target protein is closely associated with  
401 the size of the target protein and the thickness of the cryo-EM sample. Our recent results  
402 show that 50% of ~350kD membrane proteins can be found and reconstructed to pseudo  
403 atomic resolution on 50nm-diameter liposomes. We are also able to reconstruct a 1.2MD  
404 membrane protein containing only 150kD soluble domain staying on the original  
405 cellular membrane formed liposome with an average diameter of ~150nm at high  
406 resolution. Thus, isSPA method can very efficiently help to solve structures of different  
407 kinds of membrane proteins on liposome constituted by their native lipid membrane.  
408 Considering the thickness of cryo-sectioning sample and the abundance of the target  
409 protein, it is worth to try on a protein with size ~ 1MD using this method. Low SNR is  
410 the major obstacle of determining the initial parameters of protein complexes. Thus,  
411 this method will benefit greatly from hardware improvement such as better direct  
412 detectors or new generation of phase plate that improves the SNR of images at  
413 frequencies ranging from  $1/20 \text{ \AA}^{-1}$  to  $1/8 \text{ \AA}^{-1}$  (0.1~0.25 Nyquist at a pixel size of  $1 \text{ \AA}$ )  
414 or better tool to obtain a thinner cryo-sectioning.

415 When the protein is imaged in the in situ environment, the overlapping densities  
416 ruin the SNR of low frequency signals of the target protein in the images. The SPA  
417 algorithms improve the 3D structure iteratively from a low resolution 3D map to a high  
418 resolution 3D map. However, such a method become invalid when the low frequency  
419 signals of the target protein was ruined by overlapping densities. Thus, routine non-  
420 reference 2D classification and 3D classification starting from a low resolution initial  
421 model usually converge to wrong results. Thus, starting from the local refinement in  
422 both structure refinement and 3D classification is strongly suggested in isSPA. The  
423 structural information of a template with frequency above  $1/8 \text{ \AA}^{-1}$  is involved in picking  
424 the particle, therefore the FSC curves below  $1/8 \text{ \AA}^{-1}$  is not calculated from two  
425 completely independent datasets. However, the FSC beyond  $1/8 \text{ \AA}^{-1}$  is not affected by  
426 the template and can be considered as gold-standard.

427 Generally, isSPA is not a good choice for mapping protein complexes in in situ  
428 environment due to the relatively low detection efficiency. Moreover, the z height of  
429 the protein in the in situ environment is determined by its defocus value in the  
430 micrograph. Although, the per-particle CTF refinement improves the defocus value of  
431 the protein, error remains. It prevents an accurate localization of the target protein along  
432 Z direction. Thus, if the distribution of target protein is unknown, we suggest that a  
433 tomographic study prior to isSPA can be used to show the distribution of target protein  
434 in the in situ environment and provide a medium resolution template if it is necessary.  
435 Then, isSPA can be used to further improve the resolution in an efficient way.

436

## 437 **4. Materials and Methods**

### 438 **4.1.Reducing the false positive results by sorting**

439 The introduction of high resolution ( $1/8 \text{ \AA}^{-1}$ ) information in global detection will  
440 generate a large amount of false positive detections due to the especially low SNR.  
441 These false positives are presented as noise beyond  $1/8 \text{ \AA}^{-1}$  and produces biased features  
442 in the reconstruction below  $1/8 \text{ \AA}^{-1}$  which affect the alignment against the  
443 reconstruction, thus prevent from pushing high resolution in the following refinement

444 procedures. To further reduce the false positive results, we sorted the selected images  
445 based on a different score. The signals with frequency range beyond  $1/8 \text{ \AA}^{-1}$  are  
446 recommended for calculating the score since they are less affected by biased feature.

$$447 \quad \text{score} = FSC \cdot CTF \cdot \cos(\Delta\phi) \quad (15)$$

448 Due to the extremely low SNR in high frequencies range, the constant  $n$  is much smaller  
449 than  $1/SSNR$ , thus we ignored the overlapping density and used  $FSC$  and  $CTF$  as  
450 weighting factors.

#### 451 **4.2. Workflow of isSPA method**

452 IsSPA is a single-particle like method to reconstruct *in situ* structures from untilted  
453 images. The images therefore contain overlapping protein densities, and the workflow  
454 includes: (1) estimate SSNR parameters by fitting power spectrum with two  
455 exponential functions (one for noise and the other for signal); (2) for particle picking,  
456 raw images are filtered by our optimized weighting function with known SSNR  
457 parameters, then the filtered images and projections of 3D model are input to program;  
458 (3) the output of program presents a large group of locations and orientations, thus we  
459 merge the neighbor solutions according to their translations and rotations; (4) we sort  
460 solutions from each image by cc value, and the number of solutions are selected  
461 depending on the estimated particle number in each image; (5) scripts were used to  
462 extract potential particles and perform a local refine and 3D classification on these  
463 particles using different cryo-EM software packages such as RELION, JSPR and so on;  
464 (7) a sorting algorithm was provided to remove false positive particles; (8) the  
465 remaining particles are performed local refinement to achieve the final reconstruction.

#### 466 **4.3. Image acquisition of test data**

467 Images of Herpes simplex virus type 2 particles (HSV-2) (~120 nm in diameter) and  
468 Reovirus particles (~90 nm in diameter) obtained from previous studies were analyzed  
469 for testing. HSV-2 data were collected on an FEI Titan Krios microscope operated at  
470 300 kV equipped with a Falcon2 camera, with pixel size of  $1.38 \text{ \AA}/\text{pixel}$  and a total  
471 dose of  $\sim 25 \text{ e}^-/\text{\AA}^2$ . Reovirus data was acquired on an FEI Titan Krios microscope  
472 operated at 200 kV with a Gatan K2 Summit camera, and the pixel size is  $1.04 \text{ \AA}/\text{pixel}$   
473 with a total dose of  $\sim 30 \text{ e}^-/\text{\AA}^2$ .

#### 474 **4.4. Image acquisition of Bunyavirus**

475 Bunyavirus data were collected on an FEI Titan Krios EM operated at 300 kV  
476 equipped with a Gatan K2 Summit detector. Data collection was performed with Serial  
477 EM using a nominal magnification of 22,500 with a total dose of  $\sim 50 \text{ e}^-$  resulting in a  
478 pixel size of  $1.3 \text{ \AA}/\text{pixel}$ . Each micrograph was recorded as a movie composed of 36  
479 frames. A total of 1,849 such movies were collected and 10,950 viral particles were  
480 extracted from motion-corrected micrographs.

#### 481 **4.5. Image acquisition of Carboxysome**

482 Purified carboxysomes were applied to a glow-discharged Quantifoil R1.2/1.3 300-  
483 mesh copper holey carbon grid (Quantifoil, Micro Tools), blotted under 100% humidity  
484 at  $4^\circ\text{C}$  and plunged into liquid ethane using a Mark IV Vitrobot (FEI). Micrographs  
485 were acquired on a Titan Krios G2 microscope (FEI) operated at 300 kV with a K2  
486 Summit direct electron detector (Gatan). A calibrated magnification of 130k times was

487 used for imaging, yielding a pixel size of 1.04 Å on images. The defocus was set at 1.5  
488 to 2.0 µm. Each micrograph was dose-fractionated to 32 frames under a dose rate of 8  
489 e<sup>-</sup>/pixel/s with a total exposure time of 10.4 s.

490

#### 491 **Data and materials availability**

492 Data that support the findings of this study is available from the authors upon  
493 request.

494

#### 495 **References:**

496 Bäuerlein, F. J. B., Saha, I., Mishra, A., Kalemánov, M., Martínez-Sánchez, A., Klein, R., . . . Fernández-  
497 Busnadiego, R. (2017). In Situ Architecture and Cellular Interactions of PolyQ Inclusions. *Cell*,  
498 *171*(1), 179-187. doi:10.1016/j.cell.2017.08.009

499

500 Beck, M., & Baumeister, W. (2016). Cryo-Electron Tomography: Can it Reveal the Molecular Sociology of  
501 Cells in Atomic Detail? *Trends Cell Biol*, *26*(11), 825-837. doi:10.1016/j.tcb.2016.08.006

502

503 Bykov, Y. S., Schaffer, M., Dodonova, S. O., Albert, S., Plitzko, J. M., Baumeister, W., . . . Briggs, J. A. (2017).  
504 The structure of the COPI coat determined within the cell. *Elife*, *6*. doi:10.7554/eLife.32493

505

506 Chen, L., Wang, M., Zhu, D., Sun, Z., Ma, J., Wang, J., . . . Zhang, X. (2018). Implication for alphavirus  
507 host-cell entry and assembly indicated by a 3.5Å resolution cryo-EM structure. *Nat Commun*,  
508 *9*(1), 5326. doi:10.1038/s41467-018-07704-x

509

510 Chen, M., Bell, J. M., Shi, X., Sun, S. Y., Wang, Z., 0000-0003-4897-9986, O., & Ludtke, S. J. (2019). A  
511 complete data processing workflow for cryo-ET and subtomogram averaging. *Nat Methods*,  
512 *16*(11), 1161-1168. doi:10.1038/s41592-019-0591-8

513

514 Chreifi, G., Chen, S., Metskas, L. A., Kaplan, M., & Jensen, G. J. (2019). Rapid tilt-series acquisition for  
515 electron cryotomography. *J Struct Biol*, *205*(2), 163-169. doi:10.1016/j.jsb.2018.12.008

516

517 Dai, X., & Zhou, Z. H. (2018). Structure of the herpes simplex virus 1 capsid with associated tegument  
518 protein complexes. *Science*, *360*(6384). doi:10.1126/science.aao7298

519

520 Dodonova, S. O., Aderhold, P., Kopp, J., Ganeva, I., Röhlíng, S., Hagen, W. J. H., . . . Briggs, J. A. G. (2017).  
521 9Å structure of the COPI coat reveals that the Arf1 GTPase occupies two contrasting molecular  
522 environments. *Elife*, *6*. doi:10.7554/eLife.26691

523

524 Freeman Rosenzweig, E. S., Xu, B., Kuhn Cuellar, L., Martinez-Sanchez, A., Schaffer, M., Strauss, M., . . .  
525 Jonikas, M. C. (2017). The Eukaryotic CO(2)-Concentrating Organelle Is Liquid-like and Exhibits  
526 Dynamic Reorganization. *Cell*, *171*(1), 148-162. doi:10.1016/j.cell.2017.08.008

527

528 Grant, T., Rohou, A., & Grigorieff, N. (2018). cisTEM, user-friendly software for single-particle image  
529 processing. *Elife*, *7*. doi:10.7554/eLife.35383

530

531 Guo, F., & Jiang, W. (2014). Single particle cryo-electron microscopy and 3-D reconstruction of viruses.  
532 In *Electron Microscopy* (pp. 401-443): Springer.  
533

534 Guo, Q., Lehmer, C., Martínez-Sánchez, A., Rudack, T., Beck, F., Hartmann, H., . . . Fernández-Busnadiego,  
535 R. (2018). In Situ Structure of Neuronal C9orf72 Poly-GA Aggregates Reveals Proteasome  
536 Recruitment. *Cell*, *172*(4), 696-705.e612. doi:10.1016/j.cell.2017.12.030  
537

538 Himes, B. A., & Zhang, P. (2018). emClarity: software for high-resolution cryo-electron tomography and  
539 subtomogram averaging. *Nat Methods*, *15*(11), 955-961. doi:10.1038/s41592-018-0167-z  
540

541 Leigh, K. E., Navarro, P. P., Scaramuzza, S., Chen, W., Zhang, Y., Castaño-Díez, D., & Kudryashev, M. (2019).  
542 Subtomogram averaging from cryo-electron tomograms. *Methods Cell Biol*, *152*, 217-259.  
543 doi:10.1016/bs.mcb.2019.04.003  
544

545 Lučić, V., Leis, A., & Baumeister, W. (2008). Cryo-electron tomography of cells: connecting structure and  
546 function. *Histochemistry and Cell Biology*, *130*(2), 185. doi:10.1007/s00418-008-0459-y  
547

548 Ludtke, S. J., Baldwin, P. R., & Chiu, W. (1999). EMAN: semiautomated software for high-resolution  
549 single-particle reconstructions. *J Struct Biol*, *128*(1), 82-97. doi:10.1006/jsbi.1999.4174  
550

551 Mattei, S., Tan, A., Glass, B., Müller, B., Kräusslich, H.-G., & Briggs, J. A. G. (2018). High-resolution  
552 structures of HIV-1 Gag cleavage mutants determine structural switch for virus maturation.  
553 *Proc Natl Acad Sci U S A*, *115*(40), E9401-e9410. doi:10.1073/pnas.1811237115  
554

555 Mosalaganti, S., Kosinski, J., Albert, S., Schaffer, M., Strenkert, D., Salomé, P. A., . . . Beck, M. (2018). In  
556 situ architecture of the algal nuclear pore complex. *Nat Commun*, *9*(1), 2361.  
557 doi:10.1038/s41467-018-04739-y  
558

559 Pfeffer, S., Burbaum, L., Unverdorben, P., Pech, M., Chen, Y., Zimmermann, R., . . . Förster, F. (2015).  
560 Structure of the native Sec61 protein-conducting channel. *Nat Commun*, *6*, 8403.  
561 doi:10.1038/ncomms9403  
562

563 Rickgauer, J. P., Grigorieff, N., & Denk, W. (2017). Single-protein detection in crowded molecular  
564 environments in cryo-EM images. *Elife*, *6*. doi:10.7554/eLife.25648  
565

566 Rosenthal, P. B., & Henderson, R. J. J. o. M. B. (2003). Optimal Determination of Particle Orientation,  
567 Absolute Hand, and Contrast Loss in Single-particle Electron Cryomicroscopy. *4*(333), 721-745.  
568

569 Scheres, S. H. (2012). A Bayesian view on cryo-EM structure determination. *J Mol Biol*, *415*(2), 406-418.  
570 doi:10.1016/j.jmb.2011.11.010  
571

572 Schur. (2016). An atomic model of HIV-1 capsid-SP1 reveals structures regulating assembly and  
573 maturation. *Science*, *353*. doi:10.1126/science.aaf9620  
574

575 Schur, F. K. M., Hagen, W. J. H., Rumlová, M., Ruml, T., Müller, B., Kräusslich, H.-G., & Briggs, J. A. G.  
576 (2015). Structure of the immature HIV-1 capsid in intact virus particles at 8.8 Å resolution.  
577 *Nature*, 517(7535), 505-508. doi:10.1038/nature13838  
578  
579 Shaikh, T. R., Hegerl, R., & Frank, J. J. J. o. s. b. (2003). An approach to examining model dependence in  
580 EM reconstructions using cross-validation. *142*(2), 301-310.  
581  
582 Song, K., Shang, Z., Fu, X., Lou, X., Grigorieff, N., & Nicastro, D. (2019). In situ structure determination at  
583 nanometer resolution using TYGRESS. *Nat Methods*. doi:10.1038/s41592-019-0651-0  
584  
585 Turoňová, B., Schur, F. K. M., Wan, W., & Briggs, J. A. G. (2017). Efficient 3D-CTF correction for cryo-  
586 electron tomography using NovaCTF improves subtomogram averaging resolution to 3.4Å.  
587 *Journal of structural biology*, 199(3), 187-195. doi:10.1016/j.jsb.2017.07.007  
588  
589 Wan, W., Kolesnikova, L., Clarke, M., Koehler, A., Noda, T., Becker, S., & Briggs, J. A. G. (2017). Structure  
590 and assembly of the Ebola virus nucleocapsid. *Nature*, 551(7680), 394-397.  
591 doi:10.1038/nature24490  
592  
593 Yuan, S., Wang, J., Zhu, D., Wang, N., Gao, Q., Chen, W., . . . Liu, H. (2018). Cryo-EM structure of a  
594 herpesvirus capsid at 3.1 Å. *Science* 360: eaao7283. In.  
595  
596 Zhu, D., Wang, X., Fang, Q., Van Etten, J. L., Rossmann, M. G., Rao, Z., & Zhang, X. (2018). Pushing the  
597 resolution limit by correcting the Ewald sphere effect in single-particle Cryo-EM  
598 reconstructions. *Nat Commun*, 9(1), 1552. doi:10.1038/s41467-018-04051-9  
599  
600 Zivanov, J., Nakane, T., Forsberg, B. O., Kimanius, D., Hagen, W. J., Lindahl, E., & Scheres, S. H. (2018).  
601 New tools for automated high-resolution cryo-EM structure determination in RELION-3. *Elife*,  
602 7. doi:10.7554/eLife.42166  
603  
604 Zivanov, J., Nakane, T., Forsberg, B. O., Kimanius, D., Hagen, W. J., Lindahl, E., & Scheres, S. H. J. E. (2018).  
605 New tools for automated high-resolution cryo-EM structure determination in RELION-3. 7,  
606 e42166.  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618

619 **Acknowledgements**

620 We thank L. Kong for cryo-EM data storage and backup and Z. Y. Lou in Tsinghua  
621 University for offering us the cryo-EM data of Bunyavirus. Cryo-EM data collection  
622 was carried out at the Center for Biological Imaging at the Institute of Biophysics (IBP),  
623 Chinese Academy of Sciences (CAS). We thank X. Huang, B. L. Zhu, G. Ji and other  
624 staff members at the Center for Biological Imaging (IBP, CAS) for their support in data  
625 collection. The project was funded by the National Key R&D Program of China  
626 (2017YFA0504700), the Natural Science Foundation of China (31930069) the Strategic  
627 Priority Research Program of the Chinese Academy of Sciences (XDB37040101) and  
628 the Key Research Program of Frontier Sciences at the Chinese Academy of Sciences  
629 (ZDBS-LY- SM003), X.Z. received scholarships from the ‘National Thousand (Young)  
630 Talents Program’ from the Office of Global Experts Recruitment in China.

631 **Author contributions**

632 X.Z. and J.C. conceived and designed the study. J.C. developed and tested the  
633 methods. B.L. and L.S. prepared the sample of Carboxysome and collected the cryo-  
634 EM data. X.Z. and J.C. analyzed the results and wrote the paper.

635 **Competing interests**

636 All other authors declare no competing interests.

637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662

663 **Figures and Tables**

664

665 **Figure 1**

666

667

668

669

670

671

672

673

674

675

676

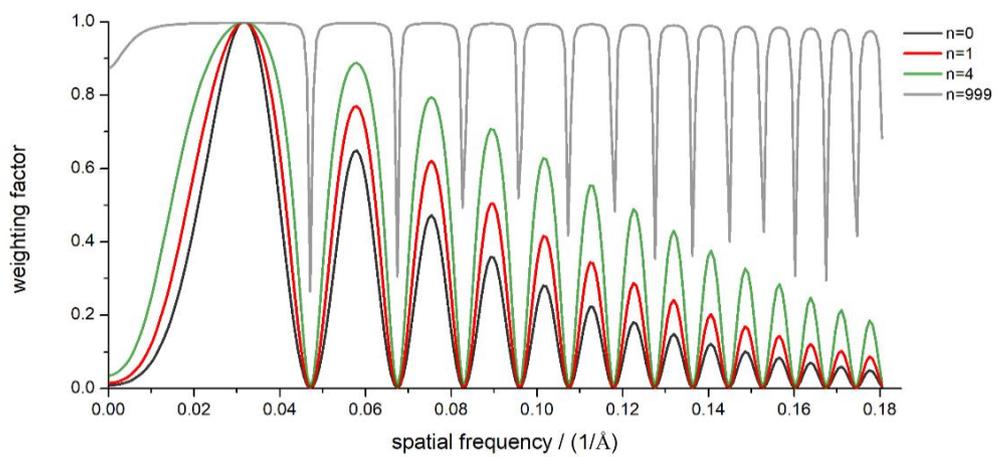
677

678

679

680

681



**Weighting functions.**

Weighting factors damp and oscillate with spatial frequency (black, red, green and gray) at different n (0, 1, 4, 999).

682 **Figure 2**

683

684

685

686

687

688

689

690

691

692

693

694

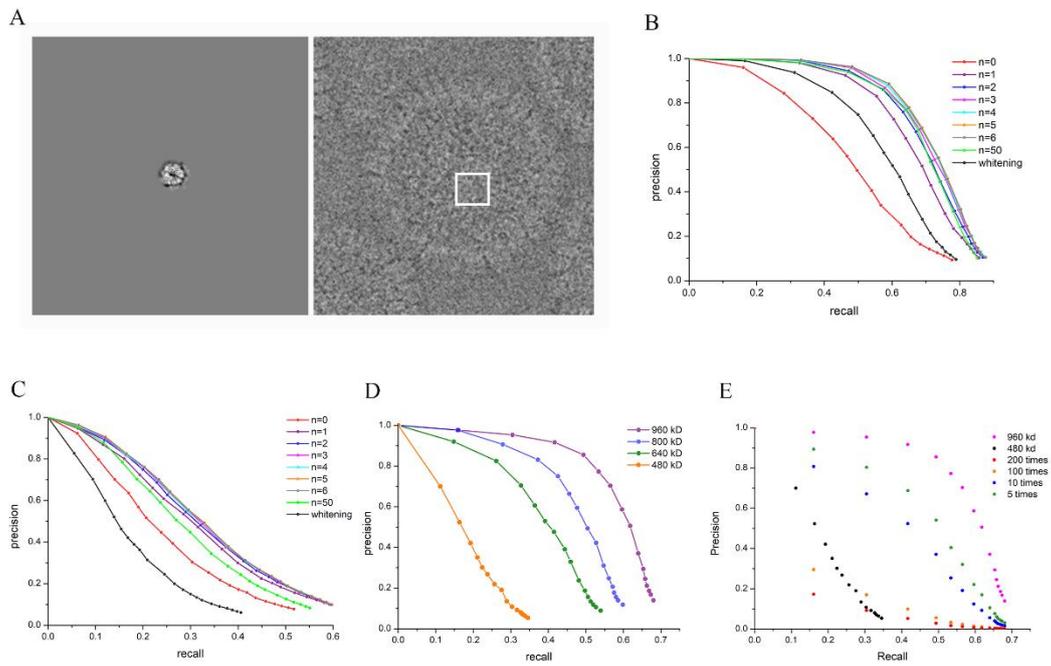
695

696

697

698

699



700 **Efficiency of particle detection.**

701 (A) Left: projection of HSV-2 hexamer model. Right: HSV-2 virus particle imaged at  $2.75 \mu\text{m}$  defocus with 25 electrons per  $\text{\AA}^2$ .

702 The location of the projection on the virus particles is indicated by a white square. (B) Precision-recall curves for detections on

703 Reovirus datasets using model in 900 kD molecular weight at  $n=0$  (red), 1 (purple), 2 (navy), 3 (magenta), 4 (cyan), 5 (orange), 6

704 (gray), 50 (green) and whitening filter (black). (C) Precision-recall curves for detections on HSV-2 datasets with the same

705 processing as in (B). (D) Precision-recall curves (purple, navy, green and orange) for detections ( $n$  equals 3) on Alphaviruses

706 using models in different molecular weights (960 kD, 800 kD, 640 kD and 480 kD). (E) Precision-recall curves for lower

707 abundance of a 960 kD protein on Alphavirus (olive for 5 times, navy for 10 times, orange for 100 times and red for 200 times

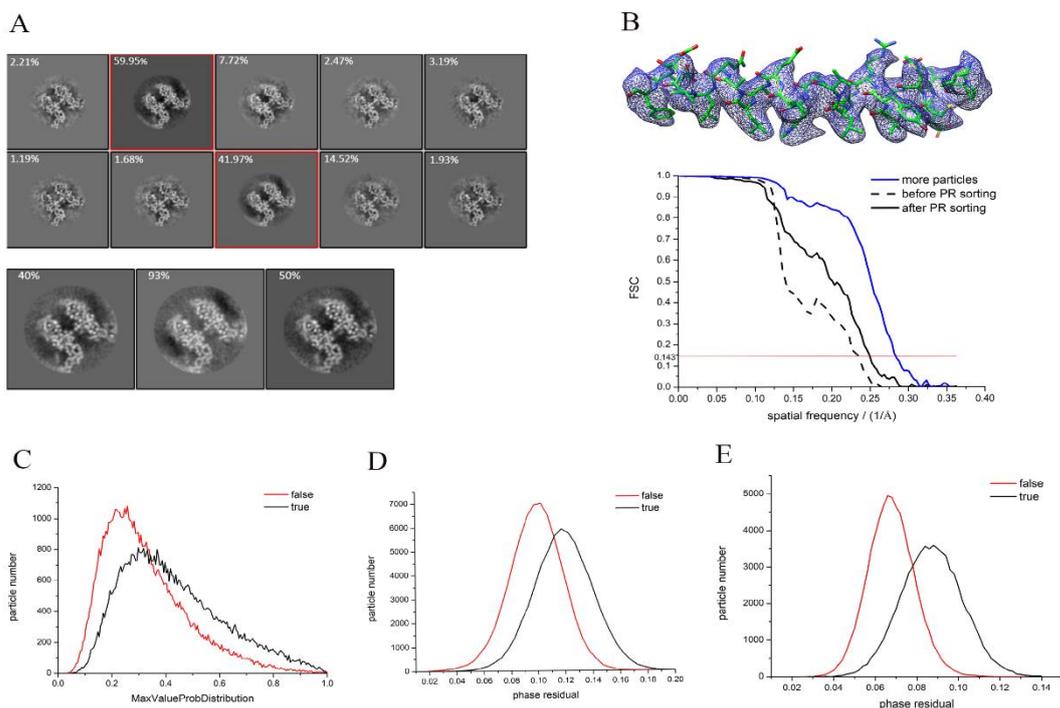
708 lower abundance), 960 kD at abundance of 60 copies per  $90 \text{ nm} * 90 \text{ nm}$  square (magenta) and 480 kD at abundance of 60 copies

709 per  $900 \text{ nm} * 900 \text{ nm}$  image (black).

710

711

712 **Figure 3**



732 **Data processing of HSV-2 hexamer.**

733 (A) Upper panel: 3D classification of raw picked positives in ten classes, the percentage of true detections in each class is shown  
 734 at the top left corner. The two selected classes are indicated by squares in red. Lower panel: 3D classification of selected particles  
 735 in three classes, and the percentage of true detections in each class is noted at the top left corner. (B) Upper panel: 3D reconstruction  
 736 of the HSV-2 hexamer. Lower panel: FSC curves denote 3 reconstructions in different conditions using 2000 viral particles before  
 737 PR sorting (dash) and after PR sorting (black), using 8000 virus particles and after PR sorting (navy). (C) True and false particle  
 738 distribution with *MaxValueProbDistribution* term in RELION. (D) True and false particle distribution with phase residual using  
 739 data from  $1/30 \text{ \AA}^{-1}$  to  $1/8 \text{ \AA}^{-1}$ . (E) True and false particle distribution with phase residual using data from  $1/8 \text{ \AA}^{-1}$  to  $1/5 \text{ \AA}^{-1}$ .

743 **Figure 4**

744

745

746

747

748

749

750

751

752

753

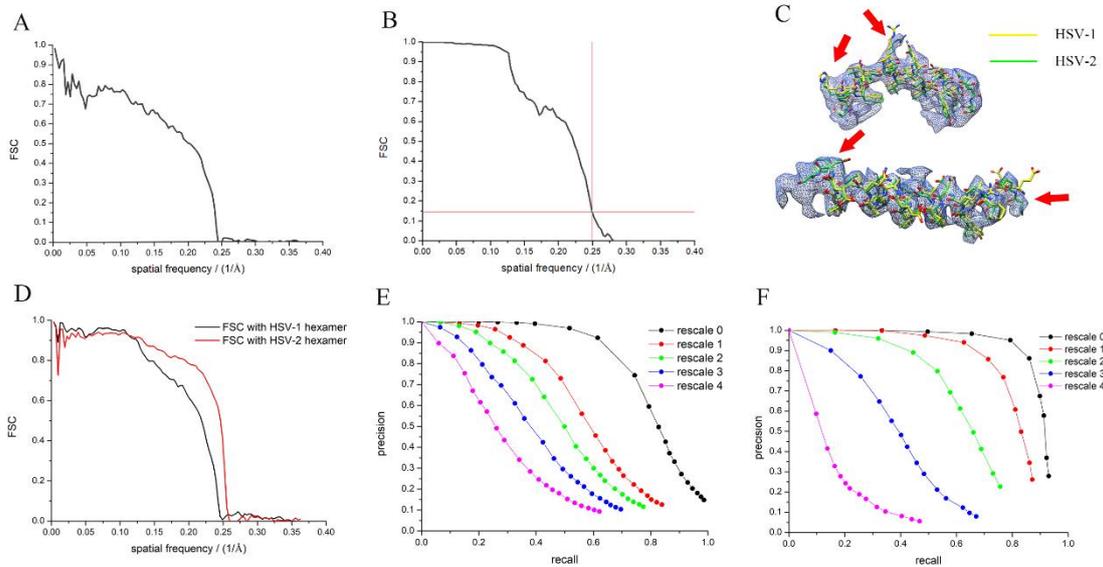
754

755

756

757

758



759 **Homologous structure as picking model.**

760 (A) FSC curve of HSV-1 hexamer with HSV-2 hexamer. (B) The resolutions determined by gold standard FSC at threshold 0.143

761 of HSV-2 hexamer reconstructed from HSV-1 hexamer. (C) PDB of HSV-1 hexamer (yellow) and HSV-2 (green) fitting to the

762 density of the 4.0 Å cryo-EM map, the disagreements of the two PDBs are pointed out by red arrows. (D) FSC curves show

763 similarities between our map and HSV-1 hexamer (black), and HSV-2 hexamer (red) respectively. (E) Precision-recall curves of

764 the ~2 MD protein of HSV-2 at a series of scales corresponding to Fig. S3. (F) Precision-recall curves of the ~1.8 MD protein of

765 Reovirus at a series of scales corresponding to Fig. S3.

766

767

768

769

770

771

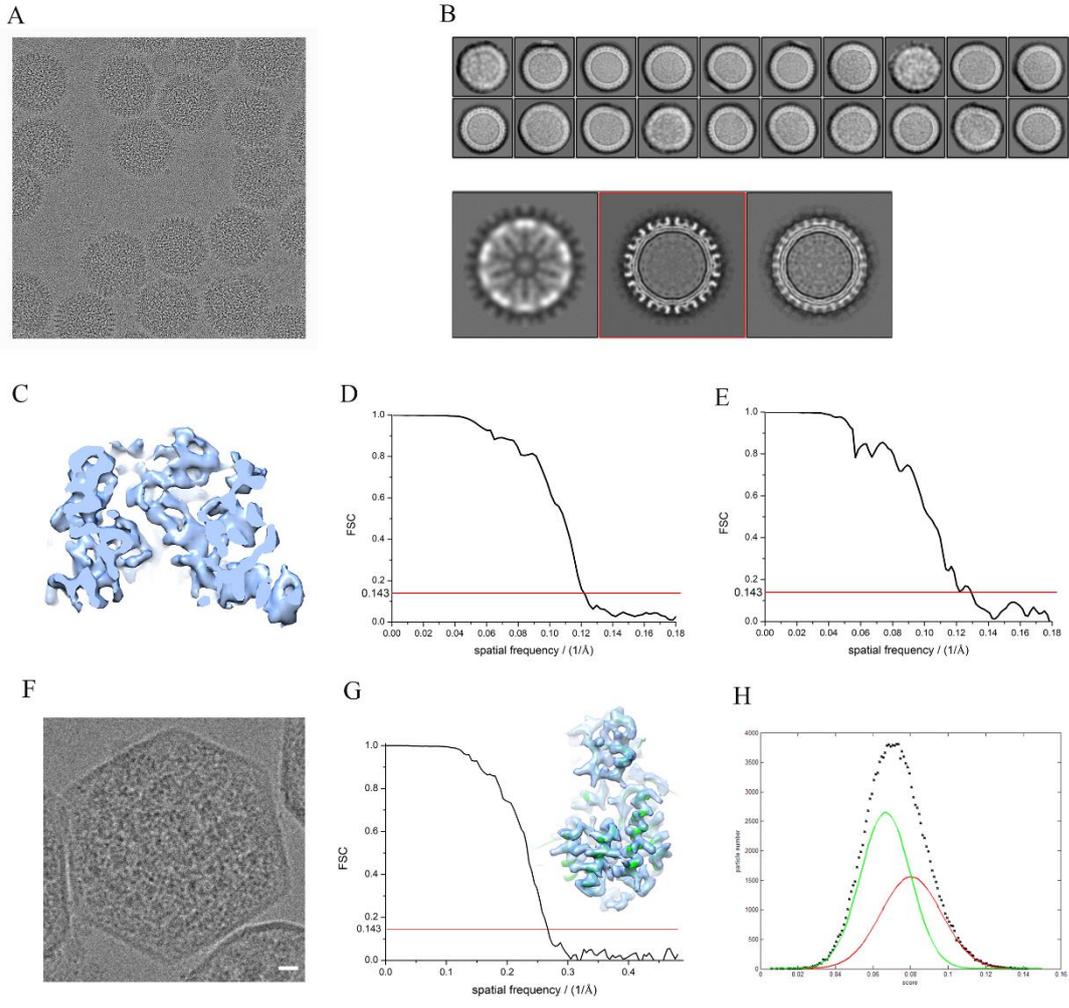
772

773

774

775

**Figure 5**



**Application.**

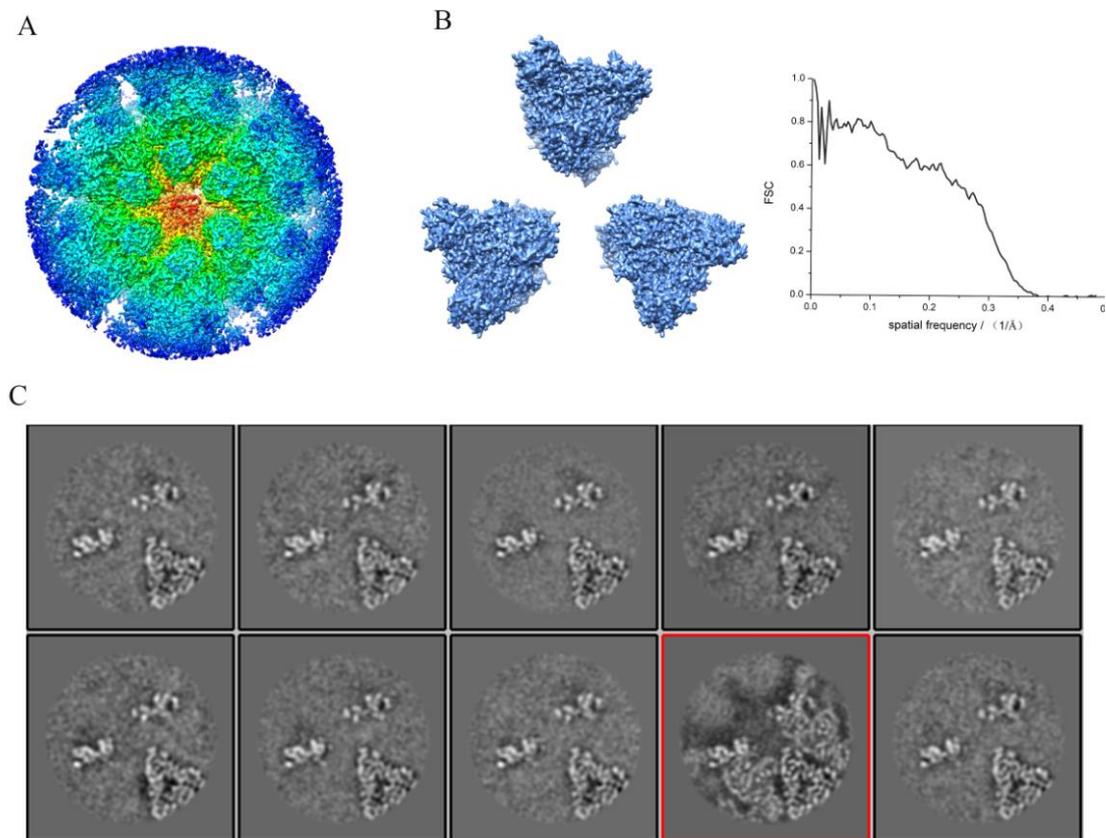
(A) A micrograph of Bunyavirus particles. (B) Upper panel: 2D classification of viral particles binned by 4. Lower panel: 3D classification of viral particles selected from 2D classification. (C) Densities in the 7.7 Å map of pentamer. (D) FSC curve shows the resolution of hexamer map. (E) FSC curve shows the resolution of pentamer map. (F) An image of a carboxysome, Rubiscos are packaged inside; scale bar represents 10 nm. (G) The 3.7 Å map of Rubisco reconstructed using our isSPA method (right) and the FSC curve showing the resolution of the 3.7 Å map at gold standard (left). (H) Distribution of particle numbers to scores fitted by two Gaussian functions.

820 **Supplementary Materials**

821

822

823 **Figure S1**



824

825 **Data processing of Reovirus.**

826 (A) The 4.0 Å reconstruction of reovirus hexamer by our isSPA methods with previous reconstruction as picking  
827 model. (B) Left: EM map produced by PDB of the crystal structure (1jmu). Right: FSC curve describing the similarity  
828 between the target and homologous model. (C) 3D classification of selected particles from global search with  
829 homologous map as picking model, and particles in class 8 (which is bolded in a red square) were most true  
830 detections.

831

832

833 **Figure S2**

834

835

836

837

838

839

840

841

842

843 **FSC weight and picking threshold on HSV-2 dataset.**

844 (A) Precision-recall curves of the ~2 MD protein rescaled to ~8 Å at FSC 0.5. We tested on HSV-2 with FSC

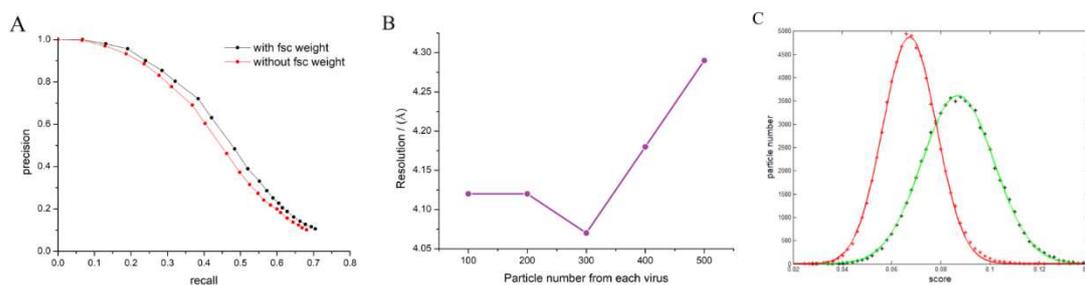
845 weighting (black) and without FSC weighting (red). (B) The relationship of resolution to threshold tested on HSV-2

846 dataset before PR sorting (black) and after PR sorting (purple). (C) Gaussian fitting of particle number to score

847 distributions by true solutions (red) and false solutions (green).

848

849



850 **Figure S3**

851

852

853

854

855

856

857

858

859

860

861

862

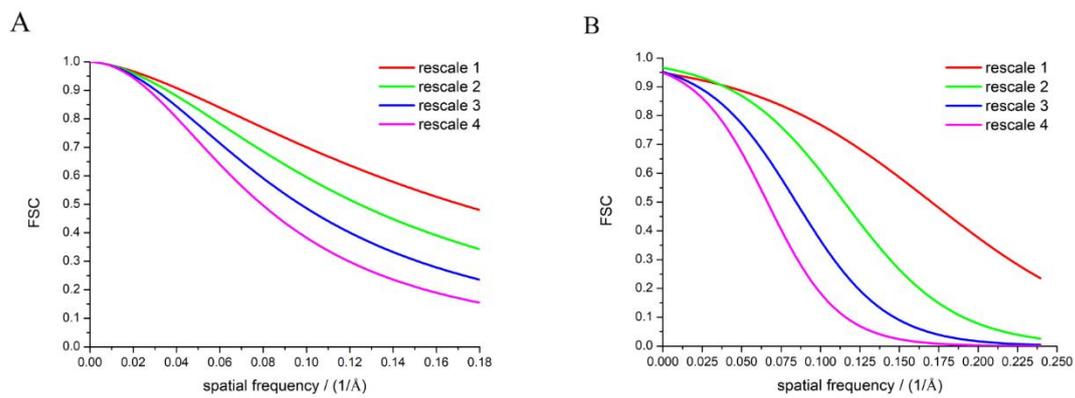
863 **FSC curves of rescaled maps with the origin maps.**

864 (A) HSV-2 hexamer rescaled at different scales. (B) Reovirus hexamer rescaled at different scales.

865

866

867



868 **Figure S4**

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

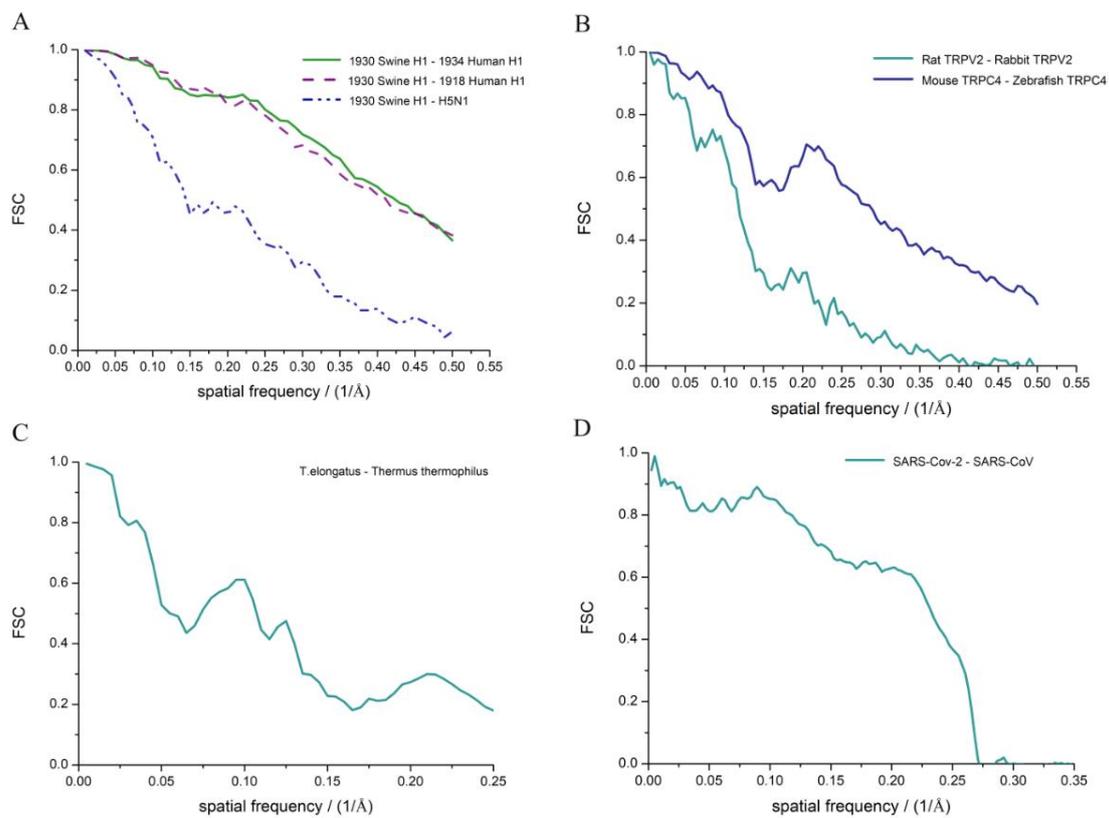
889

890

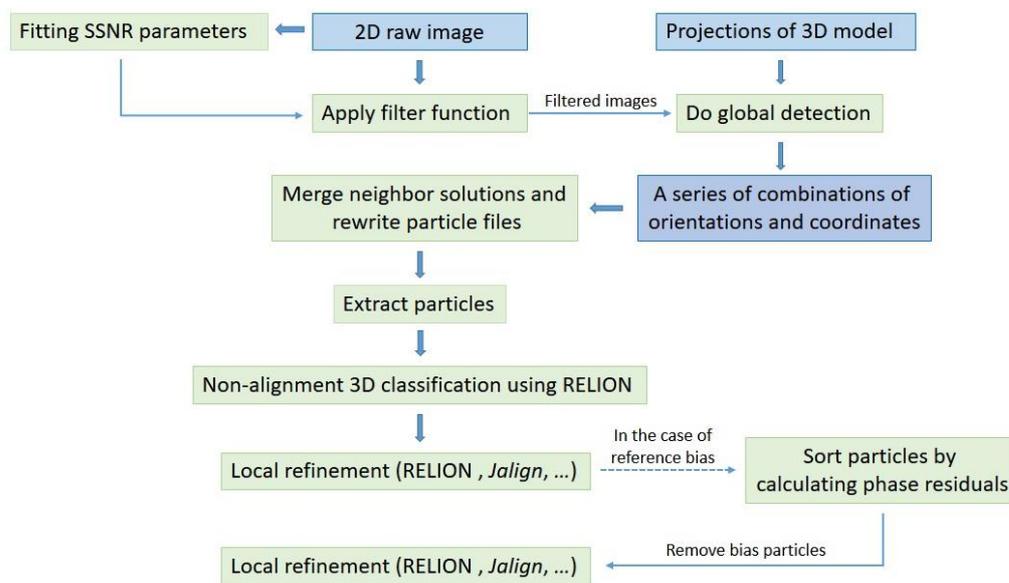
891 **FSC curves between homologous structures.**

892 (A) Glycoproteins of influenza viruses. (B) Ion channels proteins. (C) NDH. (D) Spike proteins of coronaviruses.

893



894 **Figure S5**



911 **Flowchart of isSPA method.**

938 **Figure S6**

939

940

941

942

943

944

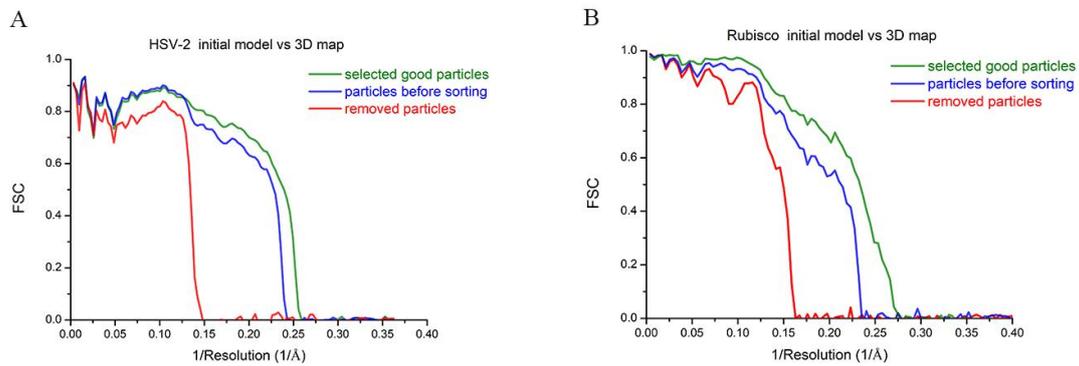
945

946

947

948

949



950 **FSC with high resolution initial model.**

951 (A) FSCs of 3D reconstructions by selected particles according to sorting (green), removed bad particles according

952 to sorting (red) and both particles before sorting (navy) with high resolution initial model. (B) the same processing

953 with Rubisco images as in (A).

# Figures

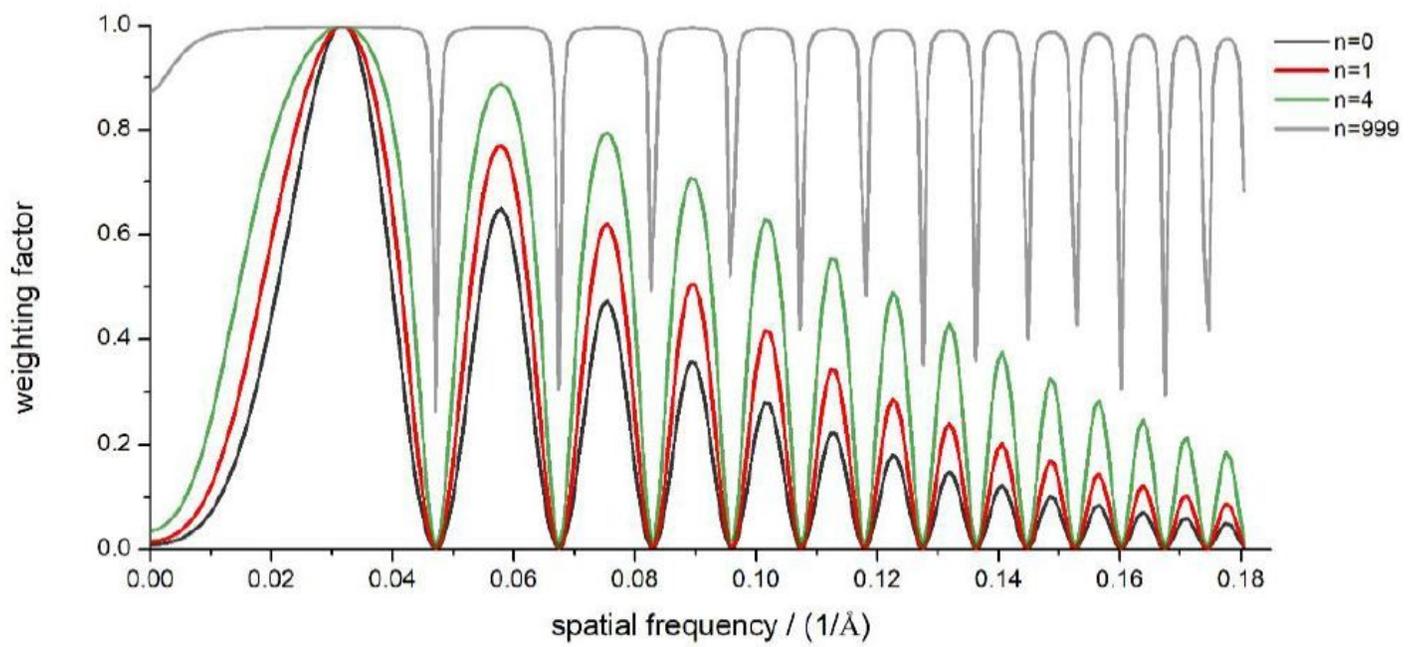
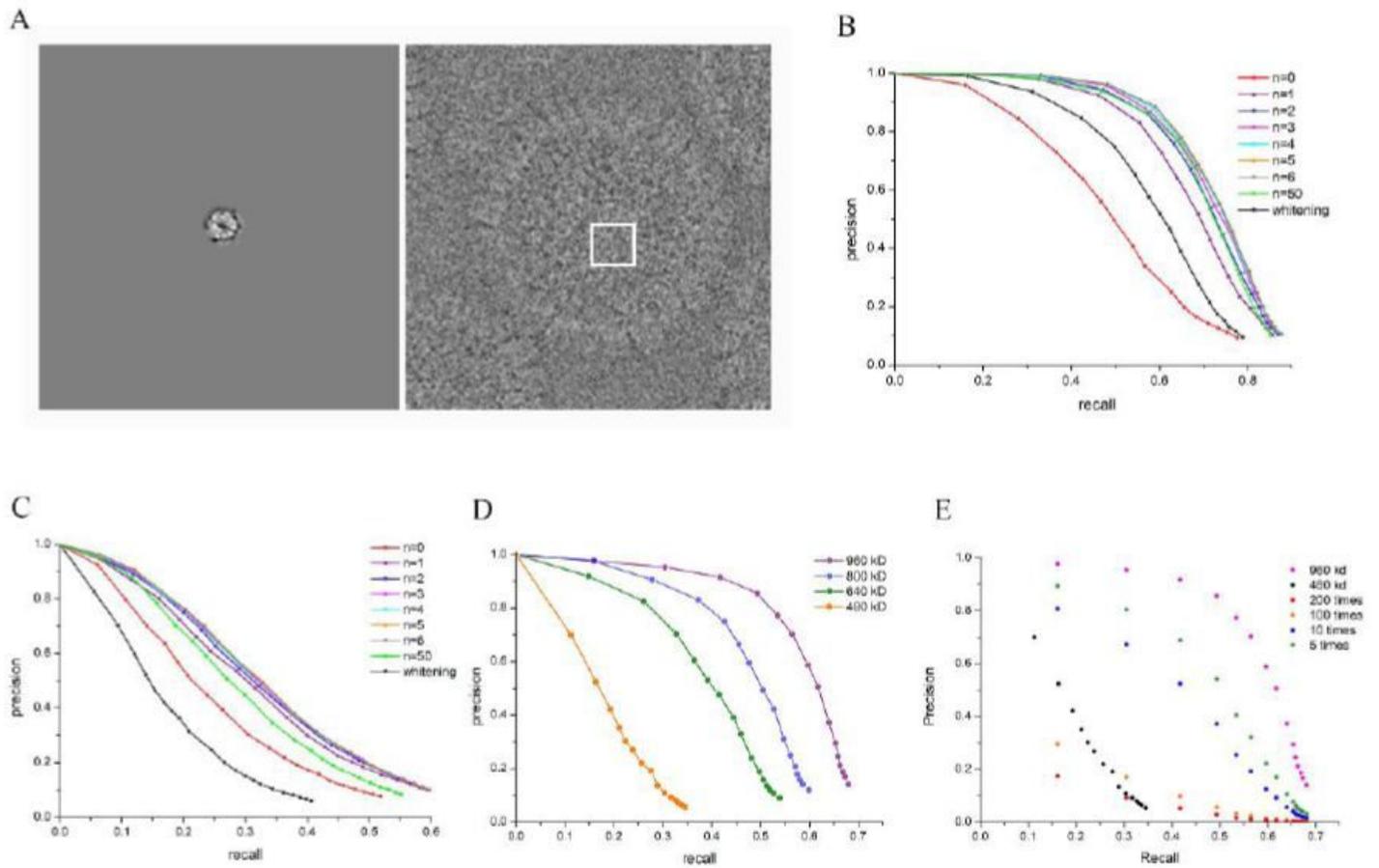


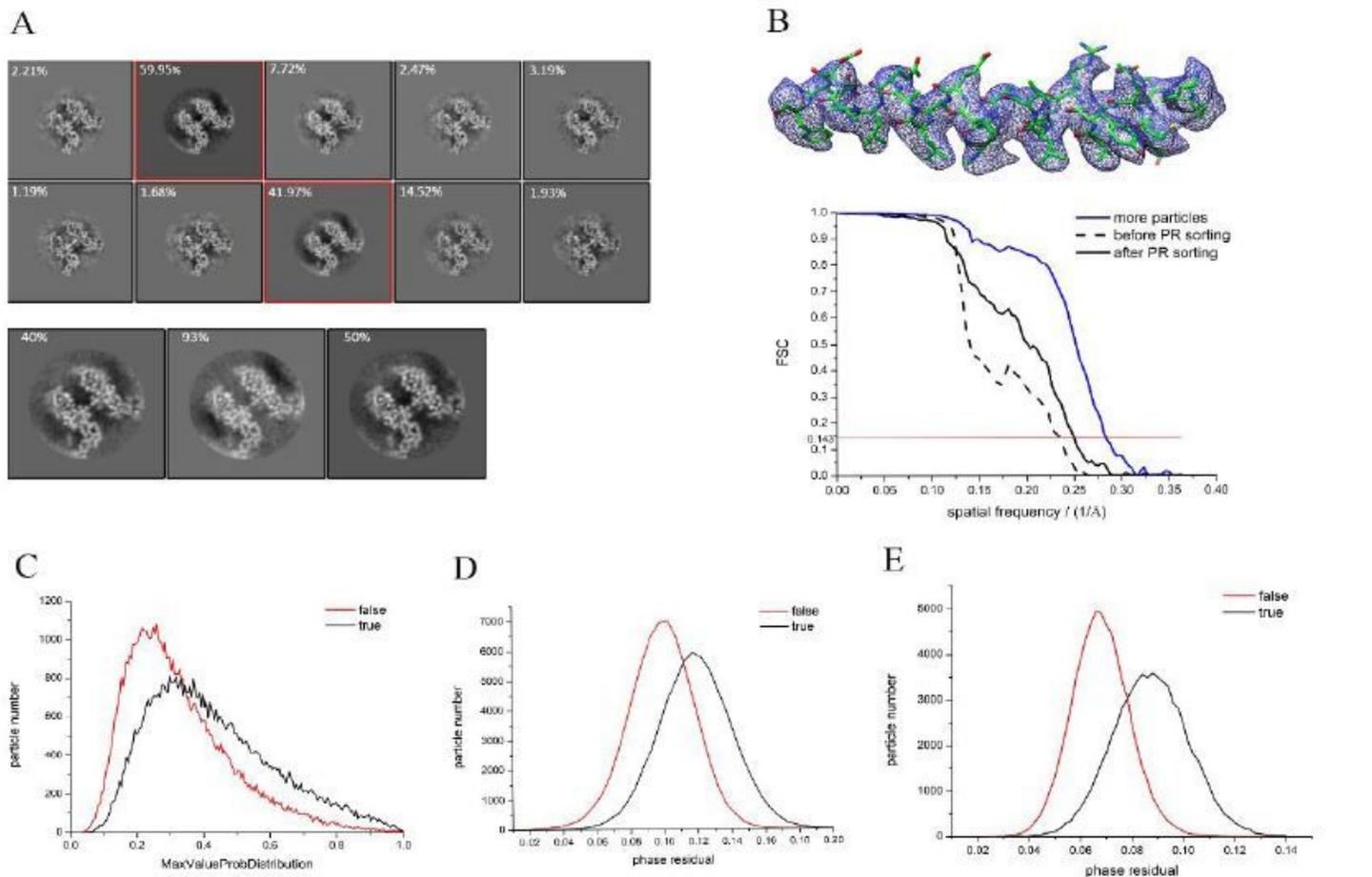
Figure 1

Weighting functions. Weighting factors damp and oscillate with spatial frequency (black, red, green and gray) at different  $n$  (0, 1, 4, 999).



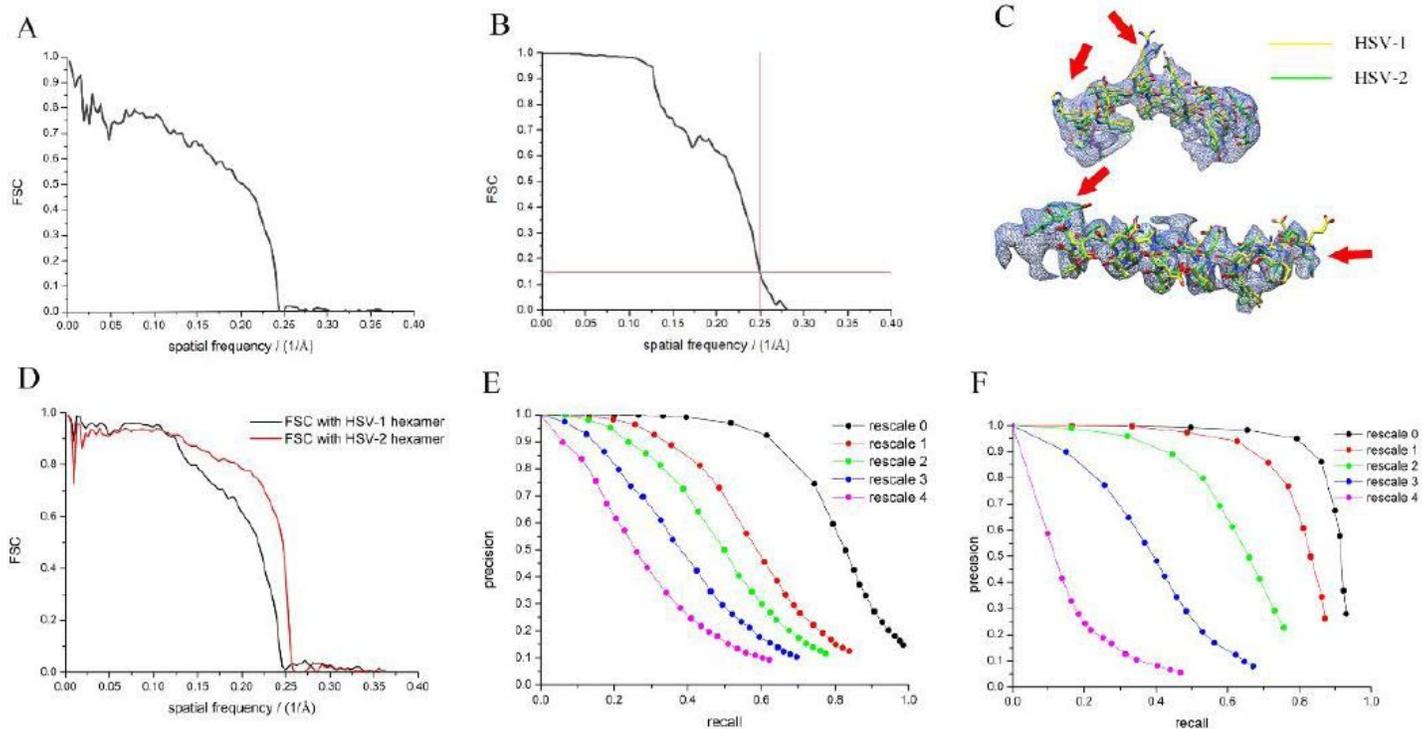
**Figure 2**

Efficiency of particle detection. (A) Left: projection of HSV-2 hexamer model. Right: HSV-2 virus particle imaged at 2.75 pm defocus with 25 electrons per k . The location of the projection on the virus particles is indicated by a white square. (B) Precision-recall curves for detections on Reovirus datasets using model in 900 KD molecular weight at  $n=0$  (red), 1 (purple), 2 (navy), 3 (magenta), 4 (cyan), 5 (orange), 6 (gray), 50 (green) and whitening filter (black). (C) Precision-recall curves for detections on HSV-2 datasets with the same processing as in (B). (D) Precision-recall curves (purple, navy, green and orange) for detections ( $l$  equals 3) on Alphaviruses using models in different molecular weights (960 kD, 800 kD, 640 kD and 480 kD). (E) Precision-recall curves for lower abundance of a 960 kD protein on Alphavirus (olive for 5 times, navy for 10 times, orange for 100 times and red for 200 times lower abundance), 960 kD at abundance of 60 copies per  $90 \text{ nm} \cdot 90 \text{ nm}$  square (magenta) and 480 kd at abundance of 60 copies per  $900 \text{ nm} \cdot 900 \text{ nm}$  image (black).



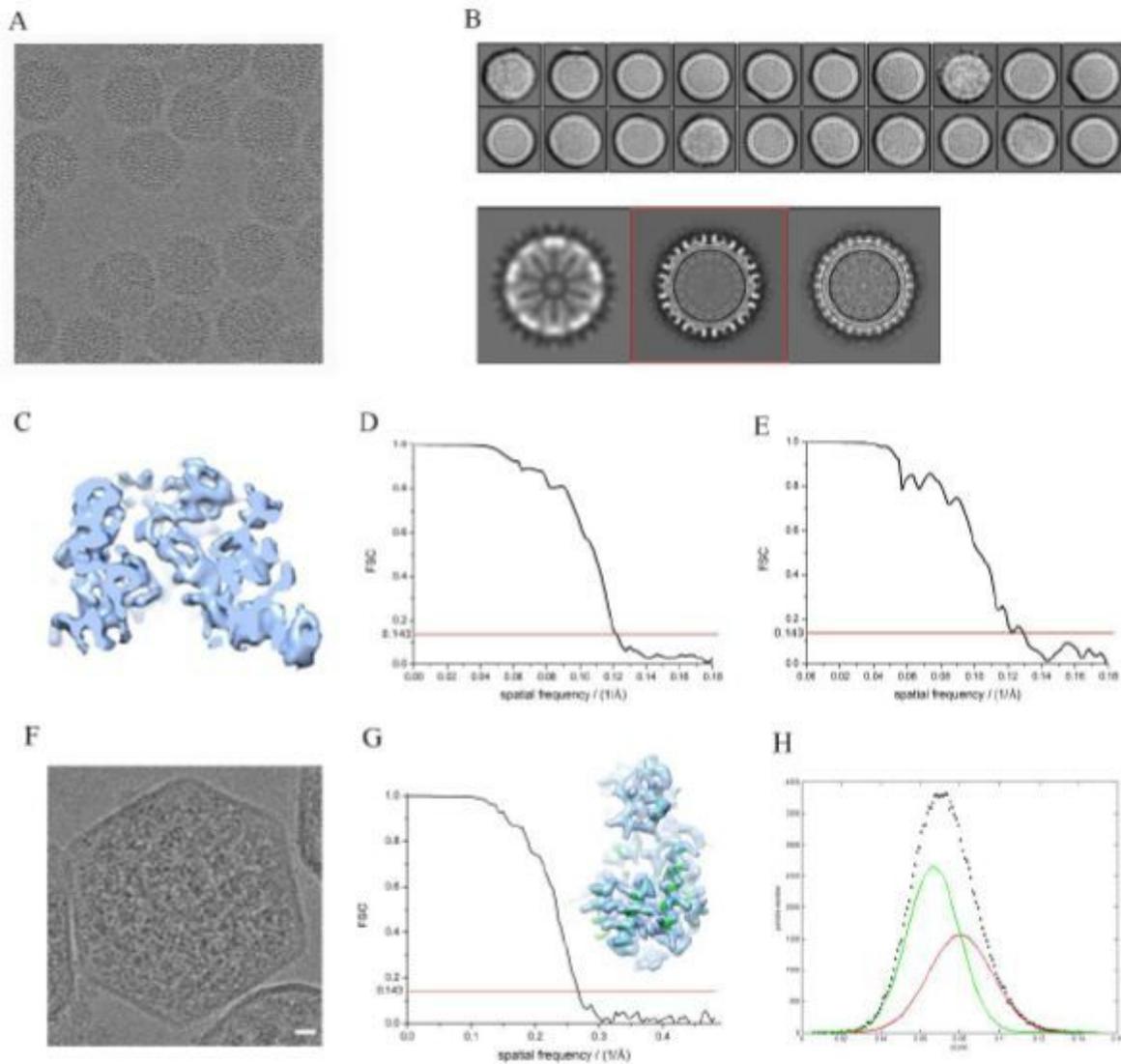
**Figure 3**

Data processing of HSV-2 hexamer. (A) Upper panel: 3D classification of raw picked positives in ten classes, the percentage of true detections in each class is shown at the top left corner. The two selected classes are indicated by squares in red. Lower panel: 3D classification of selected particles in three classes, and the percentage of true detections in each class is noted at the top left corner. (B) Upper panel: 3D reconstruction of the HSV-2 hexamer. Lower panel: FSC curves denote 3 reconstructions in different conditions using 2000 viral particles before PR sorting (dash) and after PR sorting (black), using 8000 virus particles and after PR sorting (navy). (C) True and false particle distribution with MosValueProbDistribution term in RELION. (D) True and false particle distribution with phase residual using data from 1/30 Å to 1/8 Å. (E) True and false particle distribution with phase residual using data from 1/8 Å to 115 M.



**Figure 4**

Homologous structure as picldug model. (A) FSC curve of H5V-1 hemmer with HSV-2 hemmer. (B) The resolutions determined by gold standard FSC at threshold 0.143 of H5V-2 hemmer reconstructed from H5V-1 hemmer. (C) PDB of H5V-1 hemmer (yellow) and HSV-2 (green) fitting to the density of the 4.0 Å cryo-Em map, the disagreements of the two PDBs are pointed out by red arrows. (D) FSC curves show similarities between our map and H5V-1 hemmer (black), and HSV-2 hexamer (red) respectively. (E) Precision-recall curves of the  $-2$  MD protein of HSV-2 at a series of scales corresponding to Fig. S3. (F) Precision-recall curves of the  $-LS$  MD protein of Reovirus at a series of scales corresponding to Fig. S3.



**Figure 5**

Application. (A) A micrograph of Bunyavirus particles. (B) Upper panel: 2D classification of viral particles binned by 4. Lower panel: 3D classification of viral particles selected from 2D classification. (C) Densities in the 7.7 Å map of pentamer. (D) FSC curve shows the resolution of hexamer map. (E) FSC curve shows the resolution of pentamer map. (F) An image of a carboxysome, Rubiscos are packaged inside, scale bar represents 100 nm. (G) The 3.7 Å map of Rubisco reconstructed using our isSPA method (right) and the FSC curve showing the resolution of the 3.7 Å map at gold standard (left). (H) Distribution of particle numbers to scores fitted by two Gaussian functions.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [nrreportingsummary44.pdf](#)