

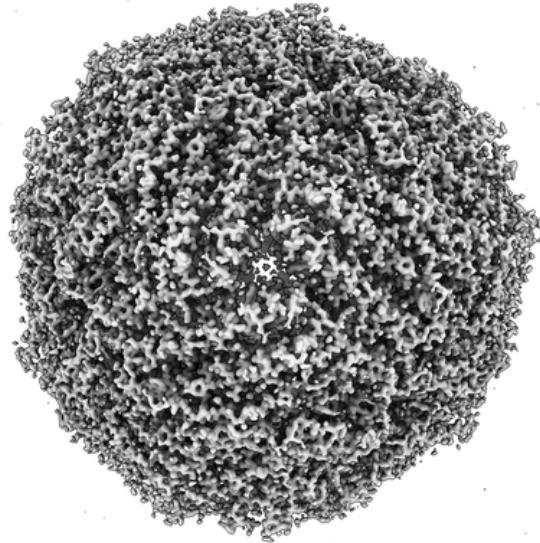


Scipion Tutorial Series

NATIONAL CENTER FOR BIOTECHNOLOGY
BIOCOMPUTING UNIT

Image Processing

November 18, 2019



Cryo-EM map of apoferritin at 1.54 Å resolution (EMD-9865)

CARLOS OSCAR S. SORZANO & MARTA MARTÍNEZ

Revision History

Revision	Date	Author(s)	Description
1.0	11.18.2019	MM	Initial draft created for the S2C2 CryoEM Image Processing Workshop held in Stanford

Intended audience

The recent rapid development of single-particle electron cryo-microscopy (cryo-EM) allows structures to be solved by this method at almost atomic resolutions. Providing a basic introduction to image processing, this tutorial shows the basic workflow aimed at obtaining high-quality density maps from cryo-EM data by using *Scipion* software framework.

We'd like to hear from you

We have tested and verified the different steps described in this demo to the best of our knowledge, but since our programs are in continuous development you may find inaccuracies and errors in this text. Please let us know about any errors, as well as your suggestions for future editions, by writing to scipion@cnb.csic.es.

Requirements

This tutorial requires, in addition to *Scipion*, *cryoSPARC2* (<https://cryosparc.com/>).

Contents

1	Introduction to image processing	4
2	Problem to solve: Apoferritin	6
3	From movies to micrographs	9
4	CTF estimation	14
5	Particle picking	20
6	Extract Particles	28
7	2D classification	31
8	Initial volume	35
9	3D classification and Refinement	47

1 Introduction to image processing

Definition

Image processing is a structure determination technique that allows to get the 3D density map from a set of cryo-EM images of a particular macromolecule. Although different structural approaches can be followed to analyze the structures of macromolecules, this tutorial focuses on cryo-EM single particle analysis (SPA). Fortunately, cryo-EM SPA is undergoing in this decade a resolution revolution that has allowed the structures of macromolecules to be solved at near-atomic resolution.

Image processing workflow

The set of successive tasks aimed to get the 3D density map is known as image processing workflow. Main steps of the general workflow are detailed from top (movies) to bottom (refined volume) in the Fig. 1. Some of the tools required to perform the respective tasks are detailed in a non-exhaustive way on the right side of the Figure. All these methods have been integrated in *Scipion* to facilitate interoperability among different software packages, data tracking and reproducibility of results.

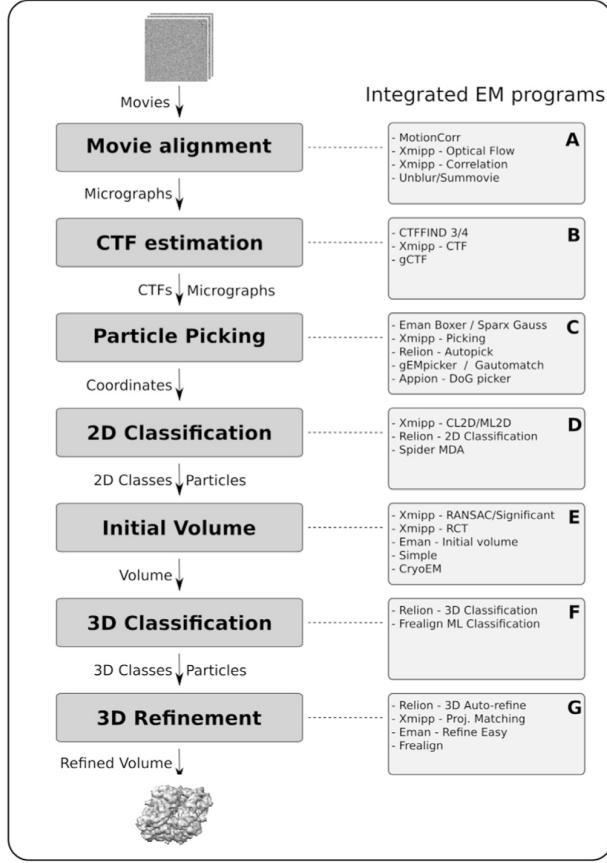


Figure 1: General image processing workflow for SPA (de la Rosa-Trevin et al., 2016).

The workflow considers as input the movie frames generated by the microscope. These movies should be global or locally aligned before computing the CTF of individual micrographs. *Scipion* allows to compare different CTF values obtained with distinct algorithms using the CTF consensus protocol. Once the CTF has been corrected, we are ready to extract individual particles of each micrograph by using different protocols of particle picking. As in the case of the CTF, we can retrieve the coordinates of each particle using different protocols of manual and automatic picking, and finally, estimate the agreement between all those methods through a consensus picking protocol. The screened particles are used for further processing. The next step involves the 2D classification of the individual selected particles. 2D

classes derived from the last procedure contribute to generate the initial 3D map. The last part of the workflow includes 3D classification and 3D refinement tasks in order to iterative refine the initial 3D map.

In this tutorial, we show all above mentioned processes of 3DEM processing, as well as the necessary tools to accomplish them, illustrating the combination of different EM software packages in *Scipion*.

2 Problem to solve: Apoferritin

Ferritins are iron storage metalloproteins ubiquitously distributed among living organisms. These proteins are involved in iron metabolism in many different types of cells, and play a relevant dual role both in iron detoxification and iron reserve. The ferritin's architecture, similar to a spherical shell, is highly conserved in bacteria, plants and animals, and it allows to accumulate high amounts of Fe(III) atoms (up to 4000 per molecule).

The highly stable iron-free shell is known as apoferritin. Mammalian apoferritins are heteromeric molecules, constituted by 24 monomers structurally equivalent that surround the central cavity. Among these monomers, variable proportions of two types of subunits with different properties, H (heavy) and L (light), can be found. The tissues involved in iron storage contain higher proportion of L chains, whereas the tissues that require higher protection against oxidation, such as heart or brain, have a higher content of H chains. Unlike L chains, H chains display ferroxidase catalytic activity, necessary to oxidize Fe(II) to Fe(III). Concerning the structure of each subunit, it is constituted by 4 long helices, a fifth smaller helix and an additional extended loop. The dinuclear iron site, or ferroxidase site, is located in the center of the four helix bundle.

This tutorial will guide us in the building process of the mouse apoferritin 3D map using the *Scipion* framework (Fig. 2). As starting input data, we are going to use the EMPIAR ID: 10248 data, obtained from mouse heavy chain apoferritin. This cryo-EM data allowed to generate the 3D map EMD-9865 at 1.54 Å resolution (Hamaguchi

et al., 2019). The most recent atomic structure of mouse apoferritin, homo 24-mer of ferritin heavy chain with octahedral symmetry, was also obtained from cryoEM data at 1.84 Å (PDB ID: 6S61). The 24 monomers of this metal binding protein are ligated to 6 Fe(III) and 24 Zn(II) ions.

Apo ferritin processing workflow in *Scipion*

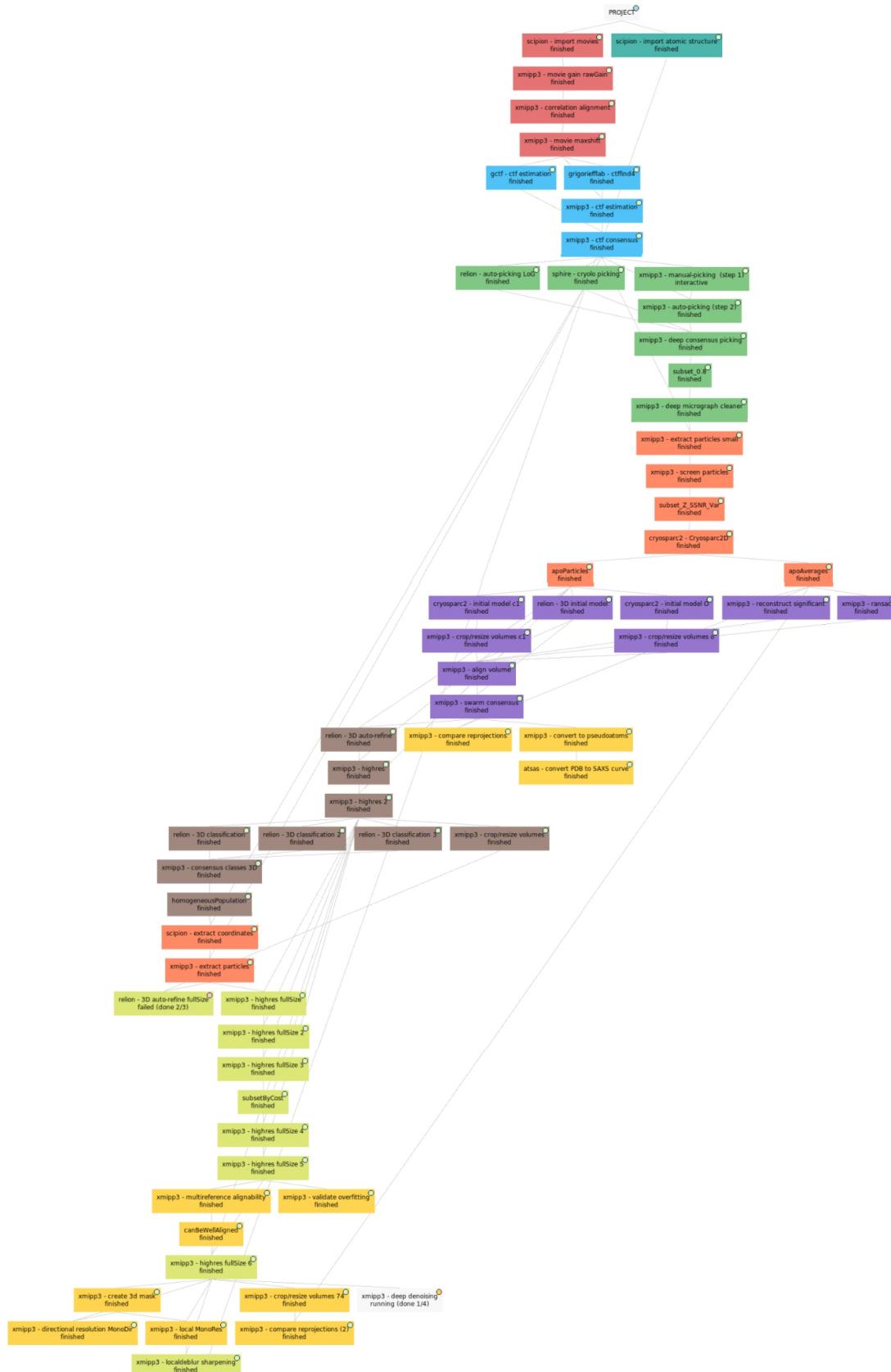


Figure 2: Apoferritin processing workflow.

3 From movies to micrographs

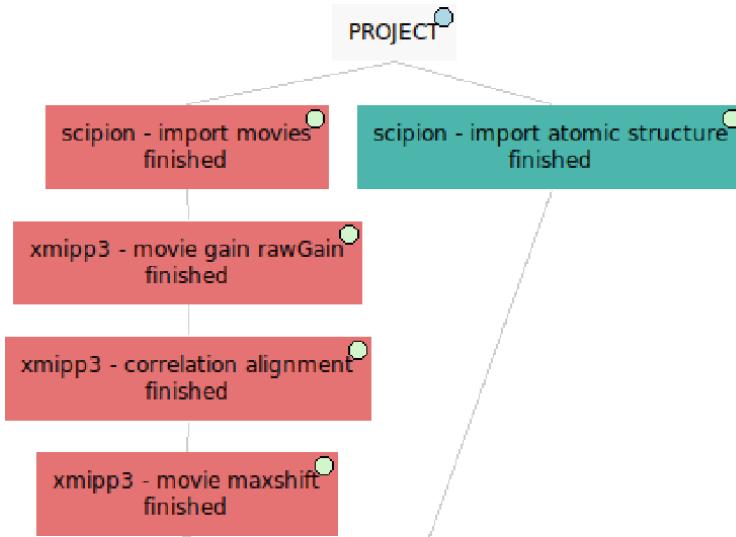


Figure 3: From movies to micrographs workflow.

Import movies

The protocol `scipion-import movies` allows to download the mouse apoferritin cryo-EM data in *Scipion*. The protocol form with parameters can be seen in Fig. 4. With this protocol, besides the set of movies, located in the Data folder, acquisition parameters such as accelerating voltage, spherical aberration and sample rate, will be registered in your *Scipion* project.

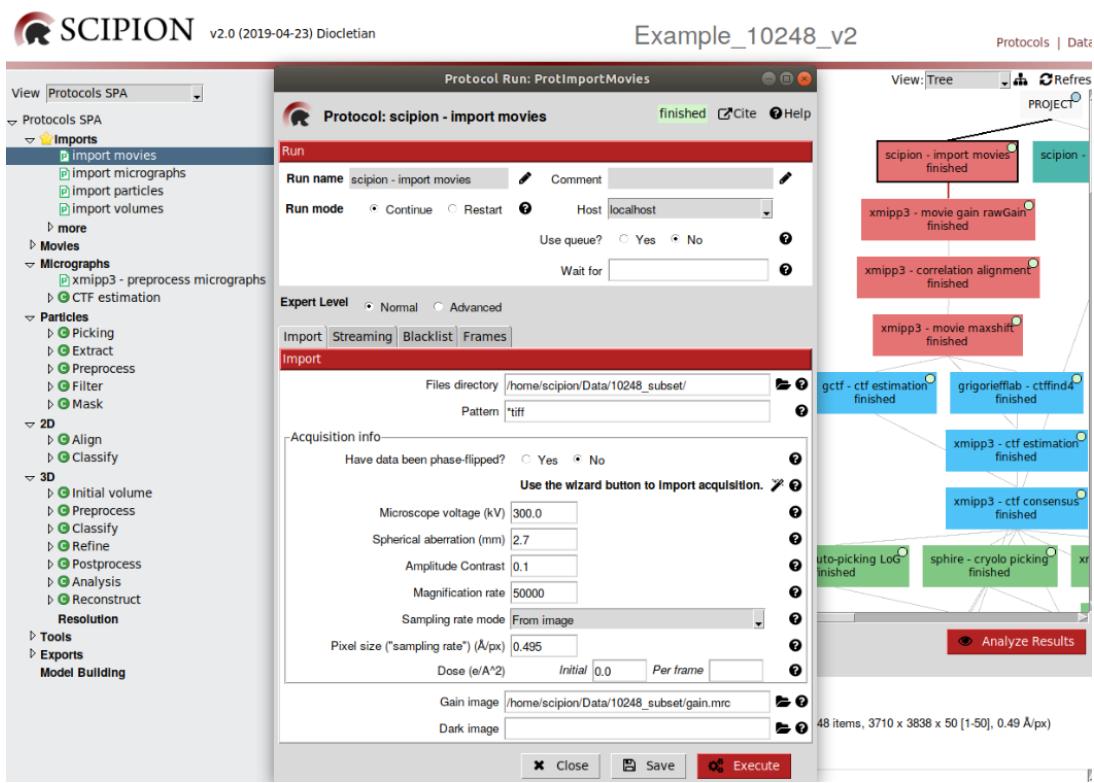


Figure 4: Filling in the protocol to import cryo-EM data.

After executing this protocol, we can visualize the list of 48 movies imported to the project by pressing **Analyze Results**. Each movie contains 50 frames (size 3710 pixels x 3838 pixels). Frames contained in each movie can be visualized by right-clicking each entry.

Computation of movie gain

The protocol `xmipp3-movie gain` is used to compute the movie gain (Fig. 5). Two movie gains will be computed: 1) Without applying the input gain, to orientate the input movie gain; 2) Applying the input gain to estimate the residual movie gain.

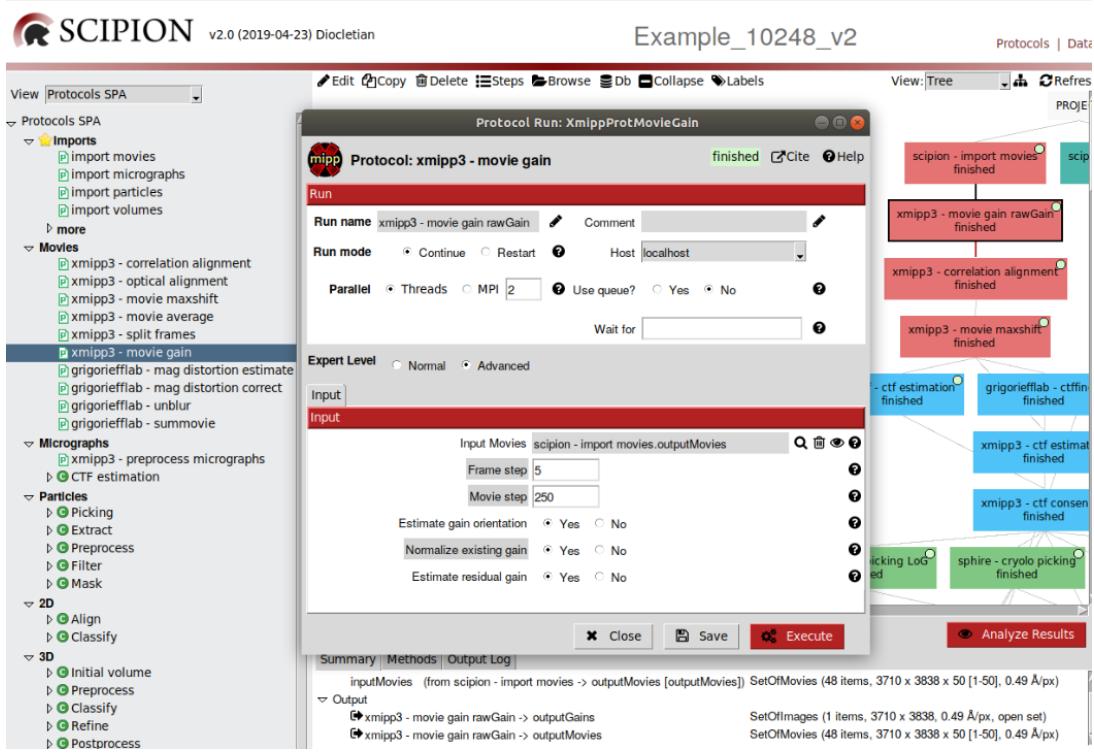


Figure 5: Completing the protocol to compute the movie gain.

After executing this protocol, by pressing **Analyze Results** we can visualize the image of the gain computed. None of the movie gains computed moves forward the protocol output.

Movie alignment

In order to correct BIM-induced image blurring and restore important high resolution information, the stack of individual frames contained in each movie needs to be aligned. Only one image will be generated, and this image is called micrograph. Although in *Scipion* we have integrated several protocols to perform global and local alignment, in this tutorial we are going to use **xmipp3-correlation alignment**. We have completed the params of this protocol in Fig. 6.

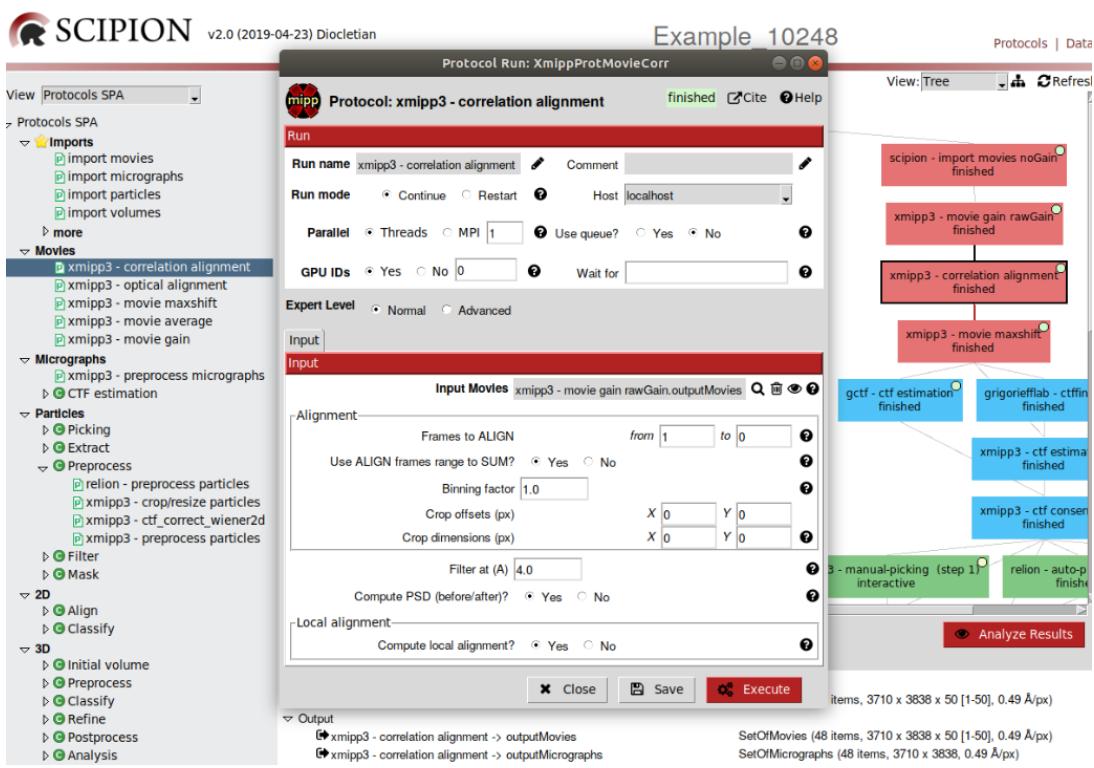


Figure 6: Filling in the protocol to align the frames of each movie.

When the execution of this protocol finishes, we can observe the list of the set of 48 resulting micrographs by pressing **Analyze Results**. The first column contains composite images with half of the PSD of the unaligned micrograph (left side) and half of the PSD of the aligned one (right side). The plots in the second column reflect accumulated shifts of the movie frames. The name of each resulting micrograph appears in the third column. Each micrograph included in the set generated can be opened for visual inspection by right-clicking its entry.

Screening of micrographs

Since some of the micrographs generated in the previous step could derive from movies with high drift among frames, we have added in the processing workflow a step to select only the micrographs originated from movies with allowed drifting

values among frames. The protocol [xmipp3-movie maxshift], completed in Fig. 7, was designed to screen micrographs according maximum shift values. Movies will be rejected if they exceed the maximum shift value among frames, the maximum travel value for the whole movie, or both previous conditions.

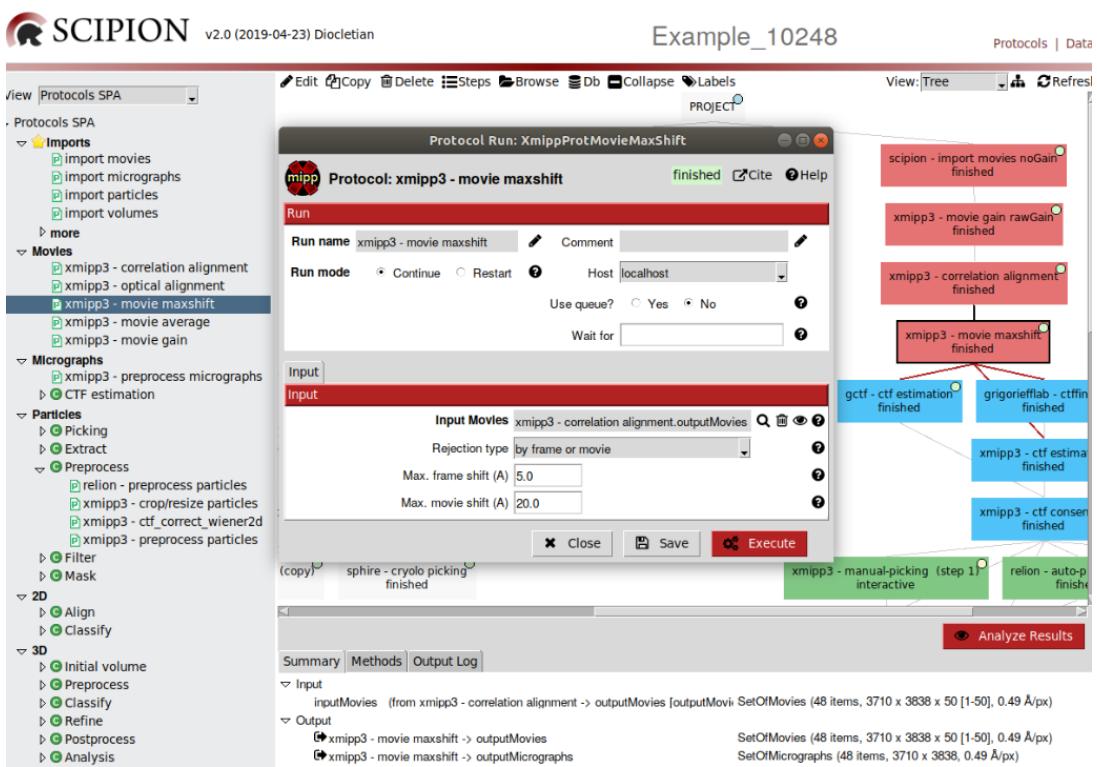


Figure 7: Completing the protocol to screen the micrographs.

Once the protocol is executed, both discarded and accepted lists of micrographs can be visualized by pressing [Analyze Results]. Each micrograph can be opened for visual inspection by right-clicking. In this case, the set of 48 input micrographs has been included in the protocol output. This set will serve as input for further processing steps.

4 CTF estimation

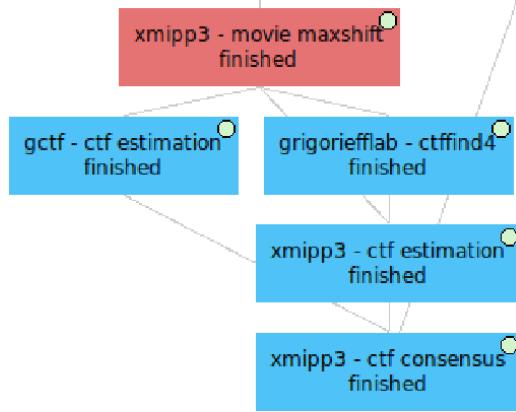


Figure 8: CTF estimation workflow (Blue color).

Since close to focus images of biological specimens embedded in vitreous ice generate very little contrast, we take our movie-frames out of focus, and retrieve systematically distorted images of our specimens. The alteration observed is due to the different transference of contrast for each frequency. In an ideal microscope all the frequencies are transferred with total contrast (+1). In a normal one, some frequencies are transferred with contrast 0 or even -1. The CTF (Contrast transfer function) indicates how much contrast is transferred to the image as a function of the spatial frequency. The estimation of the CTF is the first step to correct it, repair its negative effect and retrieve our specimens undistorted.

How to estimate the CTF?

Since part of the Fourier components are lost, attenuated or inverted, images of the specimens taken in a non-ideal microscope will appear blurry. We define this blurry effect with the PSF (Point spread function). The effect of the PSF makes that a discrete point in the specimen is reproduced in the image as a broad point with a complex shape. The PSF can be directly estimated from the micrographs and, since the PSF and the CTF are related through the Fourier Transform, by computing the

Fourier Transform of the PSF, we can directly estimate the CTF.

We count on different protocols to estimate the CTF of the micrographs in *Scipion*. In this tutorial we are going to use three different algorithms: **Gctf** (Zhang, 2016), **CTFFind4** (Rohou and Grigorieff, 2015) and **Xmipp CTF estimation** (Sorzano et al., 2013) executed with protocols **[gctf-ctf estimation]** (Fig. 9), **[grigoriefflab-ctffind4]** (Fig. 10) and **[xmipp3-ctf estimation]** (Fig. 11), respectively. Besides of estimating the CTF envelope, this last protocol improves the CTF estimation of CTFFind4. Thus, in this case an additional parameter is the defoci from a **Previous CTF estimation**.

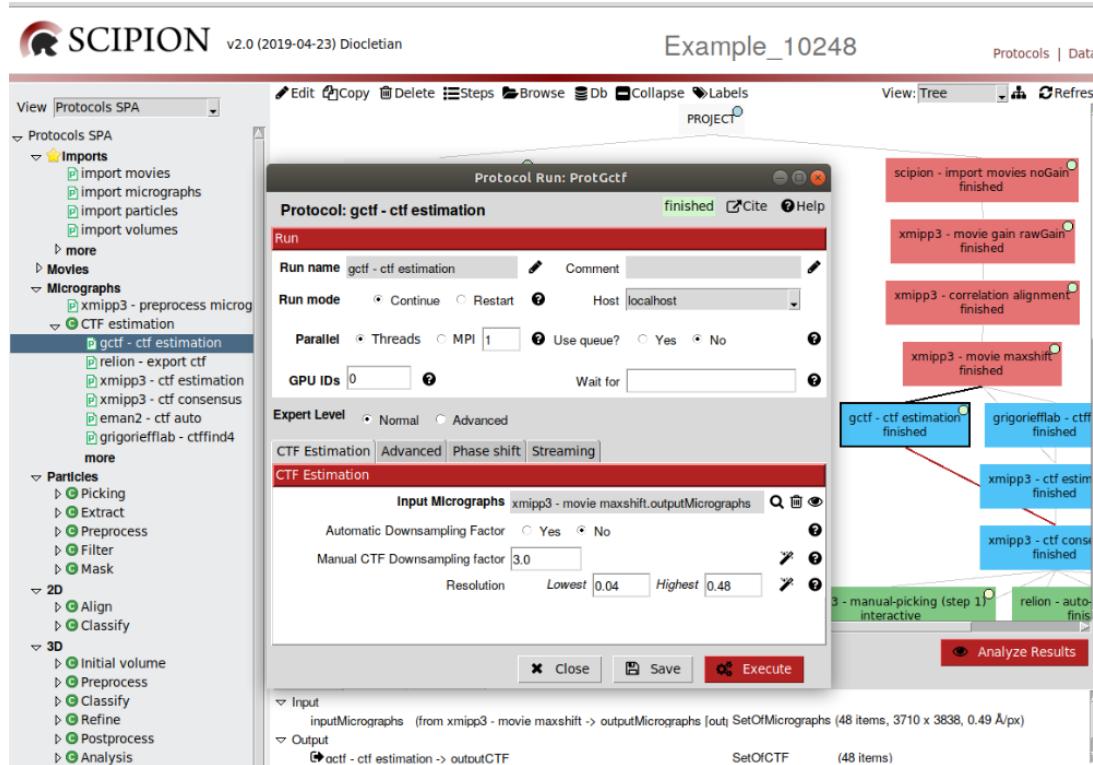


Figure 9: Protocol **[gctf-ctf estimation]** to compute the CTF.

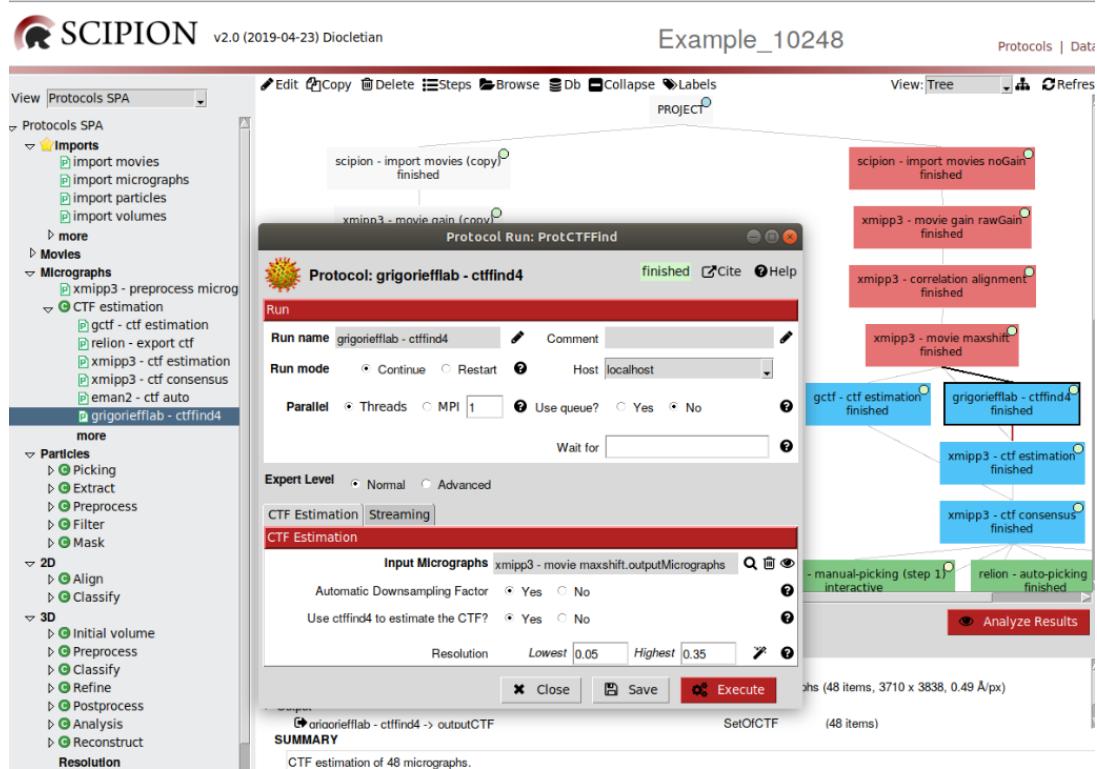


Figure 10: Protocol `grigoriefflab-ctffind4` to compute the CTF.

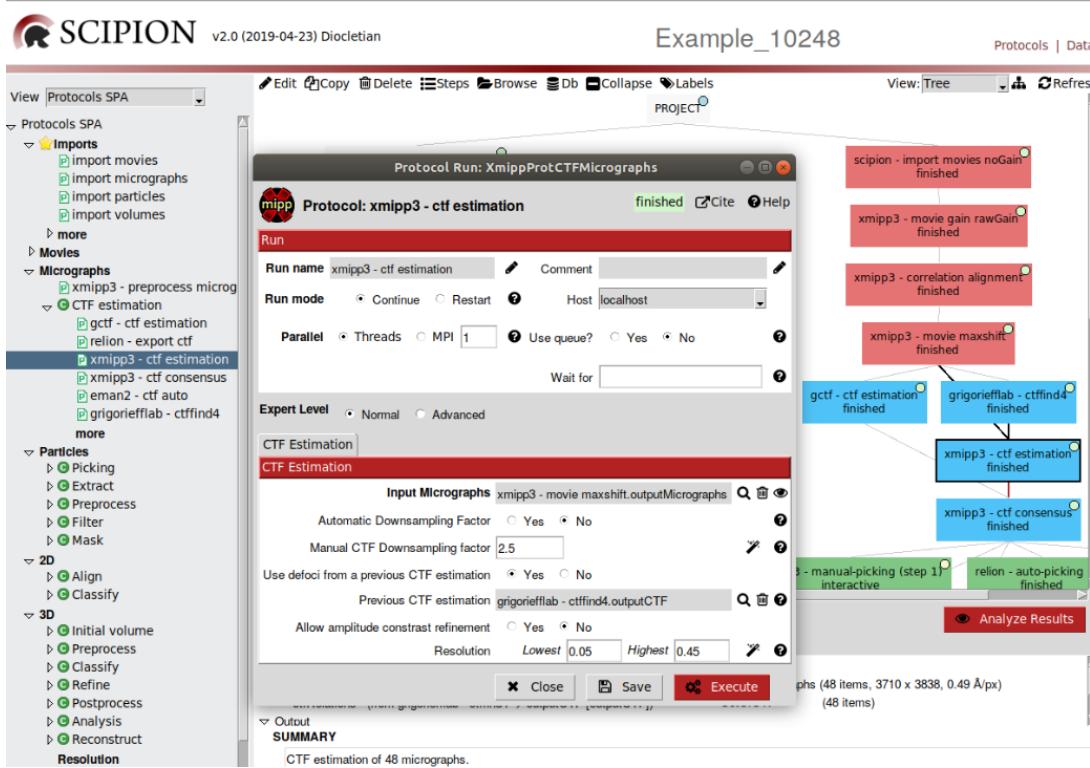


Figure 11: Protocol `xmipp3-ctf estimation` to compute the CTF.

The three algorithms estimate the PSD of the micrographs and the parameters of the CTF (defocus U, defocus V, defocus angle, etc.). They cut micrographs into many smaller images with a desired window size. After that, they compute the Fourier Transform of each image and calculate an average. The three protocols designed to apply the three respective algorithms contain very similar parameters. To estimate the CTF we need to limit the frequency region to be analyzed between the lowest and highest resolution. The frequency domain selected must include all zeros of the CTF. The wizard displayed on the right helps to choose that frequency region.

After executing each one of these three protocols, results can be observed by pressing `Analyze results`. A table will be opened showing the image of the CTF computed for each micrograph, as well as other CTF parameters. CTFs of good micrographs typically display multiple concentric rings extending from the image center towards

its edges. Bad micrographs, however, might lack rings or show very few of them that hardly extend from the image center. Micrographs like these will be discarded, likewise micrographs showing strongly asymmetric rings (astigmatic) or rings that attenuate in a particular direction (drifted). To discard a particular micrograph, select it, click the mouse right button and choose **Disable**. If you want to see the CTF profile, choose the option **Show CTF profile**, and a new window will be opened to show the CTF profile. **Recalculate CTFs** and **Micrographs** are additional options of **Analyze results** menu that can be used after selecting specific micrographs. **Recalculate CTFs** allows to estimate again the CTF when the algorithm has previously failed to find the rings, even if they can be seen by eye. The option **Micrographs** allows to create a new subset of selected micrographs.

Concerning some differences among protocols, we remark that micrographs with CTF estimated with **grigoriefflab-ctffind4** display a hole in the center because, in some cases, they have very much power and avoid appreciate what is underneath. In the particular case of **xmipp3-ctf estimation**, four different images of micrographs are shown. Besides the PSD, this last protocol displays the enhanced PSD, the CTF model by quadrants and half planes.

CTF consensus

The CTF estimation process concludes by applying a consensus protocol to assess differences among the output of the algorithms **Gctf** and **Xmipp CTF estimation**. We are going to use **xmipp3-ctf consensus** to perform this task (Fig. 12). This protocol allows to screen micrographs according meaningful CTF estimations based on defocus values, astigmatism, resolution and other **Xmipp** criteria (second tap in the protocol form), which will only be used in case that any of the CTFs computed is estimated by **xmipp3-ctf estimation**.

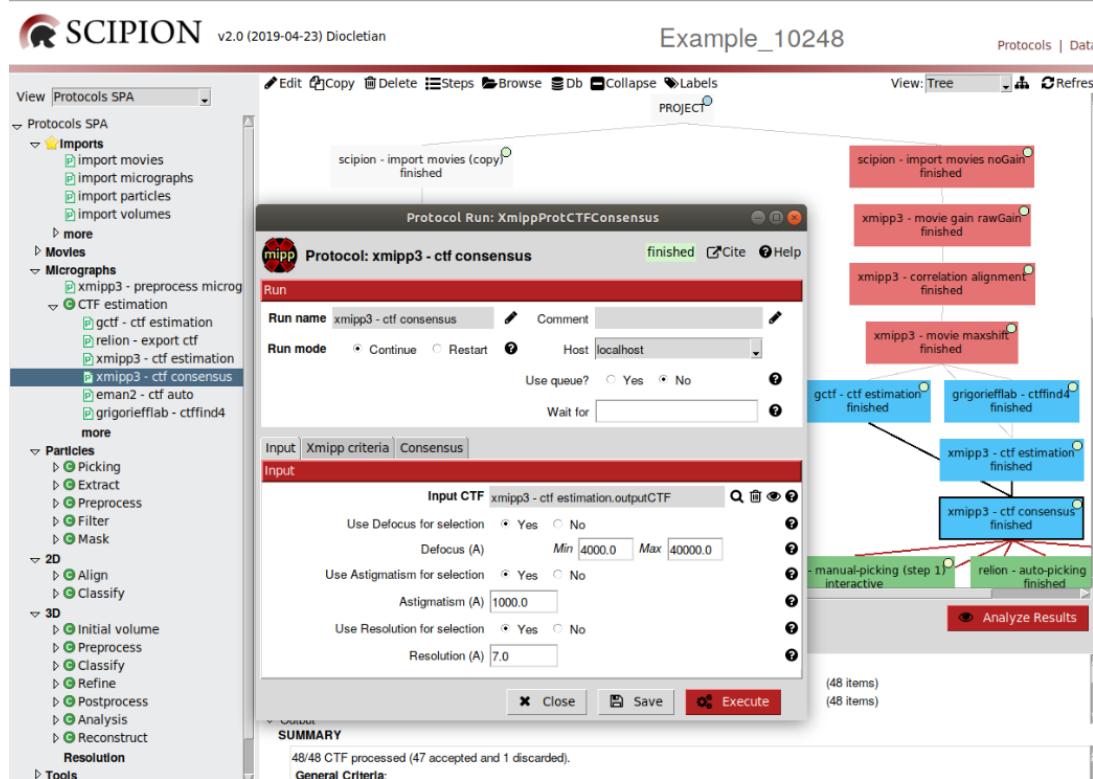


Figure 12: Filling in the consensus protocol `xmipp3-ctf consensus`.

In this case we have selected, as first input, the estimation of the CTF calculated by `xmipp3-ctf estimation` and, as second input, the estimation obtained by `gctf-ctf estimation`. By pressing `Analyze Results` both accepted (41) and rejected (1) micrographs can be visualized.

With the good micrographs we can continue the image processing. The negative effects of the CTF will be corrected in the next steps.

5 Particle picking

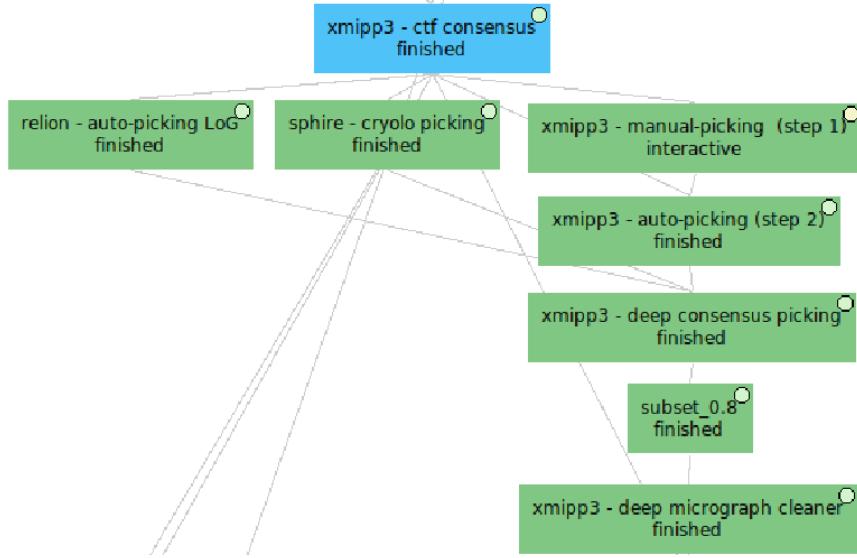


Figure 13: Particle picking workflow (Green color).

Since the reconstruction of the 3D density map of the macromolecule is based on images of individual single particles, the next step is essential in the image processing workflow. With the particle picking step we will retrieve the coordinates of each single particle.

Because manual picking can be very tedious, some tools have been designed to help in this task. We have integrated in *Scipion* different picking tools that can be used individually or in combination to obtain the final coordinates. Currently, we have tools available from **Eman2**, **Relion**, **Bsoft**, **Sphire** and **Xmipp**. In this tutorial we are going to use three different protocols that integrate tools from **Relion** (**relion- auto-picking LoG** (Zivanov et al., 2018)), **Sphire** (**sphire-cryolo picking** (Wagner et al., 2019)) and **Xmipp** (**xmipp3- manual-picking (step1)** and **xmipp3- auto-picking (step 2)** (Sorzano et al., 2013)).

The protocol **relion- auto-picking LoG** (Fig. 14) provides particle coordinates in an

automatic way. Together with the 47 input micrographs and the size in pixels for each particle, the protocol form allows to set specific parameters to compute the Laplacian of Gaussian (LoG).

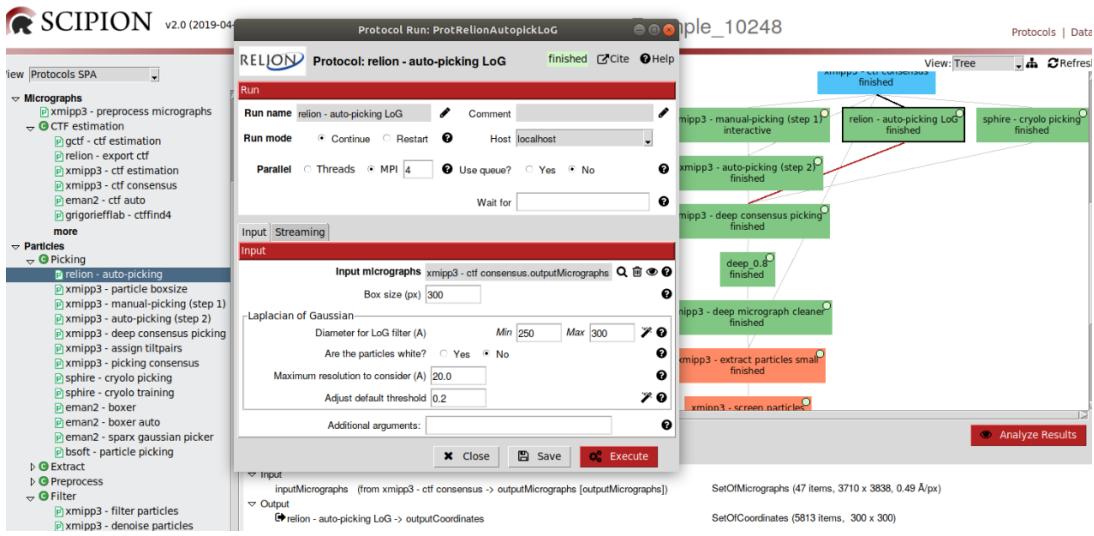


Figure 14: Filling in the protocol `relion- auto-picking LoG`.

The protocol `sphire-cryolo picking` (Fig. 15) integrates a fully automated particle picker based on deep learning. The protocol form also requires the 47 micrographs and the size of particles, and gives you the possibility of using your own network model, obtained in a previous training step, instead the general one. **Confidence threshold** allows to perform a more or less conservative picking by increasing or decreasing the value of this param.

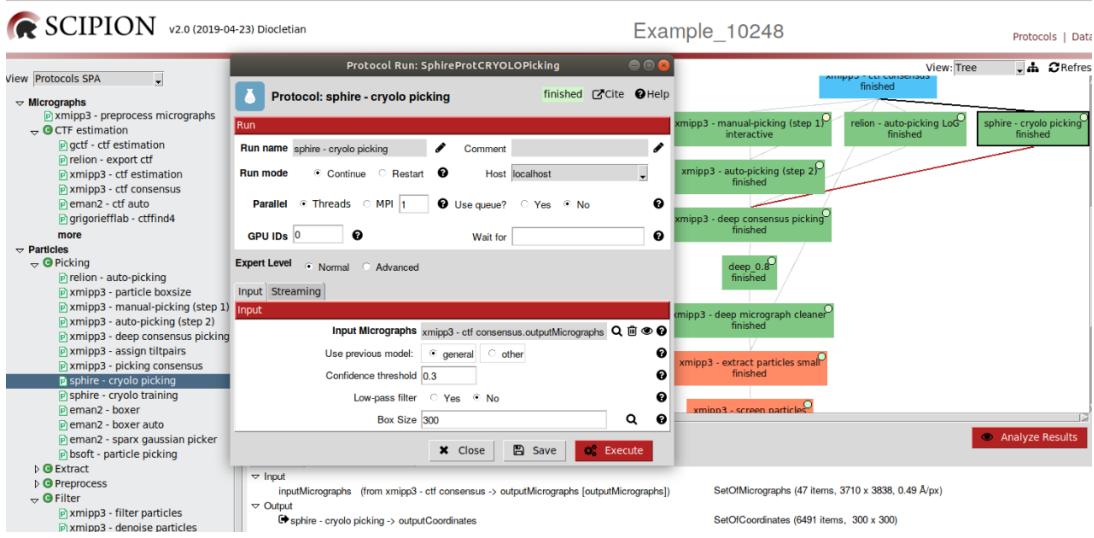


Figure 15: Completing in the protocol `sphire-cryolo picking`.

The protocol `xmipp3- manual-picking (step1)` (Fig. 16) is the first part of the Xmipp picking method, and allows to perform manual picking in a set of micrographs either manually or in a supervised mode. This protocol only requires as input the set of micrographs.

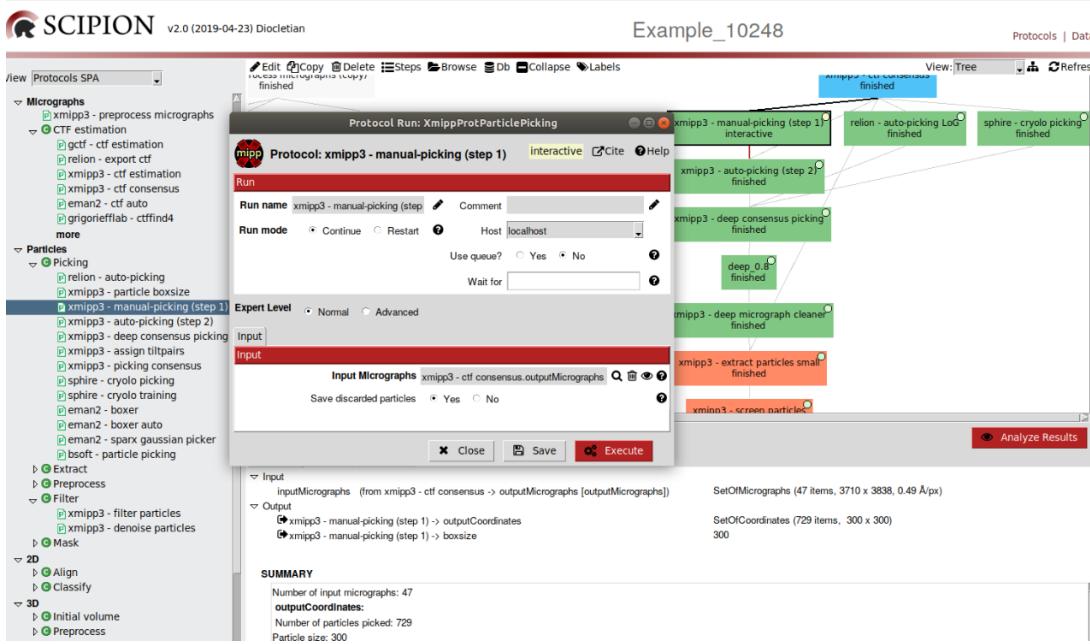


Figure 16: Filling in the protocol `xmipp3- manual-picking (step1)`.

After executing this protocol, the respective box will become light yellow because an interactive job is running and it can be relaunched at any time.

The Xmipp picking GUI contains a control panel with the list of micrographs and some other parameters. The micrograph that we are going to pick is displayed in a separate window and we can apply to it a number of filters/enhancements (like Gaussian blurring, Invert contrast, adjust histogram, etc.) just to improve the visualization of particles. Main control actions are:

- **Shift + Mouse wheel:** Zoom in and out of the overview window.
- **Mouse left button:** Mark particles. You may move its position by clicking the left mouse button on the selected particle and dragging it to a new position.
- **Shift + Mouse left:** Remove a selected particle.
- You can apply filters in the micrographs to see the particles better. Select those filters in the menu **Filters**.

In the manual/supervised step, we start picking manually a few micrographs (5 in this case) and then click the **Active training** button. At this point, the program will train a classifier based on machine learning and will propose some coordinates automatically. You can “correct” the proposal of the classifier by adding missing particles or removing wrongly picked ones. After training with a few more micrographs, we can register the output coordinates by clicking the **Coordinates** red button.

After manual picking, we can close the GUI and open the protocol **xmipp3- auto-picking (step 2)** (Fig. 17). Select as inputs the previous manual/supervised execution and all micrographs (**Micrographs to pick: Other**). The method will pick the rest of micrographs automatically. At the end, we can review the picking coordinates and we still have the chance to add/remove particles.

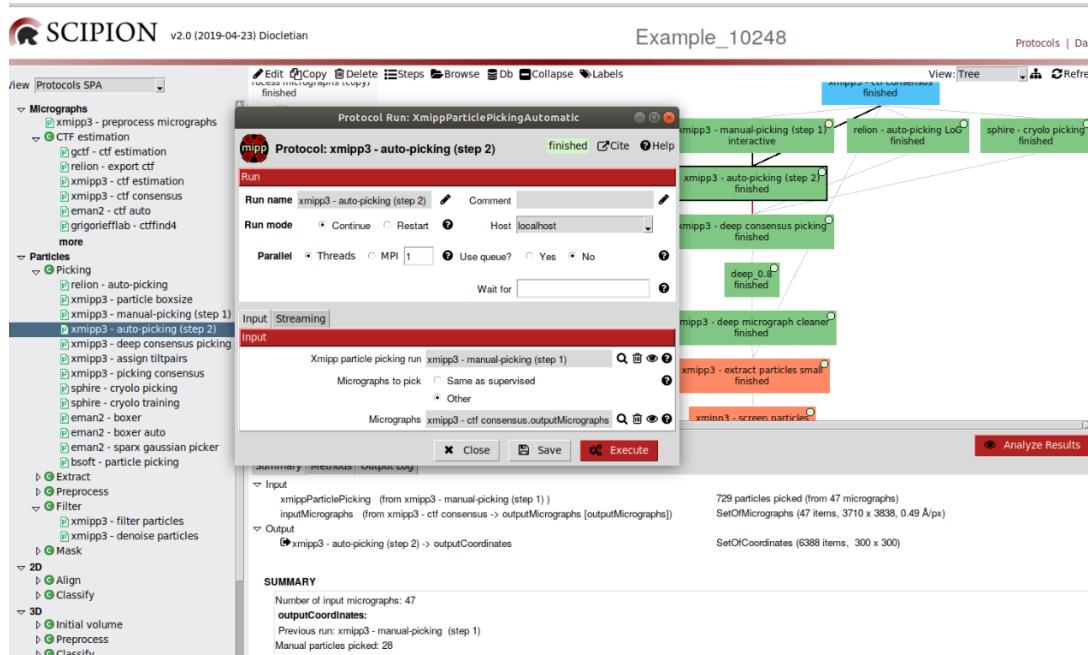


Figure 17: Completing the protocol **xmipp3- auto-picking (step 2)**.

Results of all these protocols can be observed by pressing **Analyze Results**. In all cases a table details the number of particles extracted from each micrograph. Total number of particles appear in the lower part of this table, 6,233, 6,401, 748, and

6,607 running [relion- auto-picking LoG], [sphire-cryolo picking], [xmipp3- manual-picking (step1)], and [xmipp3- auto-picking (step 2)], respectively. As a conclusion, the three algorithms devoted to particle picking give us a similar result, around 6,500 particles. However, there are some differences among programs and we would like to keep only the coordinates of the good particles selected by the three methods.

Consensus in particle picking

The protocol [xmipp3-deep consensus picking] (Fig. 18) will try to select consensus particles among different particle picking algorithms. This protocol can also be used to get the consensus of sets of coordinates retrieved after using distinct setting of parameters with the same program. In our case, the whole sets of coordinates retrieved from the three previous methods have to be included as protocol input.

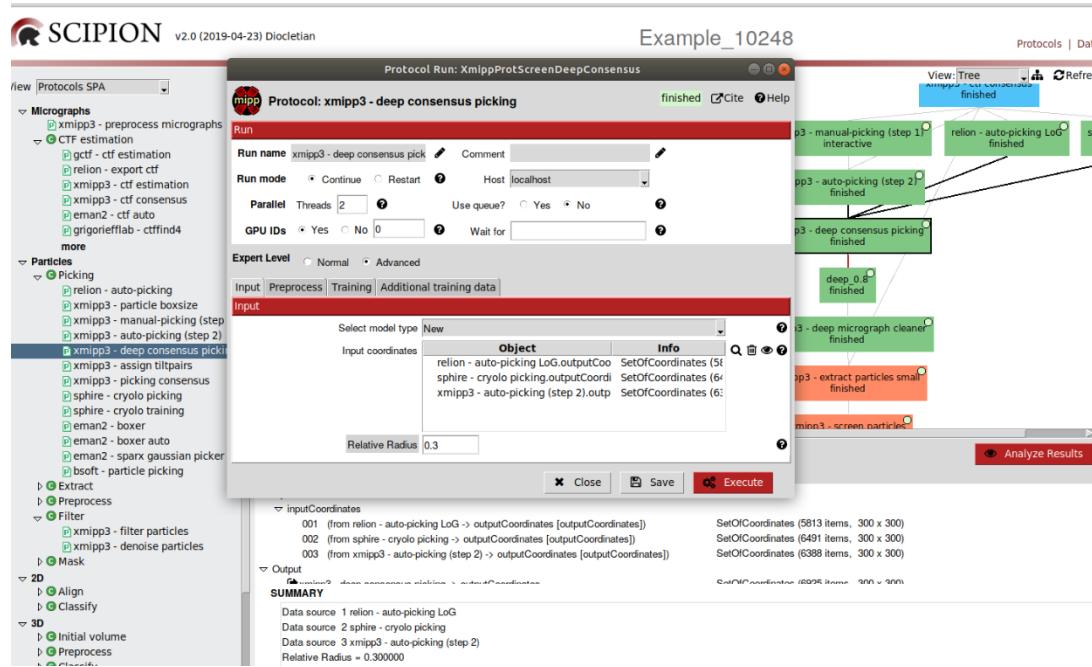


Figure 18: Completing the protocol [xmipp3-deep consensus picking].

A neural network will be trained with subsets of coordinates from particles picked and not picked. Finally, the method provides a score for each particle according to the

neural network predictions. After pressing `Analyze Results`, a menu allows to visualize a table showing an image and the value of the deep learning score of all the particles (`Select particles/coordinates with high 'zScoreDeepLearning1' values`). Considering that bad particles show scores values close to 0.00 and good particles scores close to 1.00, the threshold, automatically set to 0.50, allows to select good particles. In our example, from the total number of particles (8,406), 1,671 particles will be rejected, and around 80% of the total number of particles (6,735) will be selected. We have set, nevertheless, a more restricted threshold of 0.8. This way, 252 additional particles have been rejected and 6,483 particles still remain in the processing workflow.

An additional cleaning step, accomplished with the protocol `xmipp3-deep micrograph cleaner`, removes particles located in carbon zones or in large impurities (Fig. 19). Provide as input the previously selected set of coordinates and indicate the set of micrographs from which the computation will be performed. By default, we use the same as coordinates.

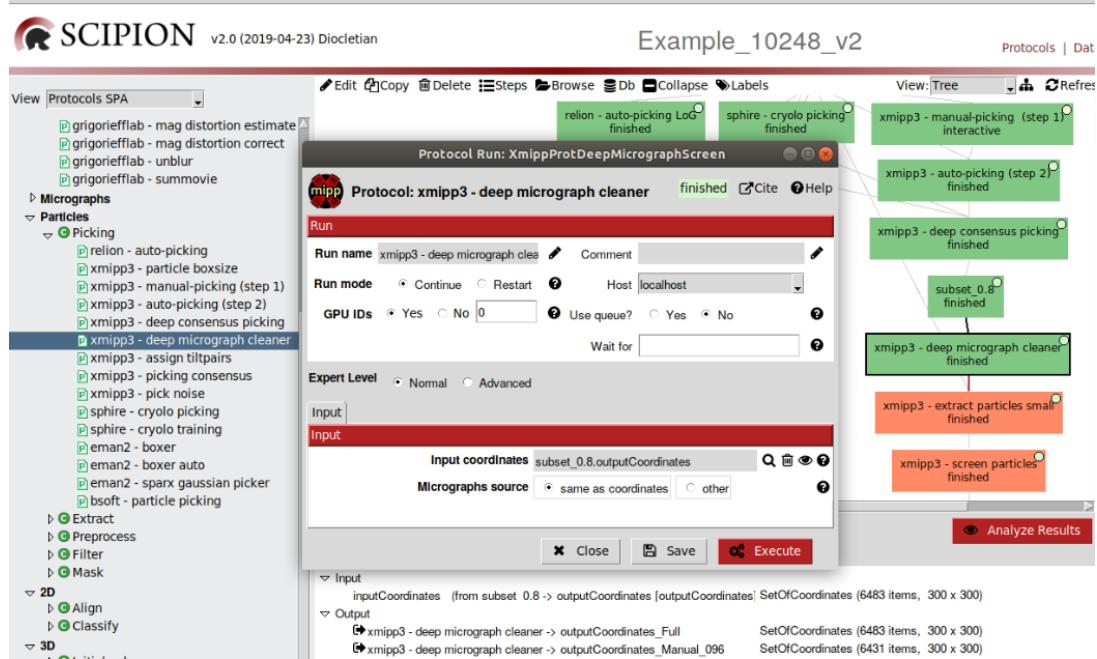


Figure 19: Completing the protocol `xmipp3-deep consensus picking`.

After this additional step of cleaning, other set of 52 particles has been rejected. The coordinates of the remaining reliable 6,431 particles are selected for further processing.

6 Extract Particles

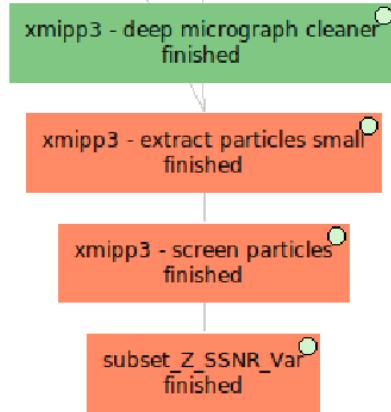


Figure 20: Extract particles workflow (Orange color).

Once we have a set of coordinates, we can proceed to extract particles with **Xmipp** protocol [xmipp3-extract particles](#) (Fig. 21). This protocol allows to extract, normalize and correct the CTF phases of the selected particles. As input, this protocol requires the set of coordinates and the consensus CTF values obtained in previous steps, and a downsampling factor. To save computing resources, include in the input the desired reduced size of the particles. In this particular case, 74 pixels.

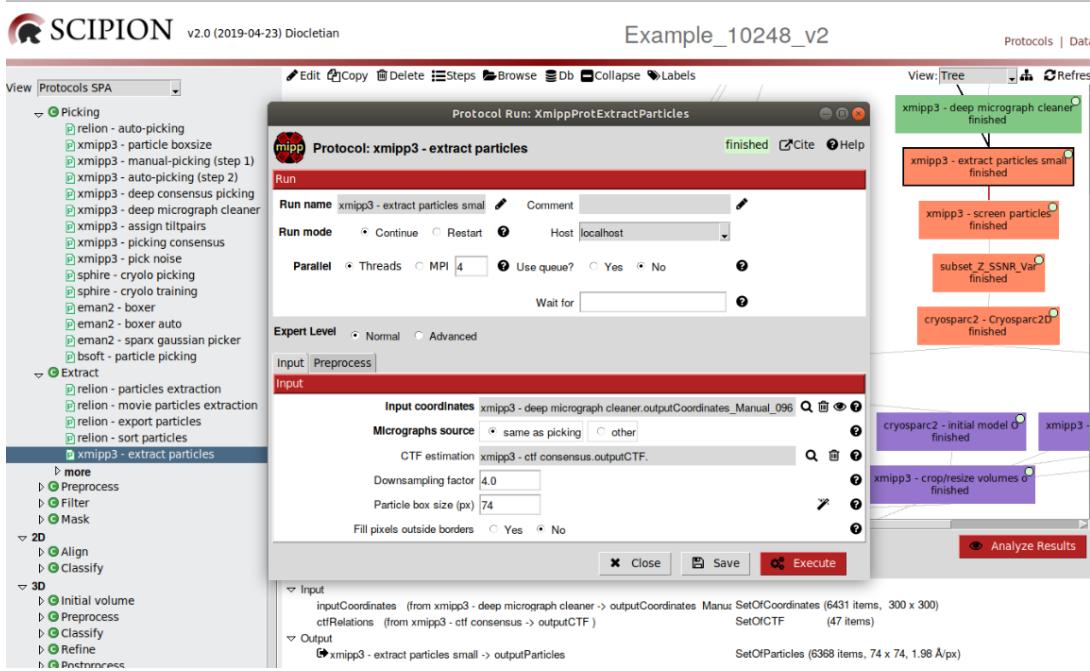


Figure 21: Filling in the protocol `xmipp3-extract particles`.

The form tap **Preprocess** gives you additional options:

- **Invert contrast:** Option **Yes** means that bright regions become dark and the other way around.
- **Phase flipping:** Option **Yes** means that the protocol corrects CTF phases of the particles.
- **Normalize:** Option **Yes** (recommended) means that the particles are normalized with zero mean and one as standard deviation for background pixels.

As output, the protocol generates a new set of 6,368 particles after discarding other 63 particles. The extracted particles have the smaller selected size and 4 times the initial sampling rate. The images of the normalized extracted particles can be seen pressing **Analyze Results**. By default, particles displayed in gallery mode can be sorted by **Zscore**. To visualize the score associated to each particle, switch the table view by pressing the top left button. If you want to remove any of the particles

showing lower score values, select them, press the mouse right button and choose **Disable**. A new subset of particles can be created by clicking on **Particles** red button.

Particle cleaning

Additional cleaning steps of bad particles can be performed with other screening protocols such as **xmipp3-screen particles** (Fig. 22). The protocol input is the subset of particles previously generated. Three different criteria for rejection can be selected, **Zscore**, **SSNR** and **Variance**. Zscore assesses the similarity of each particle with the average. SSNR evaluates the signal to noise ratio in the Fourier space. The variance is assessed for each particle in the context of particles where that particle was picked. In this particular case, we are not going to set any of the mentioned params (Zscore, SSNR and Variance). Instead, selection will be performed after visualizing the plots of those statistics.

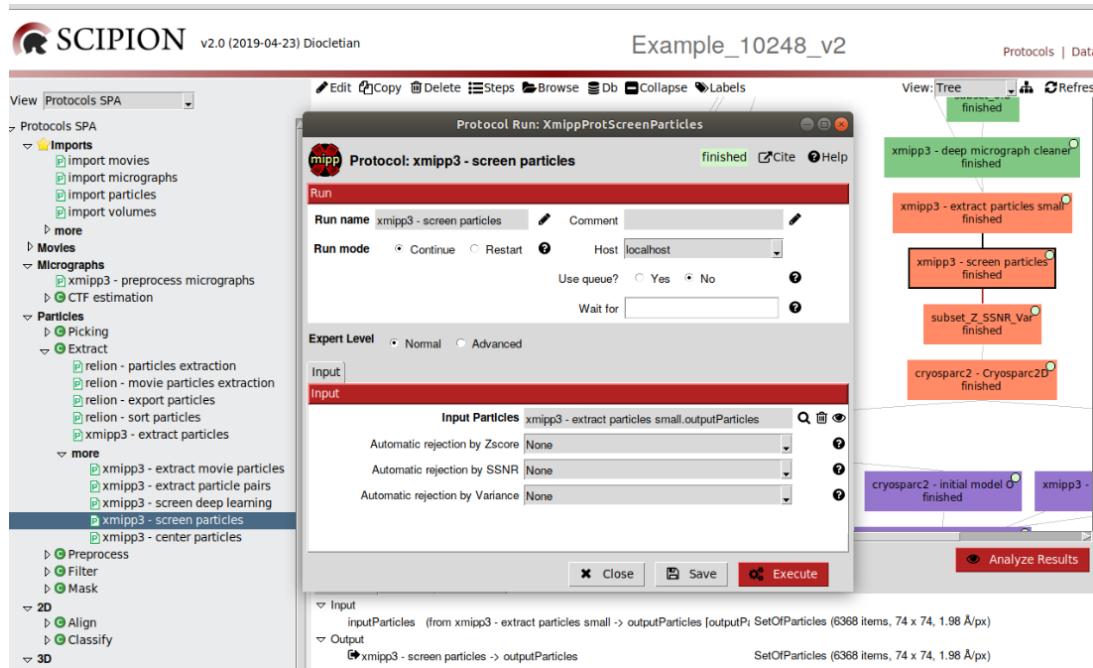


Figure 22: Completing the protocol **xmipp3-screen particles**.

After executing the protocol without discarding any particles, we press **Analyze Results**. The plot of **Zcore** and the **Variance** histogram will be open, together with the table of particles. According to the **Zcore** plot, we discard particles with **Zcore** value higher than 3.0. According to the **Variance** histogram, we reject particles with **Variance** higher than 1.21. We can also visualize the histogram of **_xmipp_cumulativeSSNR**. According to this histogram, particles with **SSNR** values lower than 2.5 and higher than 5 will be discarded. The new subset of “cleaned” particles obtained according those specific criteria (**subset_Z_SSNR_Var**) contains 5,913 particles, after discarding 455 of them (7.1%), that can be observed pressing **Analyze Results**. This new subset of selected particles is considered reliable for further image processing.

7 2D classification



Figure 23: 2D classification workflow.

The next step in image processing involves the 2D classification of the particle images to group similar ones. This process can serve as an exploratory tool of your data and might also be used to throw away bad particles. In addition, by overlapping similar images we can obtain the average images or 2D classes. Since these classes are the projections of the 3D object that we try to reconstruct, they can also be used in the reconstruction of the 3D object.

Although there exist several 2D classification algorithms, in this tutorial the 2D classes will be created with the *cryoSPARC* (Punjani et al., 2017) 2D classification method (Punjani et al., 2016), integrated in the protocol **cryosparc2- 2d classification**.

We used as input the subset of particles previously selected. The 2D classification parameters can be observed in the central tap of the protocol form in the Fig. 24. Remark that we choose in advance the number of classes, 16 in this case.

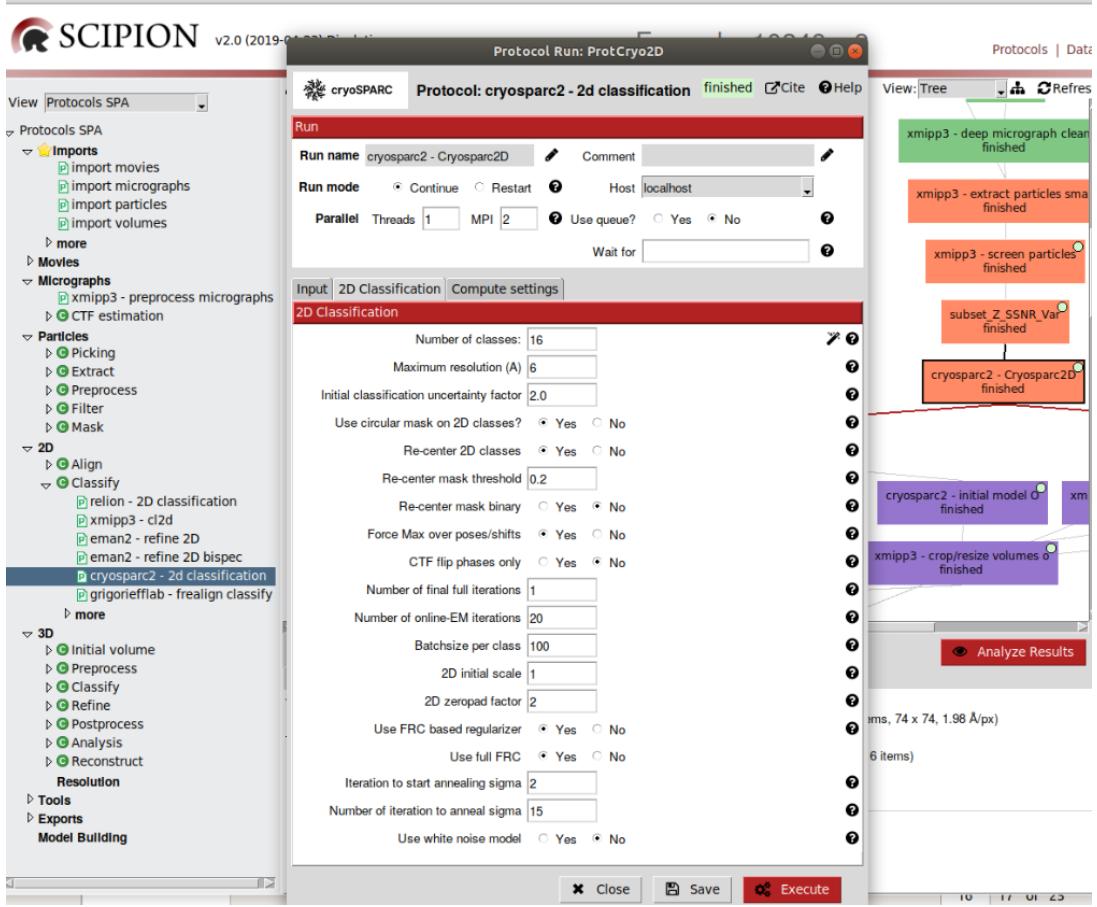


Figure 24: Filling in the second tap of the protocol [cryosparc2-2d classification](#).

After running the protocol, particle classes can be visualized selecting any of the two options of the menu opened with Analyze Results, the common *Scipion* viewer or the *cryoSPARC* GUI. The final number of classes also appears in the Summary output detailed in the Fig. 25.

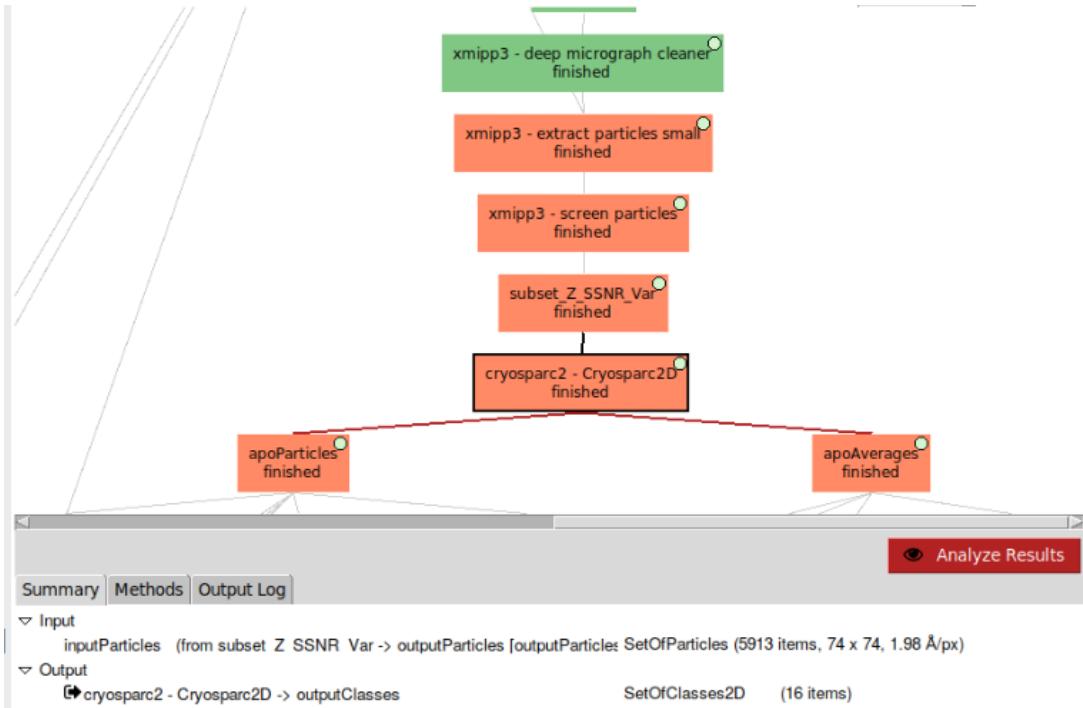


Figure 25: Summary of *cryoSPARC* 2D classification results.

Two branches derive from the `cryosparc2- 2d classification` protocol box (Fig. 25), pointing to two boxes, `apoParticles` and `apoAverages`, which include the whole set of particles and the 2D classes, respectively. These two branches can be obtained by clicking in the lower part of the Summary (`cryosparc2 - Cryosparc2D -> outputClasses`). A panel will be displayed with the 2D classes. By selecting the classes that we are interested in (14 from 16) and pressing `Particles`, a new set of 5,673 aligned particles will be created, included in the box `apoParticles` (see Summary output of the Fig. 26).

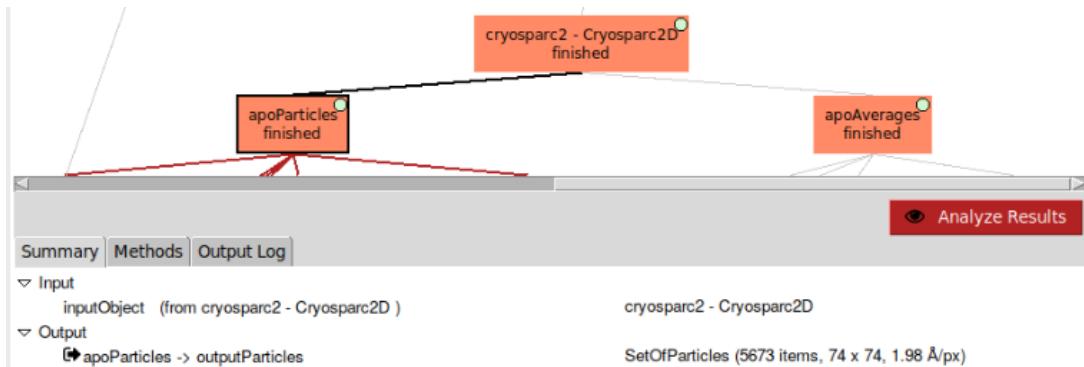


Figure 26: Summary of particle selected in **apoParticles** box.

If, instead, we press **Averages**, a new set of 14 class representative particles will be created, included in the box **apoAverages** (Fig. 27, Summary output).

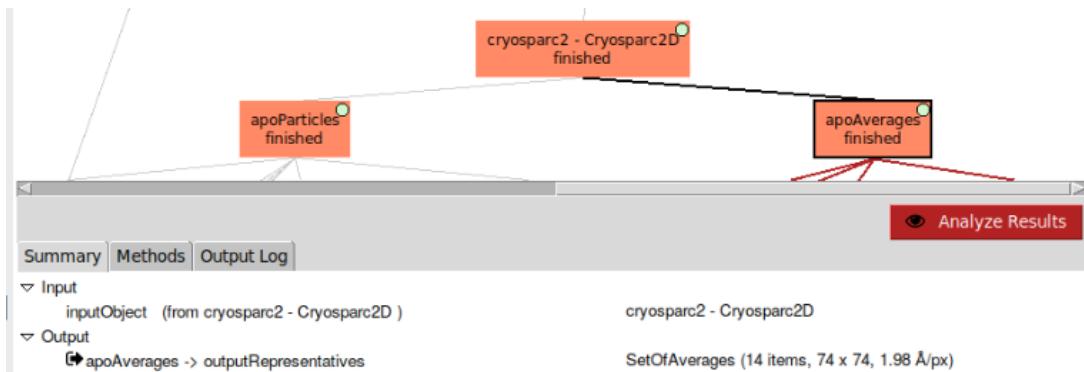


Figure 27: Summary of particle selected in **apoAverages** box.

Finally, if we press **Classes** after selecting some of the classes, a new set of classes with the number of the selected classes will be created.

The selected elements in the two mentioned branches, individual particles and class representative particles, will be used in the next step in the processing workflow to generate the initial volume.

8 Initial volume

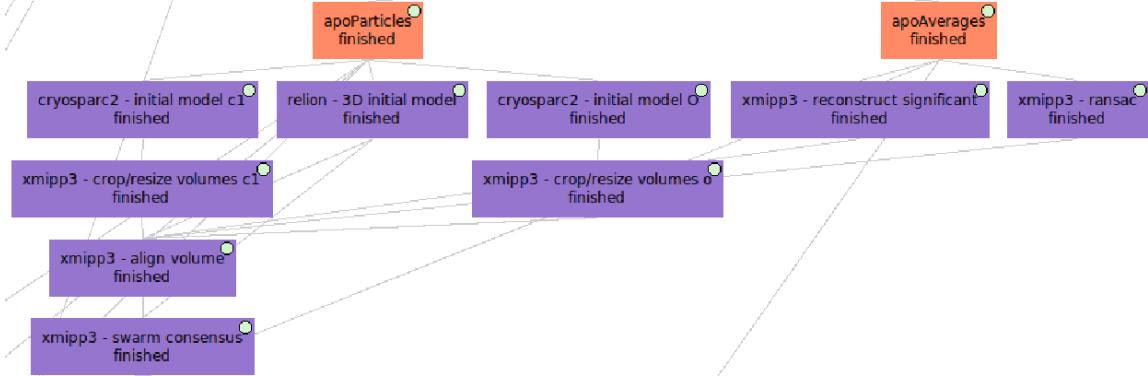


Figure 28: Initial volume (Purple color).

Right angles of each projection image are needed to reconstruct the 3D map. In 3DEM, however, these angles are unknown and we have to estimate them. The most popular way of estimating them is comparing the projections of a volume similar to ours, known as initial volume, with the images obtained from the microscope. Since the 3D map reconstruction process requires an approximate low resolution map to be refined in further steps according to the projection images of particles, in this tutorial we are going to generate a *de novo* initial map model combining the results obtained by different algorithms: First, to compute the initial volume using the set of aligned particles as input, we have used *cryoSPARC Stochastic Gradient Descent* (SGD) and *Relion Stochastic Gradient Descent* (SGD) (Fig. 28, left). Second, to generate the initial volume from the class representative particles, we have run *Xmipp reconstruct significant* and *Xmipp RANSAC* (Fig. 28, right). *Xmipp reconstruct significant* sets the map in a **Weighted Least Square** framework and calculates weights through a statistical approach based on the cumulative density function of different image similarity measures. *Xmipp RANSAC* is based on an initial non-lineal dimension reduction approach with which selecting sets of class representative images that capture the most of the structural information of each particle. These reduced sets will be used to build maps starting from random orientation assignments. The best map will be selected from these previous assumptions

using a random sample consensus (RANSAC) approach.

cryoSPARC Stochastic Gradient Descent (SGD)

The algorithm *cryoSPARC* Stochastic Gradient Descent (SGD) has been implemented in the protocol `[cryosparc2-initial model]`, that we execute twice with two different values for params Number of Ab-Initio classes and Symmetry (Ab-initio reconstruction tap of Fig. 29 and Fig. 30). The set of particles selected in the previous step is used as input in both cases (see the Input tap of the protocol form).

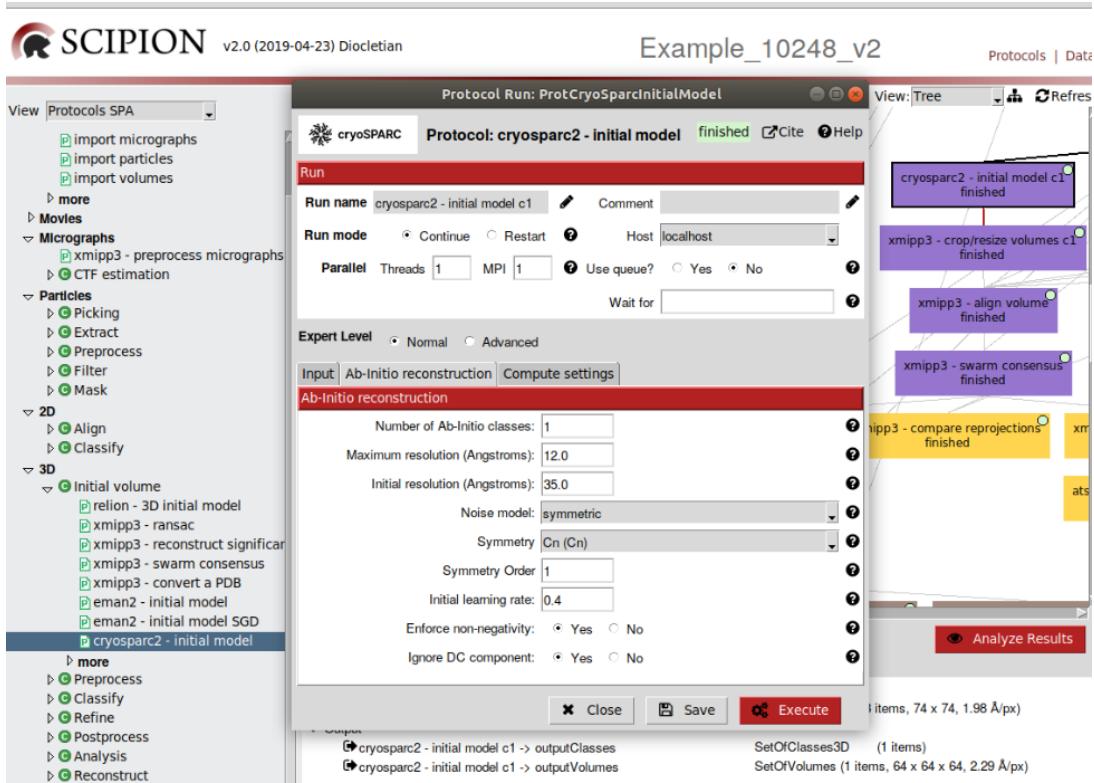


Figure 29: Completing the second tap of the protocol `[cryosparc2-initial model]` with one Ab-Initio class and no Symmetry.

In the first case, only one Ab-Initio class has been selected and C1 Symmetry

has been considered (Fig. 29). In the second one, instead, three Ab-Initio classes and octahedral Symmetry have been selected (Fig. 30). Each Ab-Initio class will be randomly initialized, unless an initial map is provided. In that case, each class will be a random variant of the initial map. Regarding symmetries, enforcing symmetry above C1 is not recommendable for *ab-initio* reconstructions.

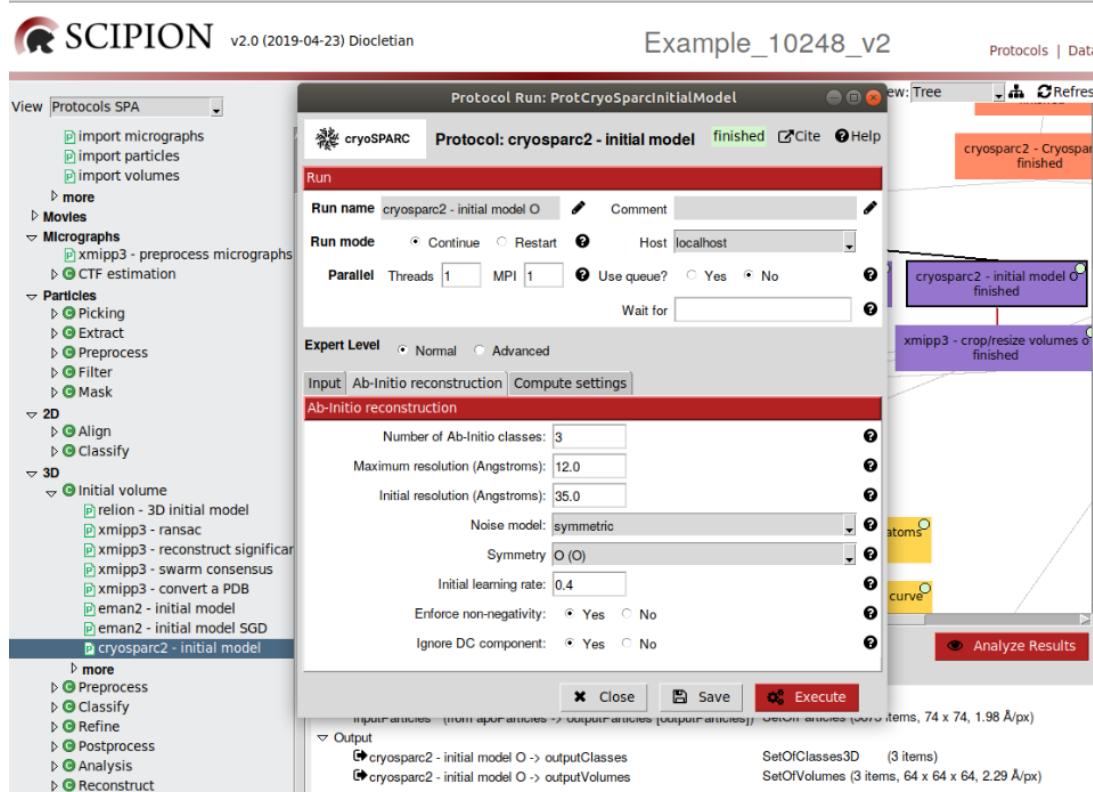


Figure 30: Filling in the second tap of the protocol `cryosparc2-initial model` with three Ab-Initio classes and octahedral Symmetry.

Maps or 3D classes generated can be appreciated by pressing `Analyze Results`. One map has been generated in the first case and three in the second one, although only one of these three maps has been built with most of particles (5,643).

Since the size and sampling rate of maps generated with `cryosparc2-initial model`

differ from the size and sampling rate of the input particles, a resizing intermediate method has to be applied to recover these dimensions. Protocol `xmipp3-crop/resize volumes` will help us with this task (Fig. 31). As input, select the output volumes of the previous protocols, `Dimensions` for `Resize option`, and 74 pixels as `New image size`.

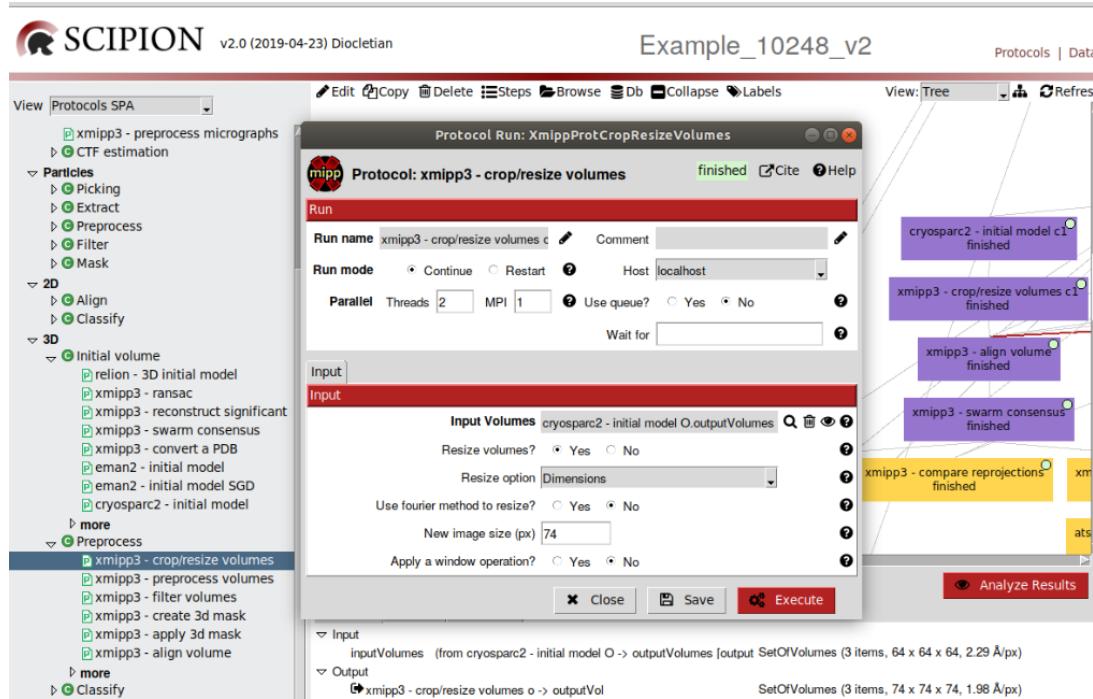


Figure 31: Completing the protocol `xmipp3-crop/resize volumes`.

Relion Stochastic Gradient Descent (SGD)

Relion Stochastic Gradient Descent (SGD) has been implemented in the protocol `relion-3D initial model`. Fill in the param values and execute it (Fig. 32).

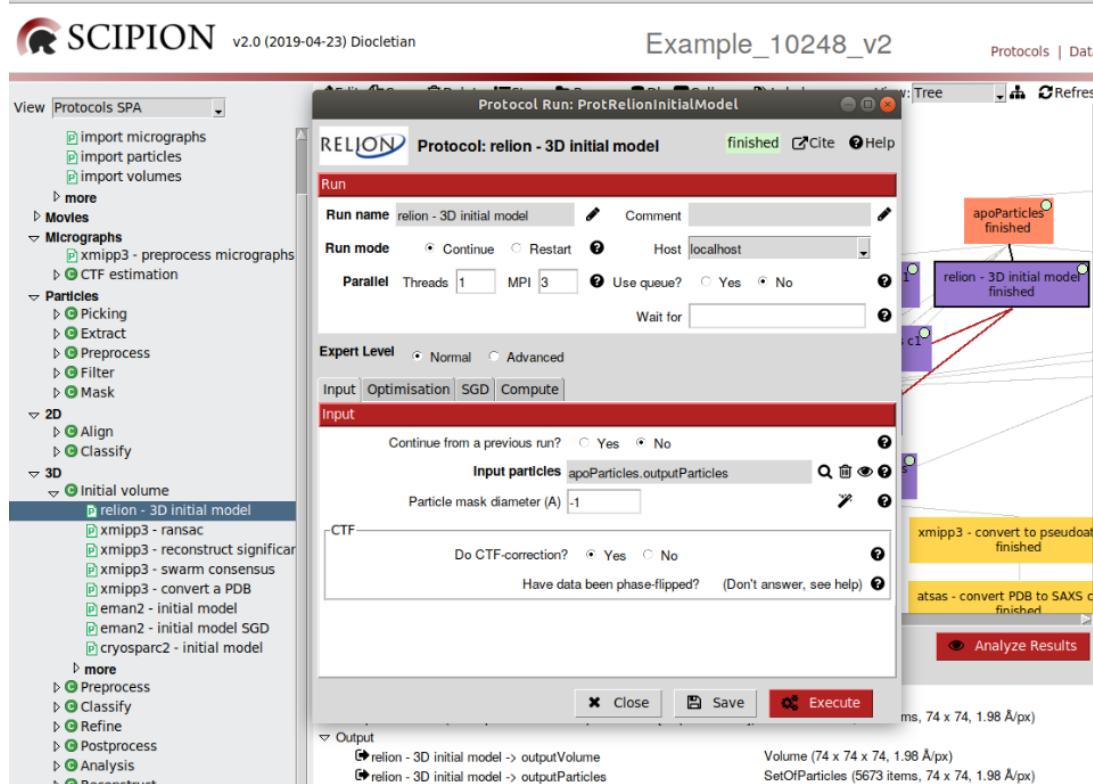


Figure 32: Filling in the Input tap of the protocol `relion-3D initial model`.

Only one volume has been generated with this protocol that keeps size and sampling rate of the input particles. You can visualize it with *Chimera* in 3D by pressing `Analyze Results` and selecting in the Volumes box `Display volume with chimera`.

Xmipp

Using the 14 class representative particles as input, as well as the type of symmetry (octahedral), the protocol `xmipp3-reconstruct significant` (Fig. 33) also generates one initial volume and preserves the size and sampling rate of the input 2D classes.

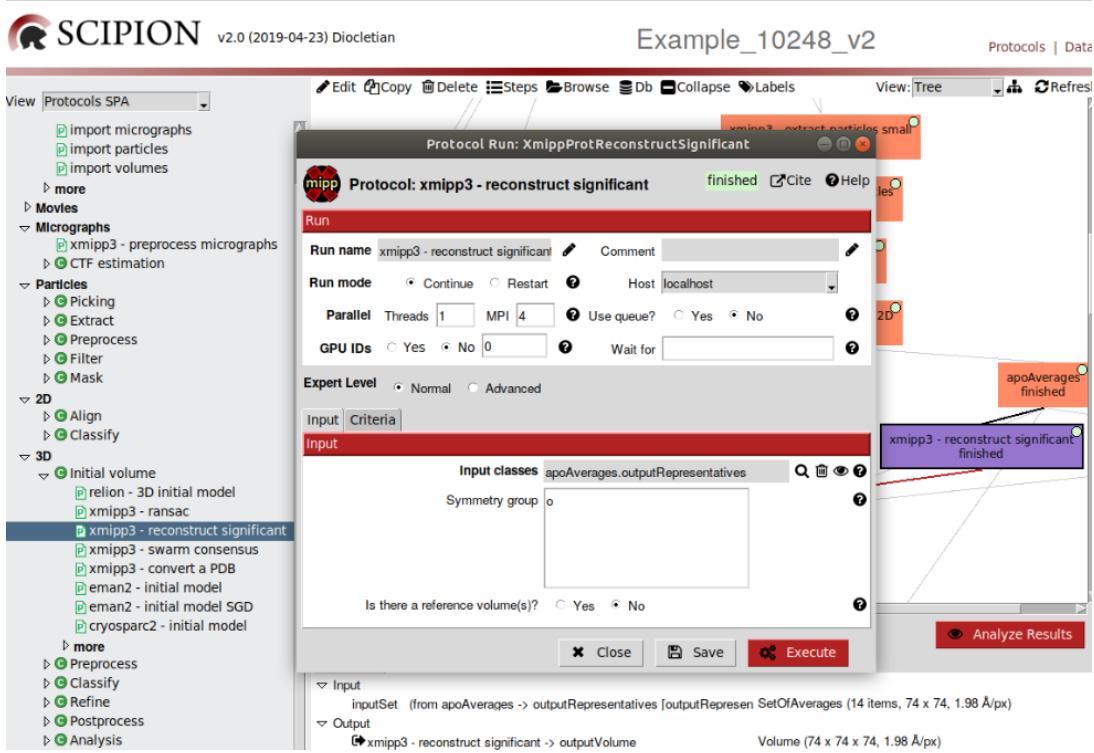


Figure 33: Completing the Input tap of the protocol **xmipp3-reconstruct significant**.

Xmipp RANSAC algorithm, implemented in the protocol **xmipp3-ransac** (Fig. 34), although starts from the same input than *Xmipp* **reconstruct significant**, generates 10 different maps and preserves the size and sampling rate from the input 2D classes. You can choose a different number of output maps in the advanced param **Number of best volumes to refine**.

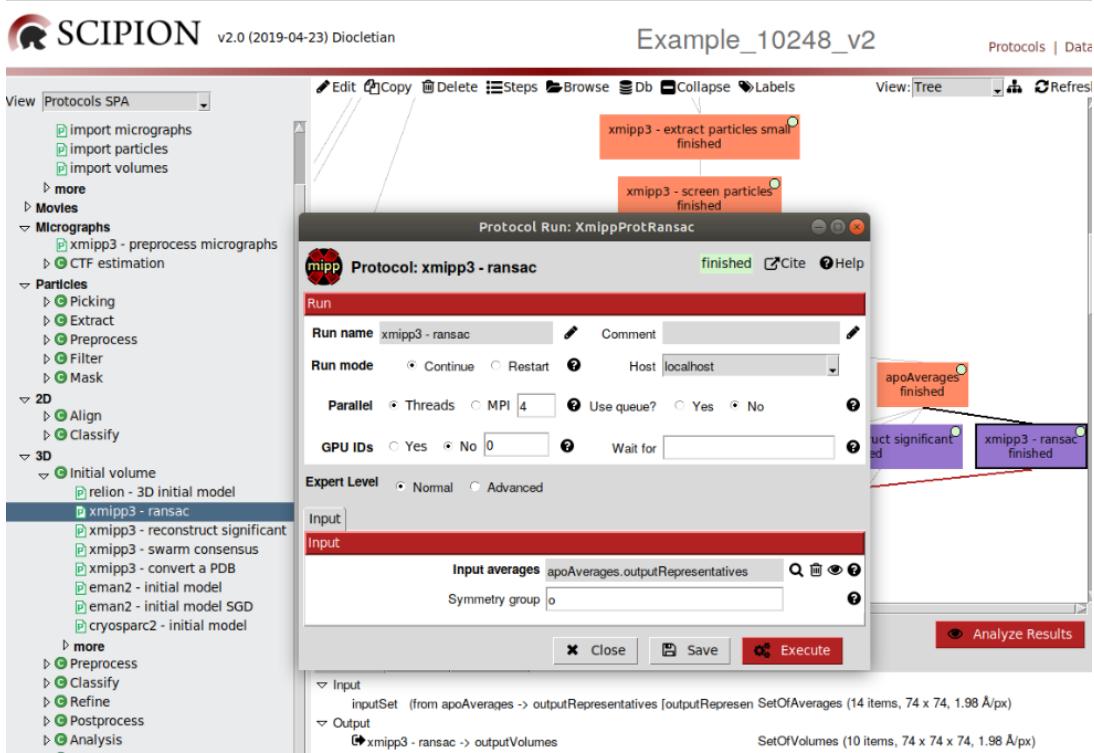


Figure 34: Filling in the params of the protocol `xmipp3-ransac`.

Map alignment and swarm consensus

Next, we perform a local alignment of the 16 maps generated, starting both from particles and 2D classes, using the protocol `xmipp3-align volume` (Fig. 35). As Reference volume we select the initial volume obtained by the *Relion* algorithm.

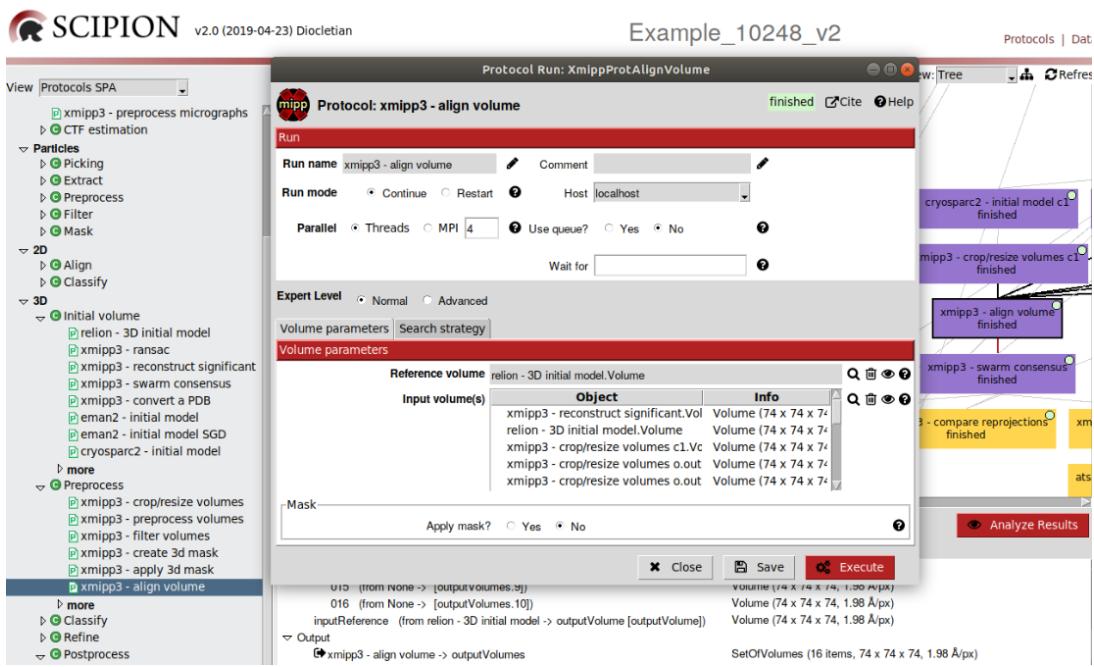


Figure 35: Completing the params of the protocol **xmipp3-align volume**.

A new set of 16 maps has been created keeping the same size and sampling rate shown by the initial particles. These maps can be visualized by pressing **Analyze Results**.

Next, in order to have only one initial volume partially refined against the selected set of particles, we use the protocol **xmipp3-swarm consensus** (Fig. 36). The inputs of this protocol are the set of 16 maps and the set of 5,673 *Relion* extracted particles, previously generated. In this case, maps and particles have the same size and sampling rate. The program try to optimize the correlation between the swarm of volumes and the set of particles. Only a fraction of the particles are used to update this stochastic maximization.

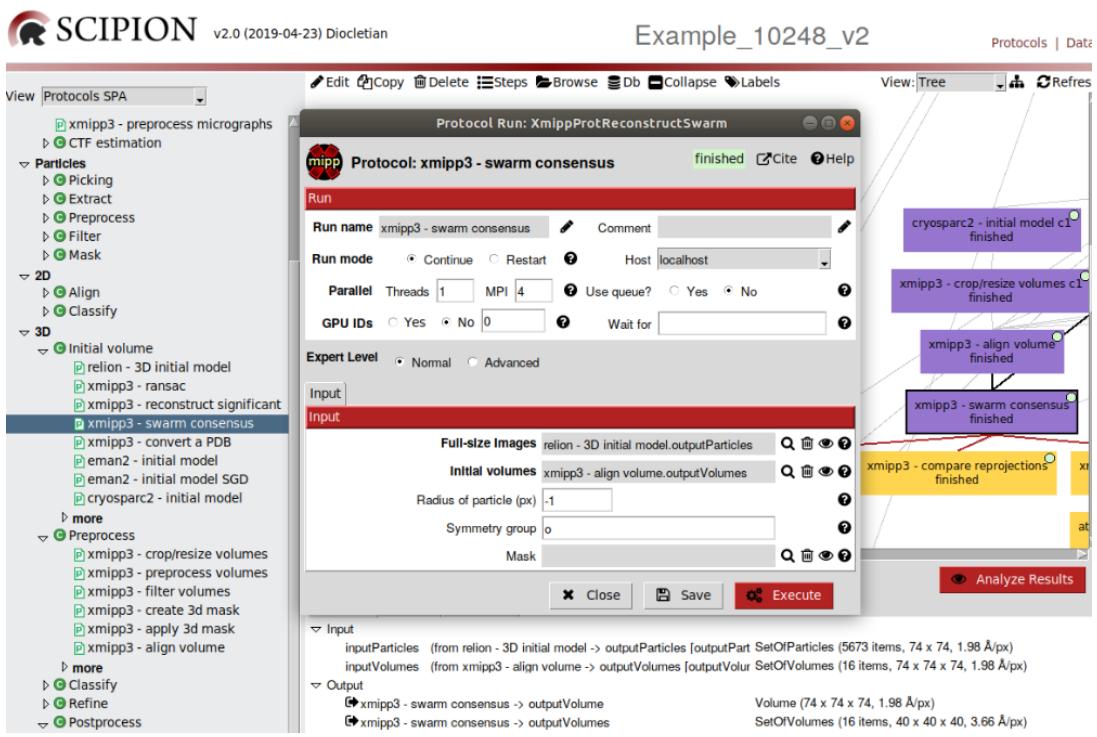


Figure 36: Filling in the params of the protocol `xmipp3-swarm consensus`.

After 15 iterations, this protocol generates two outputs, 16 downsampled maps and a map with the size and sampling rate of the inputs. This map will be the initial volume that will be further refined.

Validation of the initial volume

To check the reliability of this initial volume, we can compare its projection images with the set of representative 2D classes using the protocol `xmipp3-compare reprojections` (Fig. 37). Symmetry group has to be also included in the input form. This 3D validation protocol computes residues, *i.e.*, differences between the experimental images and reprojections, as well as other statistics. The values of these statistics might suggest the presence of outliers.

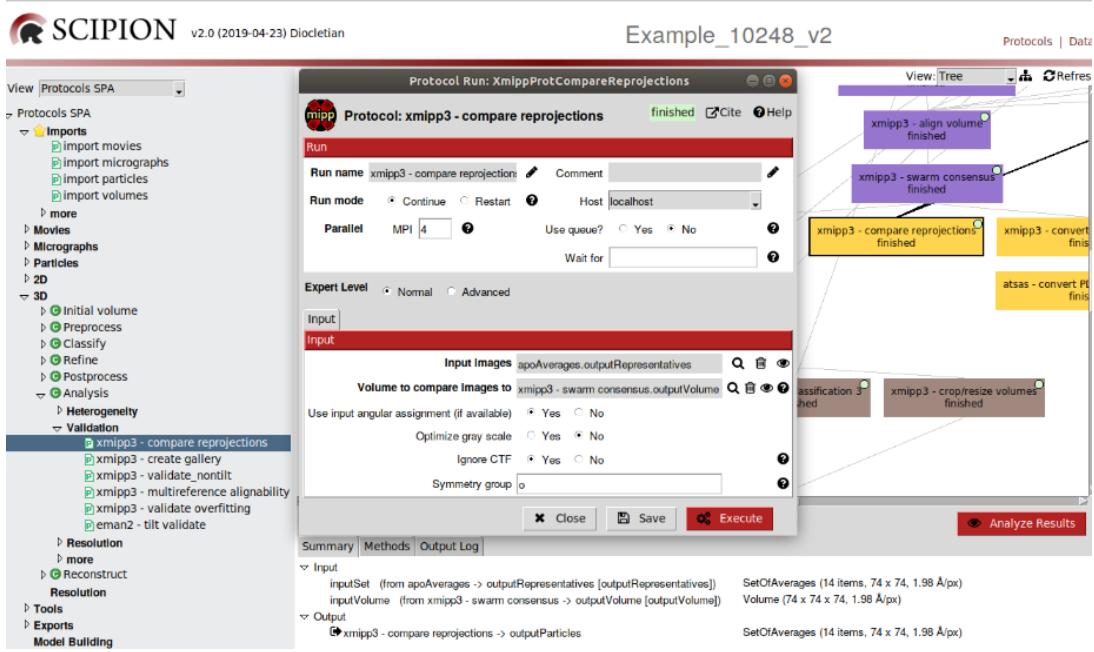


Figure 37: Completing the params of the protocol `xmipp3-compare reprojections`.

In our example, values of `_xmipp_cost` between 0.8115 and 0.8725 (press `Analyze Results`) suggest quite good similarity among initial volume projection images and the set of 14 2D classes.

From the initial volume to a pseudoatom structure

The initial map can also be used to generate an initial pseudoatom structure to be used for inferring the SAXS (Small-Angle X-ray Scattering) curves. Two consecutive executed protocols perform these two tasks, `xmipp3-convert to pseudoatoms` and `atsas-convert PDB to SAXS`.

The protocol `xmipp3-convert to pseudoatoms` (Fig. 38) requires the initial map, previously computed, and the `Pseudoatom radius (vox)` as inputs.

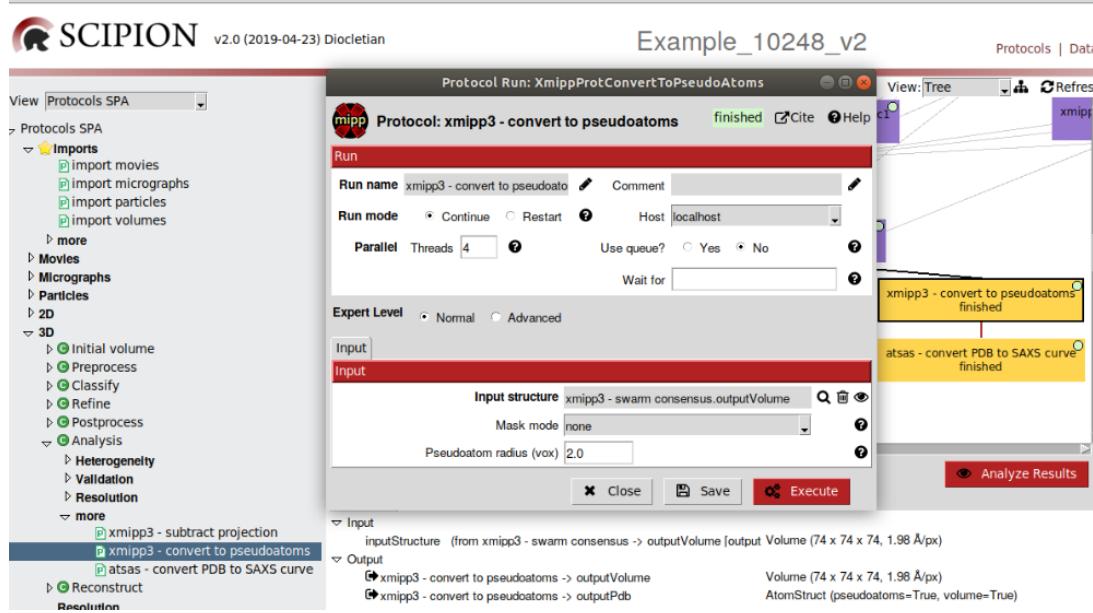


Figure 38: Filling in the params of the protocol `xmipp3-convert to pseudoatoms`.

The pseudoatomic structure obtained can be visualized with *Chimera* (press **Analyze Results**). This structure can be used as input of the protocol `atsas-convert PDB to SAXS` (Fig. 39) to generate the respective SAXS curves by running the program *Crysol* from *Atsas* (Svergun et al., 1995). Press **Analyze Results** to visualize the curves in solution and in vacuo.

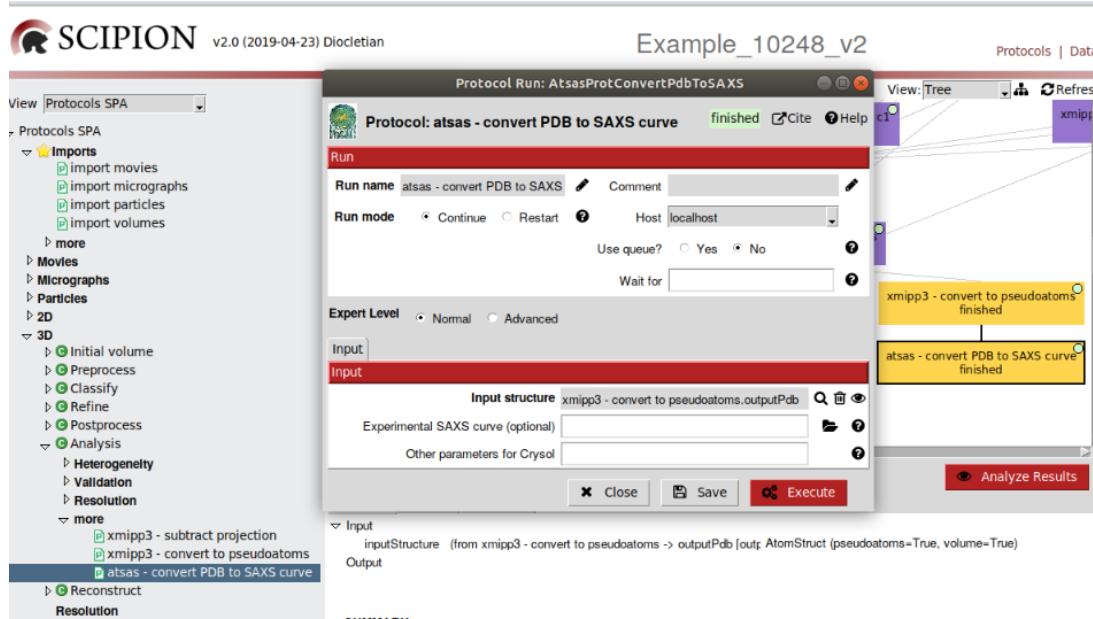


Figure 39: Completing the params of the protocol **atsas-convert PDB to SAXS**.

9 3D classification and Refinement

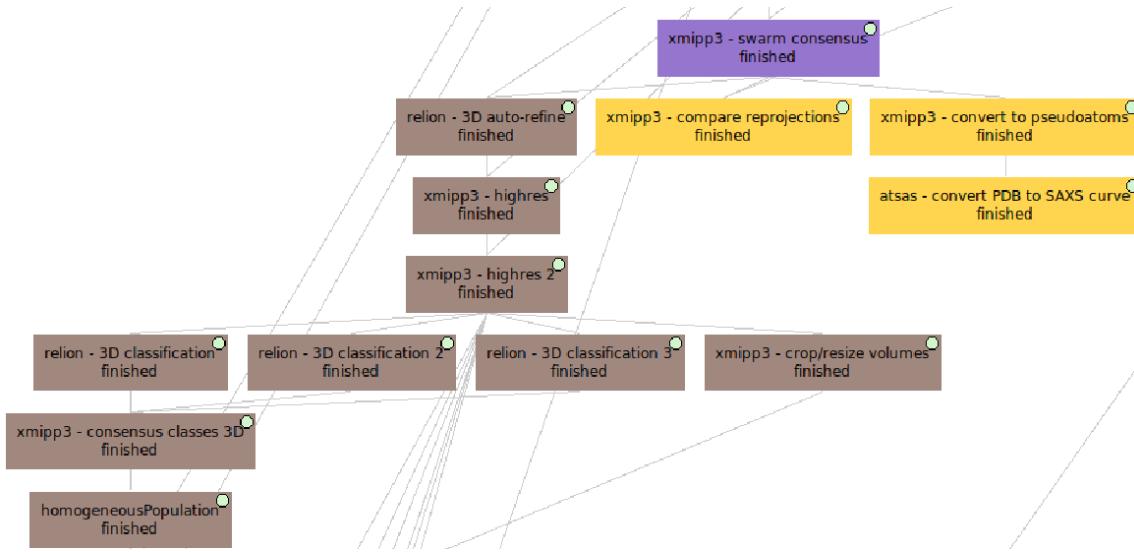


Figure 40: Refinement and 3D classification (Brown color).

3D classification and Refinement are the two last overlapping steps in image processing. They consume the most time and resources with the aim of obtaining a 3D map at the highest possible resolution. This is only feasible if data are homogeneous enough, *i.e.*, if data represent a unique conformation of the specimen.

Refinement of the initial map

Before starting with the 3D classification properly, three consecutive steps of refinement will be performed with our initial map. The first approach to get a high resolution map in a fully automated manner was performed with the algorithm *Relion auto_refine*, based on an empirical Bayesian approach. This procedure employs the so-called gold-standard Fourier Shell Correlation (FSC) to estimate the resolution. Combined with a novel procedure to estimate the accuracy of the angular assignments, the algorithm converges. We have implemented it in the protocol `relion- 3D auto-refine` (Fig. 29). In the Input tap of this protocol form we include the

subset of homogeneous particles selected previously. The initial volume will be included in the Reference 3D map tap, as well as a Initial low pass-filter (A) of 60.0. This tap gives you the possibility of using Reference mask (optional) and, in some cases, *e.g.* non-empty icosahedral viruses, a Second reference mask (optional).

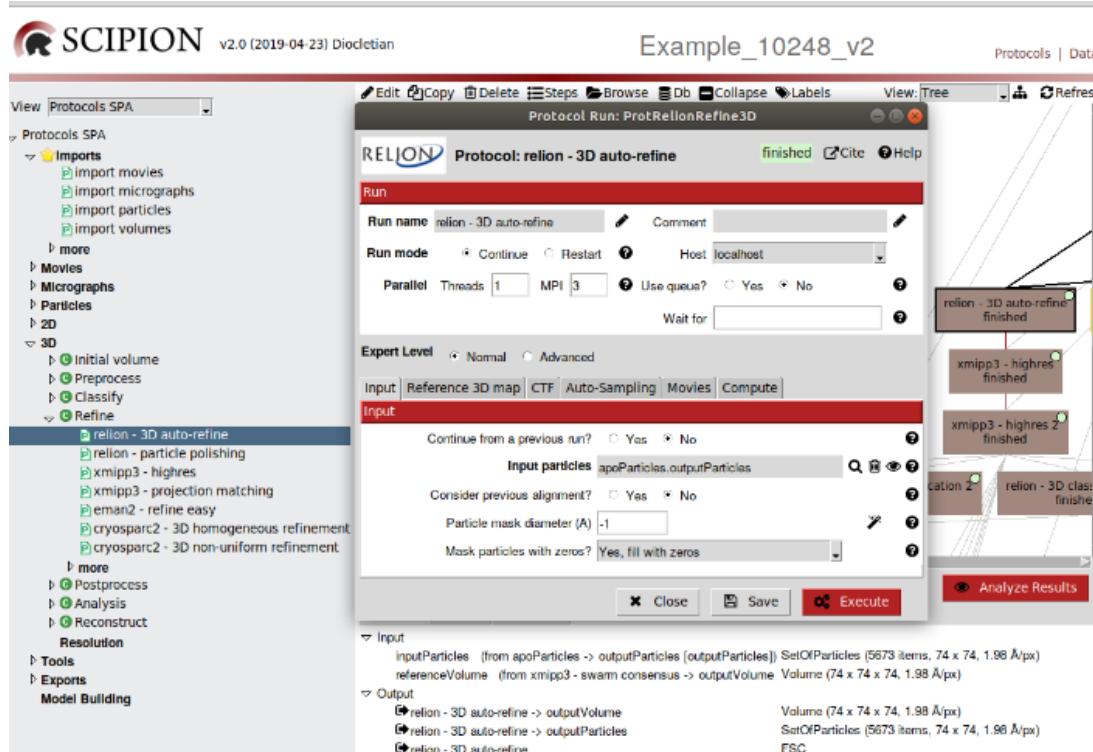


Figure 41: Completing the params of the protocol
relion- 3D auto-refine.

There are three questions in the tab CTF:

- Do CTF correction?, set to Yes to perform full phase + amplitude CTF correction.
- Has reference been CTF-corrected?, set to No because the Fourier transforms of the reference projections are not multiplied by the CTF in the first iteration.

- Do manual grouping ctfs?, set to No because we have enough number of particles that we do not need to group them.

The Angular sampling interval (deg) option in the tab Auto-Sampling will be used only in the first few iterations. Later, the algorithm will automatically increase its value until convergence. For symmetries lower than octahedral or icosahedral we use the default values of Angular sampling interval (deg) and Local search from auto-sampling (deg).

Movies tab allows to align movie-particles of each frame and execute later a protocol of particle polishing.

After executing 7 iterations, a refined map of 5.05 Å of final resolution was obtained as output, with the same size and sampling rate that we had in the inputs. Press **Analyze Results** and visualize any iteration or the last one by default. Concerning particles, their angular assignment and the `_optimiser.star` file, with general information about the refinement process, can be shown. Different volumes can be 2D or 3D visualized, such as each half map, both, or the final one. SSRN and resolution FSC plots are also available.

The next two following steps of refinement have been performed with the *Xmipp* algorithm `highres` (Sorzano et al., 2018) that we have implemented in the protocol `xmipp3-highres` (Fig. 42). This method computes a weight for each particle and performs both global and local alignment. Iterations can be performed one by one, removing particles that worse fit the map from one iteration to the next one. This 3D refinement protocol uses as input the refined map obtained from *Relion auto-refine* and the same set of particles used by this algorithm. The Symmetry group has been also included in the Input tap. In the Angular assignment tap we choose Global as Image alignment and 1 Number of iterations with 4 Å as Max. Target Resolution.

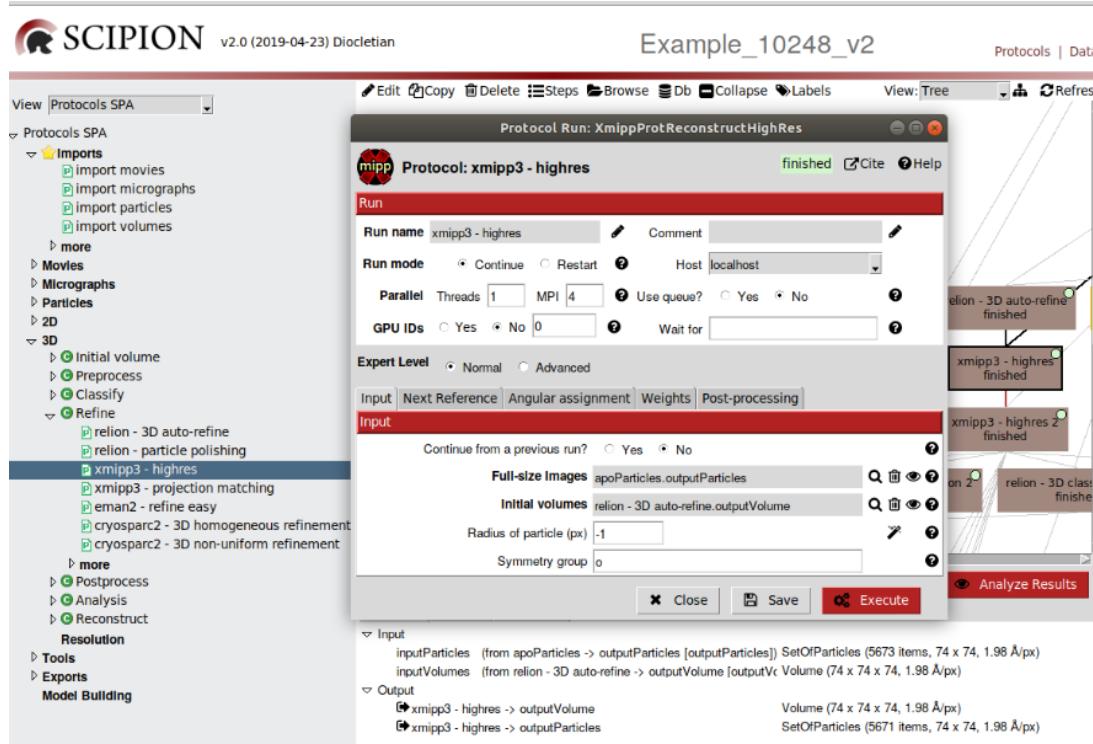


Figure 42: Filling in the params of the protocol `xmipp3-highres`.

This protocol also generates one map as output with the initial size and sampling rate. Press `Analyze Results` to check the results of the iteration 1. Particles and map can also visualized. 2 particles have been rejected.

The second time that we execute the protocol `xmipp3-highres` we select the map obtained in the previous step and the same set of particles as input. In the tap `Angular assignment` we replace `Global` by `Local` in the `Image alignment` param. The rest of the params remain unchanged since we have selected `Yes` in the param `Continue from a previous run?` of tap `Input`. In this case, another map has been generated as output with the same size and sampling rate. 6 particles have been discarded this time.

3D classification

To continue with the refinement process and to obtain a better resolution, we start executing three independent times the same algorithm of *Relion 3D classification* that we have implemented in the protocol `relion-3D classification` (Fig. 43). In the tap **Input** we include the particles derived from executing the previous protocol `xmipp3-highres`. The map derived from this protocol will be the **Input volume(s)** in the tap **Reference 3D map**. The optimization params appear in the tap **Optimization**: 3 Number of classes and 25 Number of iterations. As Regularisation parameter T values as 3-4 are common for 3D classification.

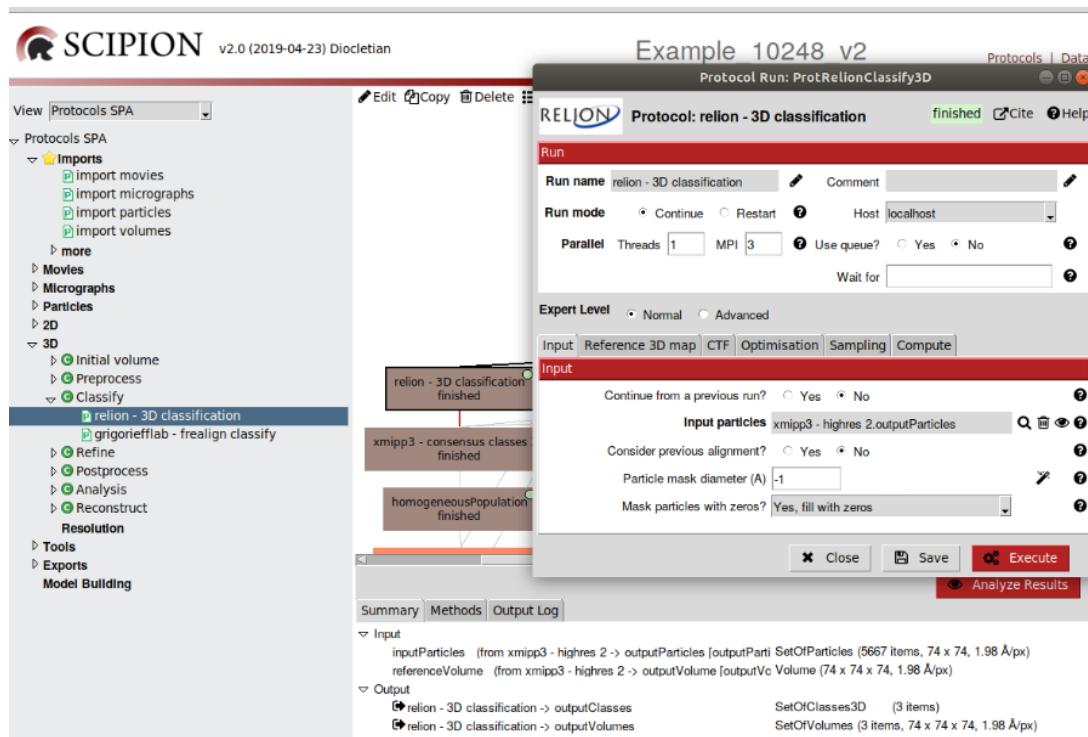


Figure 43: Completing the params of the protocol `relion-3D classification`.

The output of each one of these three `relion-3D classification` protocols are 3 maps with the initial size and sampling rate, reconstructed from different groups of reclassified particles. By pressing `Analyze Results` and `Particles/ Show classification`

in Scipion, a table will be opened showing the projection representative of each map and the number of particles contributing to its reconstruction:

- 4,713, 886 and 68 in the first classification.
- 2,854, 2,739 and 74 in the second one.
- 4,964, 583 and 120 in the last one.

The results of the first and third 3D classifications are similar because in both cases the first class contains most of the particles. However, in the second classification the first two classes are quite similar regarding the number of particles. In order to have a consensus of these three results, we execute the protocol `xmipp3-consensus classes 3D` that compares several sets of 3D classes and return the intersection of the input classes (Fig. 44).

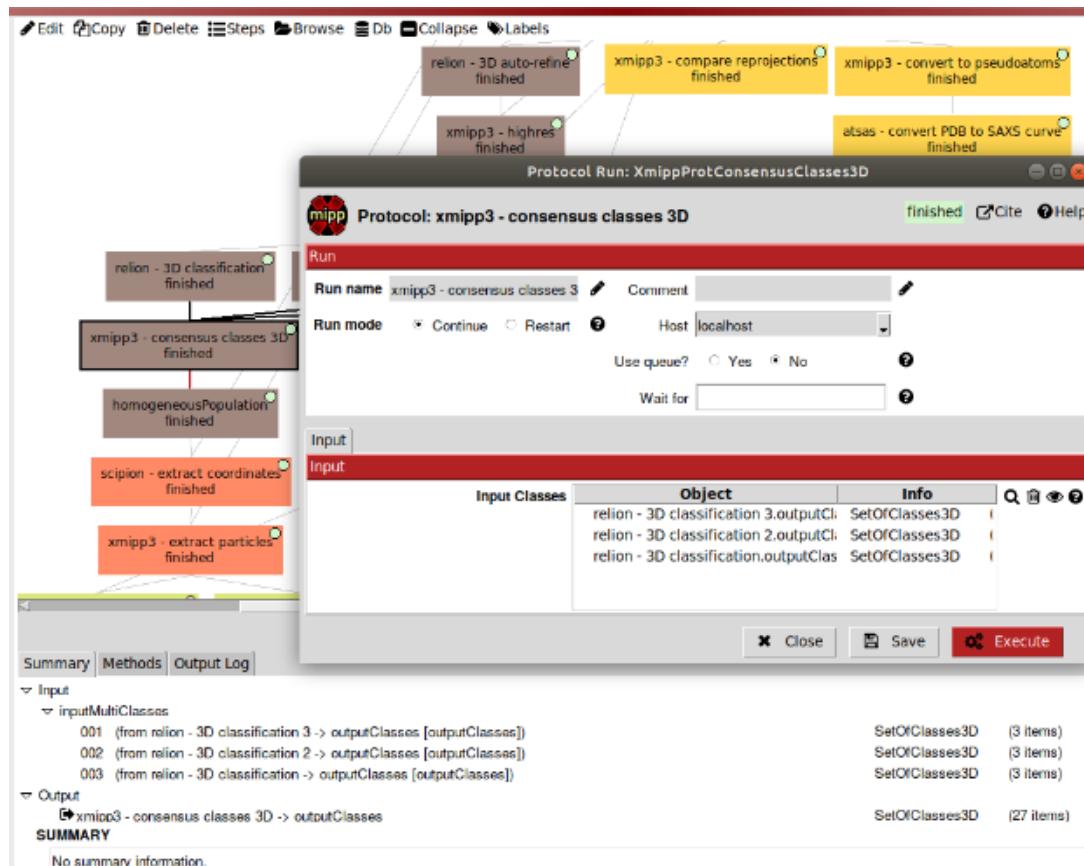


Figure 44: Filling in the params of the protocol **xmipp3-consensus classes 3D**.

By pressing **Analysis Results** you can visualize the 27 intersection 3D classes with the number of particles assigned to each one. Two of these classes derive from about 2,000 particles (aprox. 75% of total particles), 4 classes from 199-466 particles, 8 classes from 11-81 particles, and the other 13 classes from less than 10 particles.

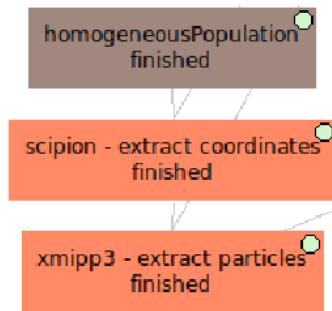


Figure 45: Extraction of re-classified particles (Orange color).

From this global population of reclassified particles (in brown in Fig. 45) we are going to extract again particle coordinates from the CTF-corrected micrographs, using the protocol `xmipp3-extract coordinates` (Fig. 46). This protocol allows to re-extract coordinates from particles with their original dimensions and visualize those particles in their locations on micrographs. The homogeneous population of particles previously obtained and the CTF consensus micrographs are the protocol input params.

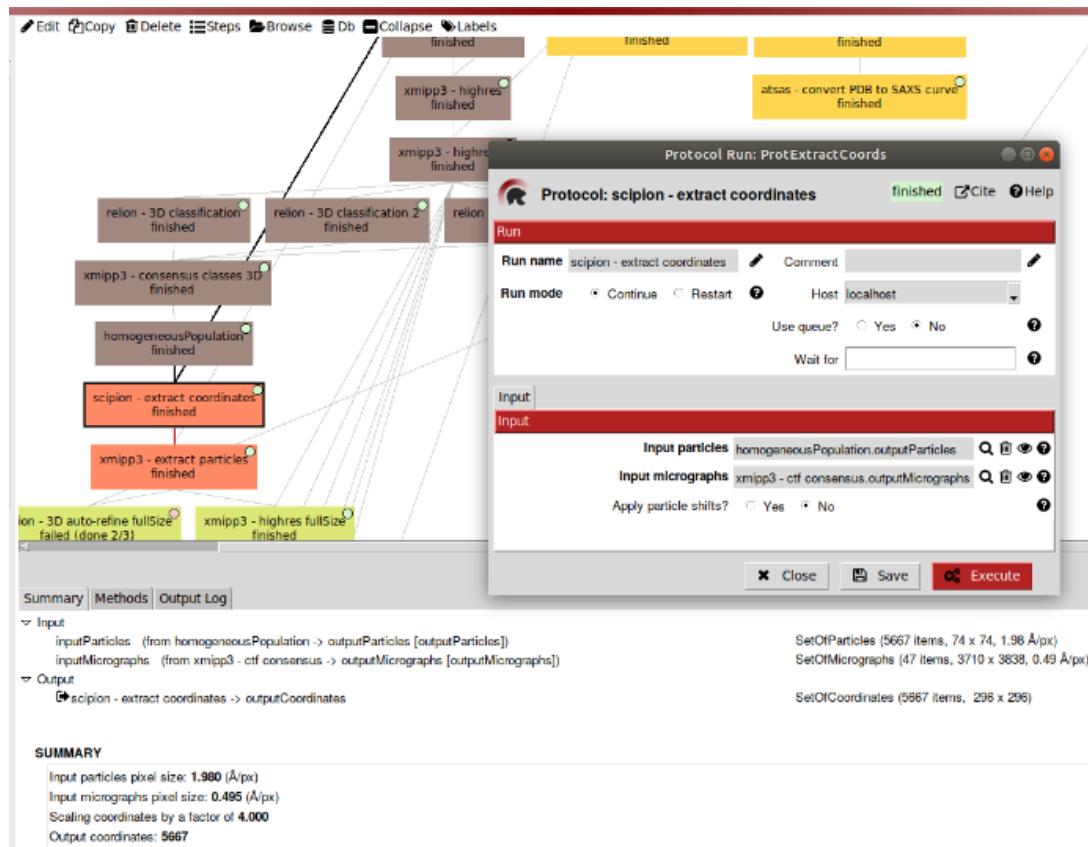


Figure 46: Completing the params of the protocol **xmipp3-extract coordinates**.

As output, **xmipp3-extract coordinates** generates the 5,667 particle coordinates scaled by a factor of 4.0. Press **Analyze Results** and observe the particles in each micrograph. These particles will be extracted with the above used protocol **xmipp3-extract particles** (Fig. 47). The particle coordinates previously extracted and CTF-consensus micrographs will be inputs of the protocol, as well as the Downampling factor of 1.0 and the Particle box size (px) of 450.

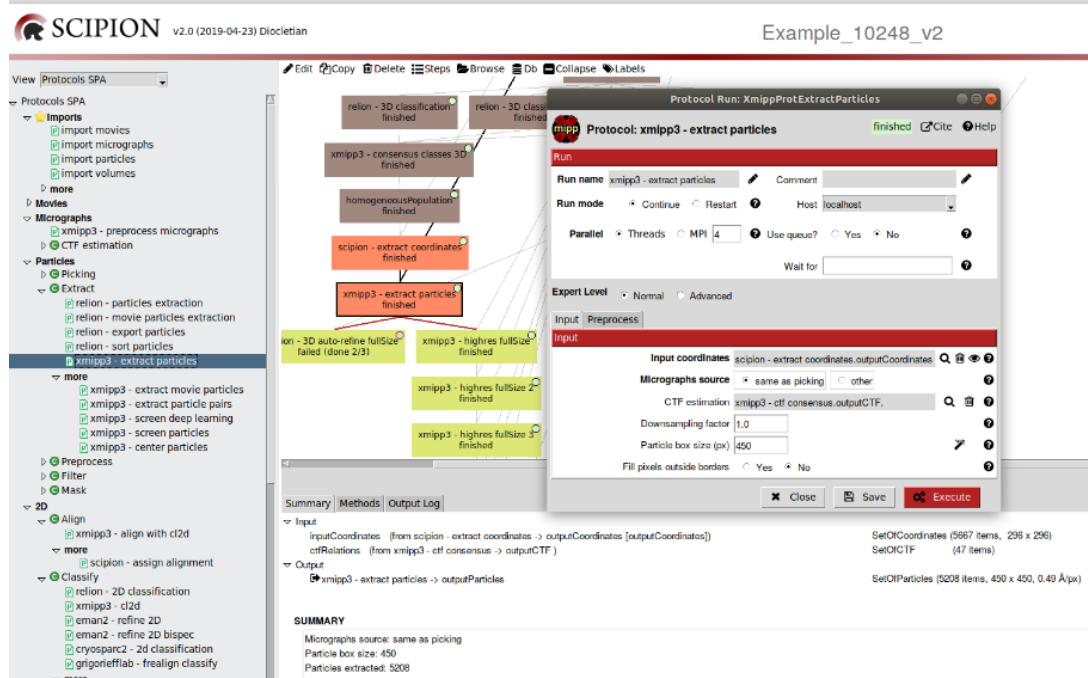


Figure 47: Filling in the params of the protocol `xmipp3-extract particles`.

The output of the protocol includes 5,208 particles (459 less than in the input) with the selected size and the starting sampling rate. Press **Analyze Results** to observe the table with the new set of particles.

Final refinement iterations with *Xmipp* highres

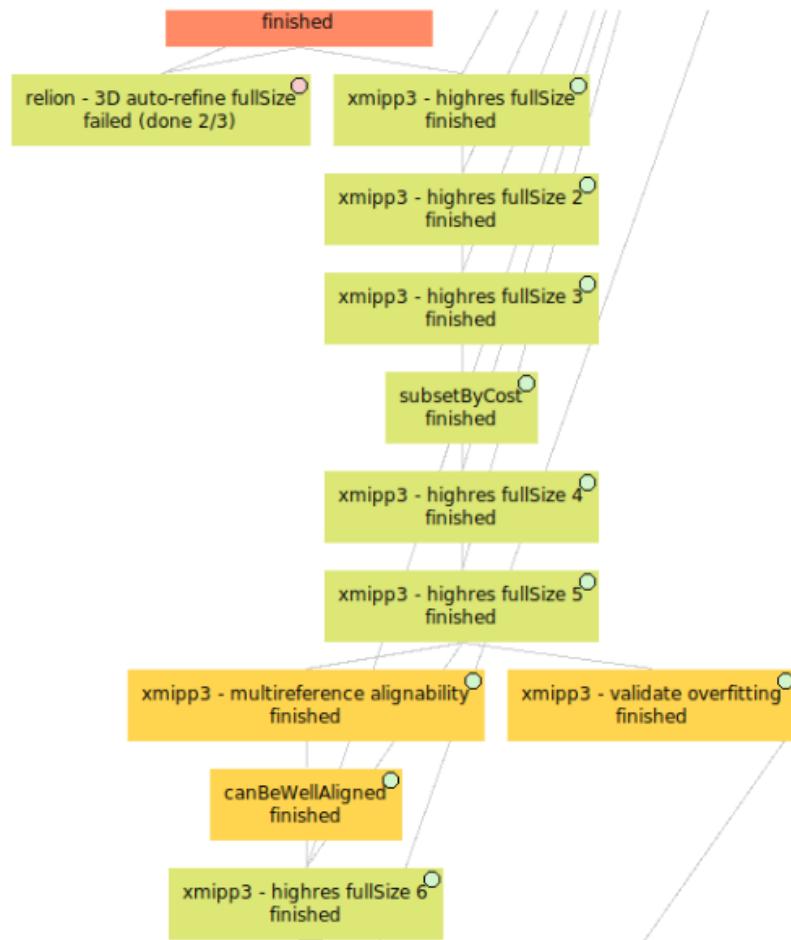


Figure 48: Map final refinement (Light green color).

From now ahead several steps of refinement will be accomplished using the above mentioned protocol `xmipp3-highres`. The input of the first round of refinement includes the particles extracted with the previous protocol and the volume generated by the same algorithm before performing the 3D classification step (Fig. 49). In the **Angular assignment** tap, we select **Global** for the **Image alignment** param, 1 as **Number of iterations** and 3 as **Max. Target Resolution**.

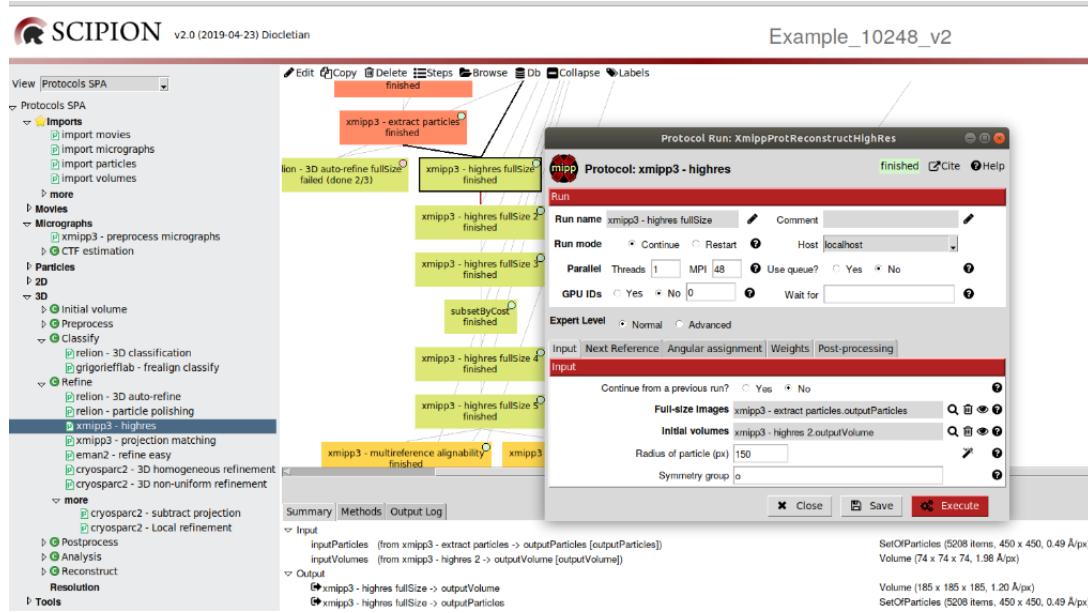


Figure 49: *Xmipp* `highres` map global refinement (Iteration 1).

The output resampled volume generated can be seen by pressing **Analyze results**, as well as the particles from which it derives.

The second run of refinement continues from the previous one, as it is selected in the **Input** tap. The particles derived from the first round of refinement are included as input params. The same params have been selected in the **Angular assignment** tap (Fig. 50).

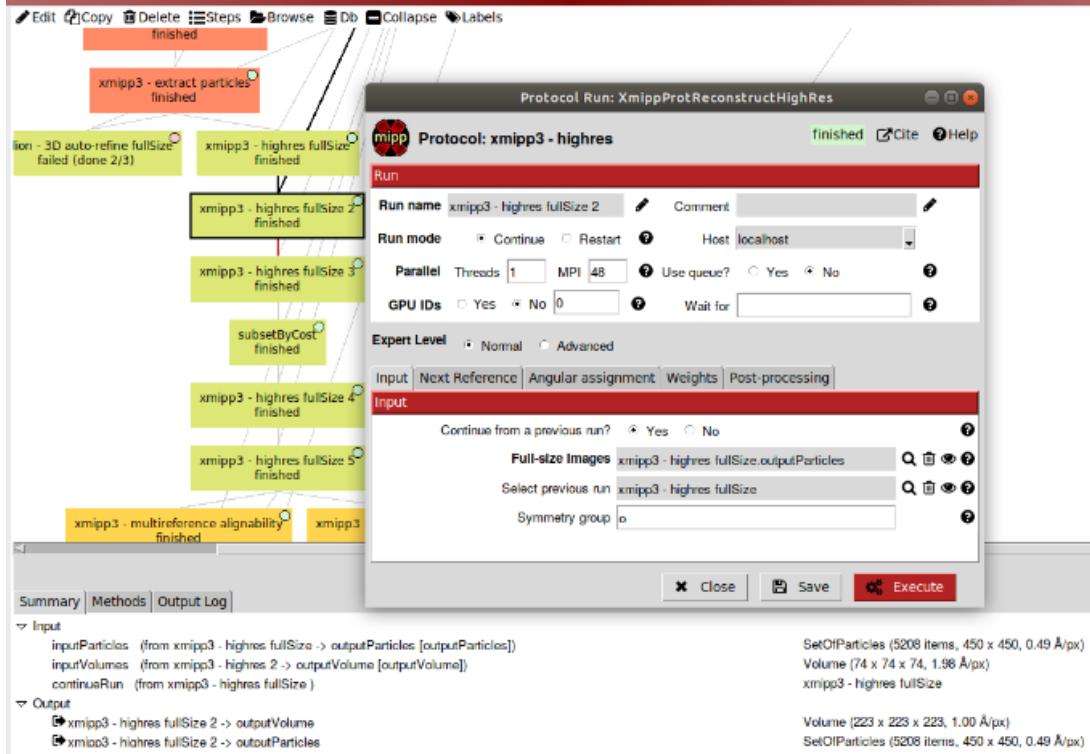


Figure 50: *Xmipp* highres map global refinement (Iteration 2).

A new output resampled volume has been obtained that move from 1.20 Å/px to 1.00 Å/px).

Once we have finished the global refinement, we continue with the local refinement in the third round of refinement (Fig. 51). Again, we use the particles derived from the previous iteration. In this case, we select Local for the Image alignment param and 2.5 for the Max. Target Resolution of tap Angular assignment. Shifts and angles will be also optimized.

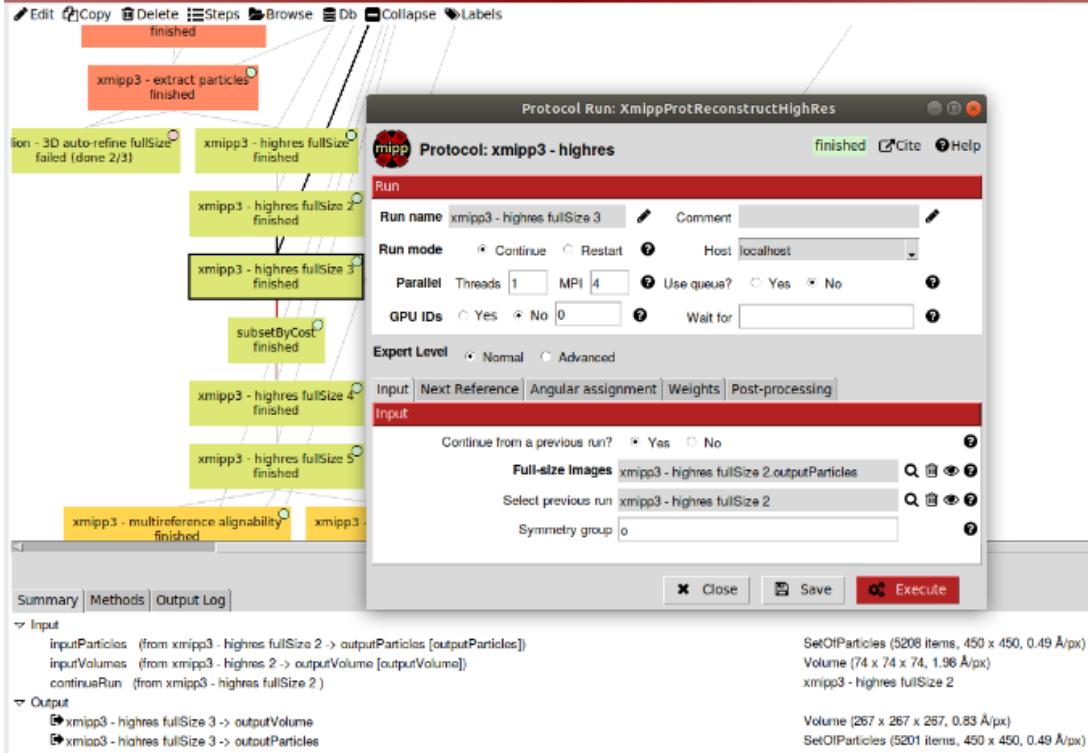


Figure 51: *Xmipp highres* map local refinement (Iteration 3).

The output resampled volume obtained ($0.83 \text{ \AA}/\text{px}$) derives from a set of particles slightly smaller (5,201). The table of particles can be also observed by pressing **Analyze results**. We can select particles according to the value of the `_xmipp_cost` param. Choosing values higher than 0.15, a total of 696 particles (13.4%) of the input set has been removed. The remaining 4,505 particles will be used as inputs of the fourth round of refinement (Fig. 52). Params from **Angular assignment** tap will remain unchanged except the Optimization ones. In addition to **shifts** and **angles**, **scale** and **defocus** will be also optimized.

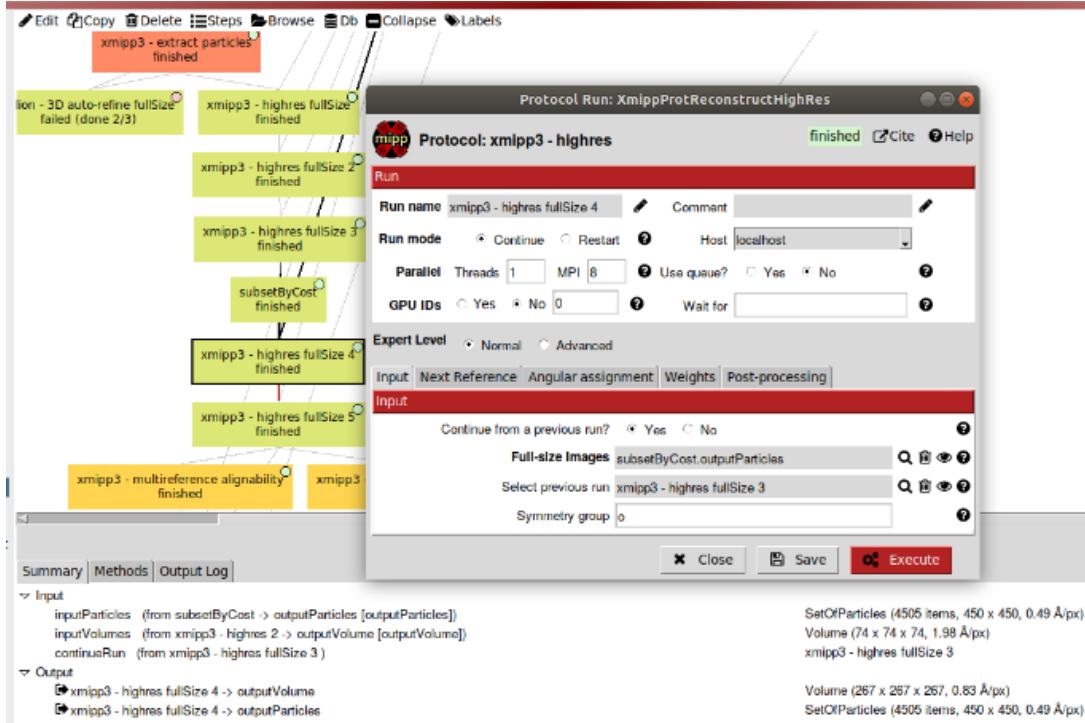


Figure 52: *Xmipp* highres map local refinement (Iteration 4).

The new map, based on the last set of particles, appears in the output. These particles are included in the input of the fifth round of local refinement (Fig. 53). This time, we reduce the value of the `Max. Target Resolution` param to 2.25 and set to Yes all the Optimization params.

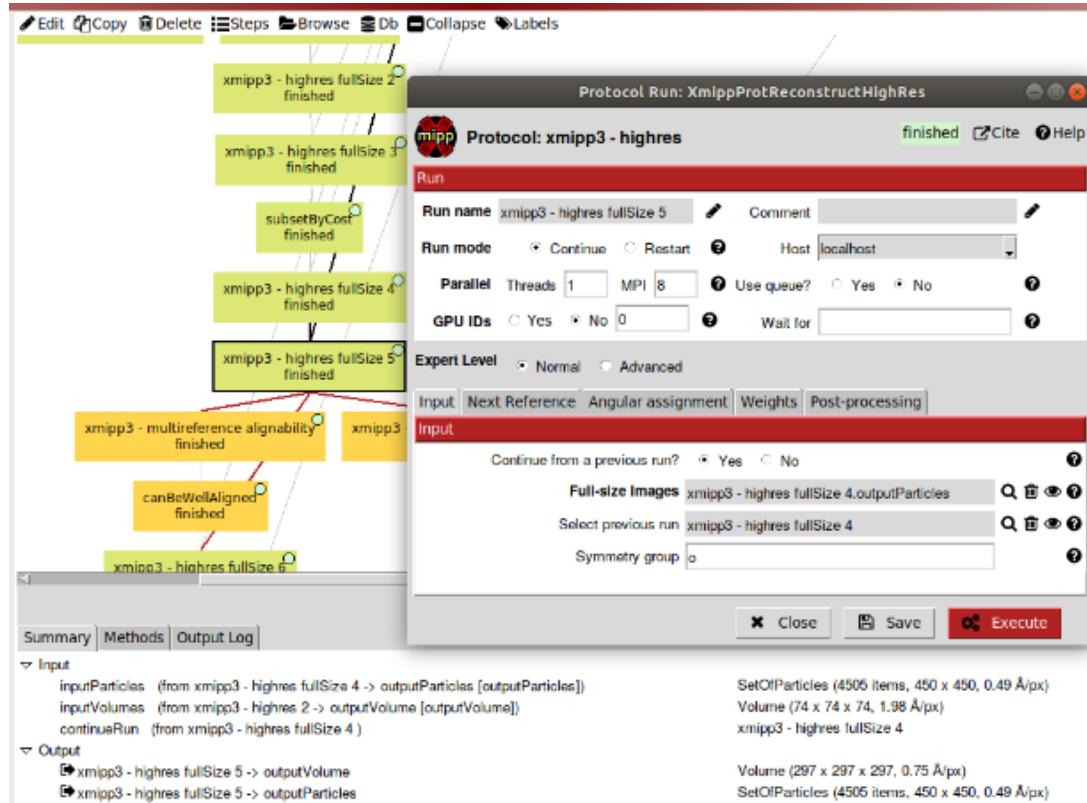


Figure 53: *Xmipp highres* map local refinement (Iteration 5).

Before continuing with the sixth round of refinement we are going to assess the output resampled map ($0.75\text{\AA}/\text{px}$) regarding soft alignability and overfitting of particles and 3D map. Two protocols are going to be independently executed: `xmipp3-multireference alignability` (Fig. 54) and `xmipp3-validate overfitting` (Fig. 55). The input of both protocols requires map and particles generated in the last refinement iteration.

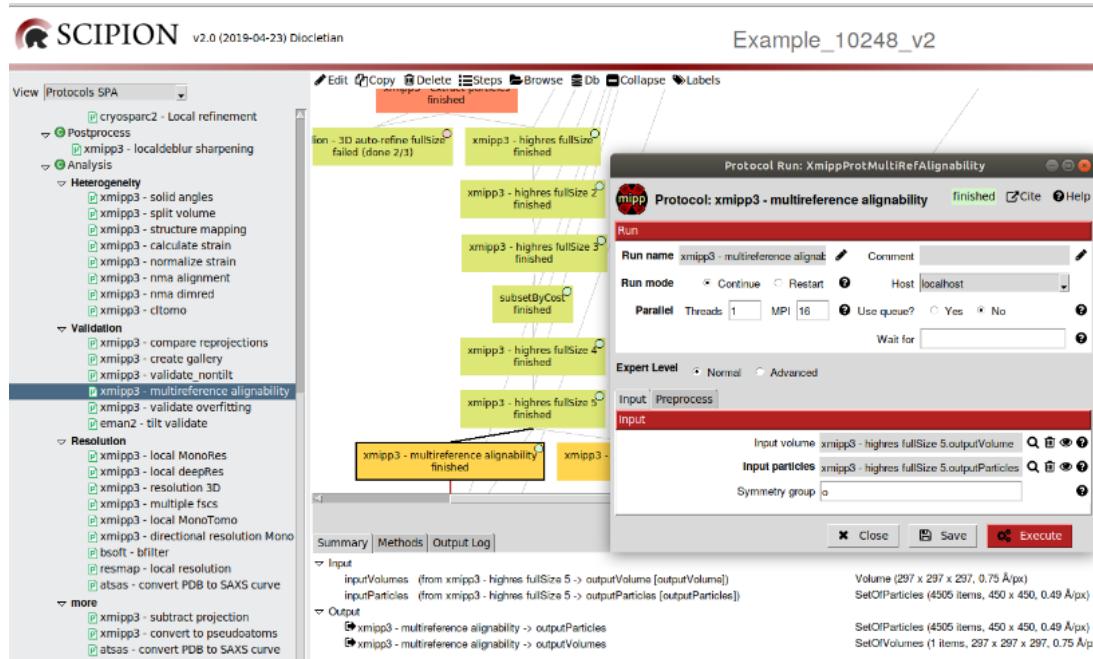


Figure 54: Completing the form of the protocol **xmipp3-multireference alignability**.

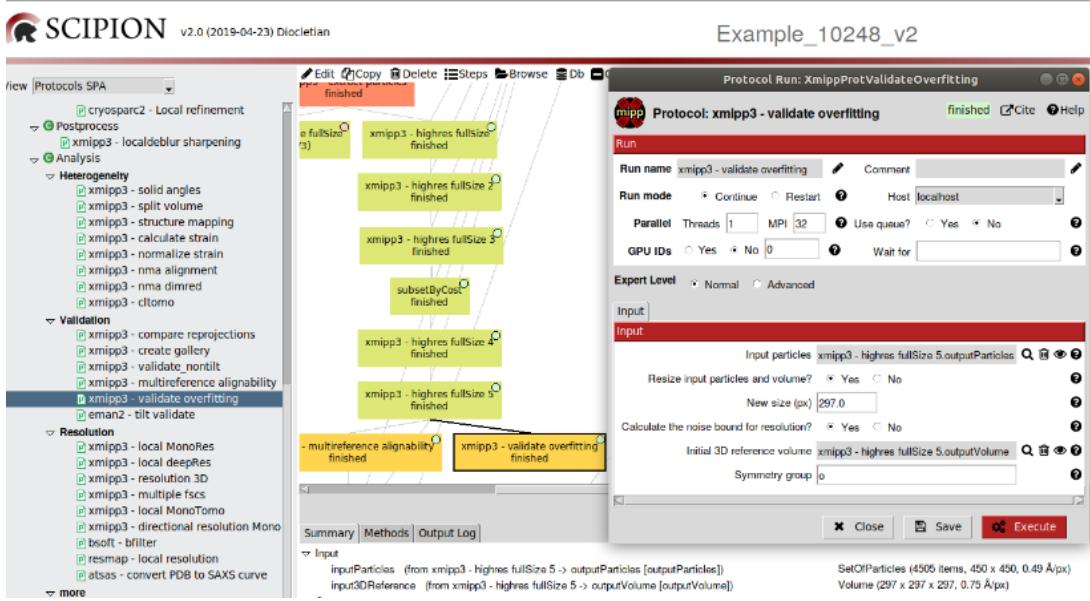


Figure 55: Filling in the form of the protocol `xmipp3-validate overfitting`.

The output values of particle alignment, precision and accuracy, generated by `xmipp3-multireference alignability` (press `Analyze Results` to check table columns `_xmipp_scoreAlignabilityAccuracy` and `_xmipp_scoreAlignabilityPrecision`) allow us to discard particles with worse alignment. In this case, 1,020 particles (22.6% of the total input) are rejected. In order to improve the refined map resolution, the remaining 3,485 particles will be used to perform the sixth local refinement iteration with *Xmipp highres* algorithm. The protocol params will remain unchanged compared with the fifth iteration.

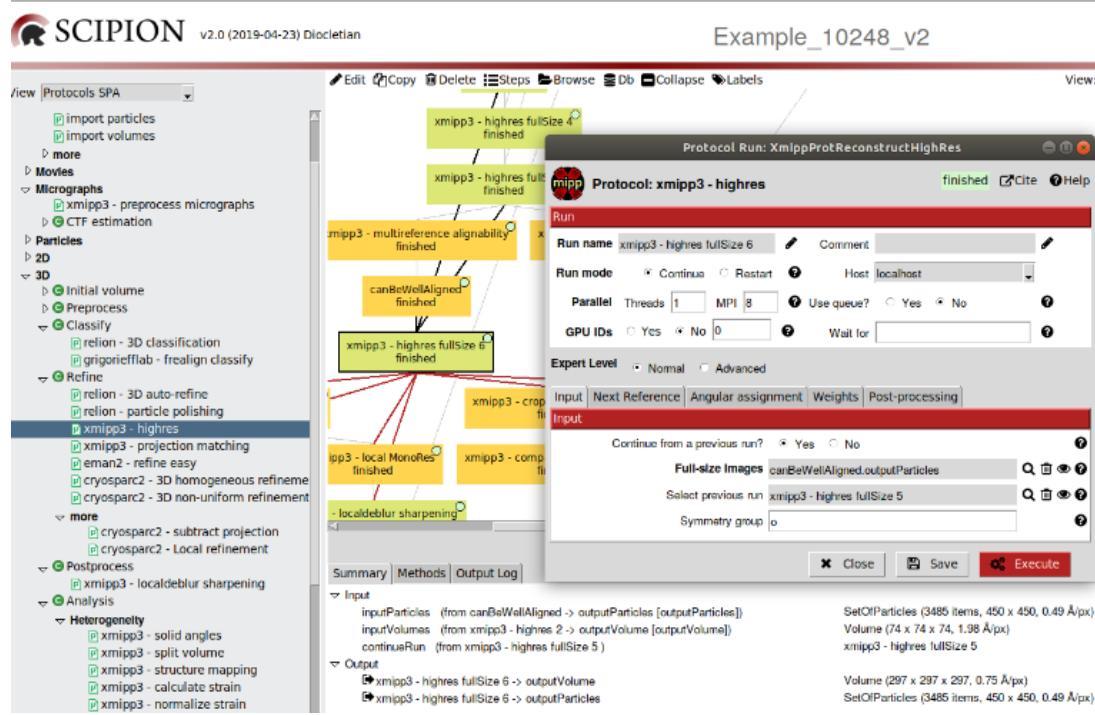


Figure 56: *Xmipp highres* map local refinement (Iteration 6).

The sampling of the output map is the same than in the previous iteration, despite the selection of best aligned particles. We have thus achieved convergence and we can compute the global or local resolution. The local resolution can be calculated with the protocol `xmipp3-local MonoRes` (Vilas et al., 2018). To have an overview of protocol and function of MonoRes see our *Scipion* tutorial in Model Building (download from https://github.com/I2PC/scipion/wiki/tutorials/tutorial_model_building_basic.pdf).

References

de la Rosa-Trevin, J. M., Quintana, A., Del Cano, L., Zaldivar, A., Foche, I., Gutierrez, J., Gomez-Blanco, J., Burguet-Castell, J., Cuenca-Alba, J., Abrishami, V., Vargas, J., Oton, J., Sharov, G., Vilas, J. L., Navas, J., Conesa, P., Kazemi, M., Marabini, R., Sorzano, C. O., Carazo, J. M., 07 2016. Scipion: A software frame-

- work toward integration, reproducibility and validation in 3D electron microscopy. *J. Struct. Biol.* 195 (1), 93–99.
- Hamaguchi, T., Maki-Yonekura, S., Naitow, H., Matsuura, Y., Ishikawa, T., Yonekura, K., Jul 2019. A new cryo-EM system for single particle analysis. *J. Struct. Biol.* 207 (1), 40–48.
- Punjani, A., Brubaker, M. A., Fleet, D. J., 2016. Building proteins in a day: Efficient 3d molecular structure estimation with electron cryomicroscopy. *IEEE transactions on pattern analysis and machine intelligence* 39 (4), 706–718.
- Punjani, A., Rubinstein, J. L., Fleet, D. J., Brubaker, M. A., 2017. cryosparc: algorithms for rapid unsupervised cryo-em structure determination. *Nature methods* 14 (3), 290.
- Rohou, A., Grigorieff, N., Nov 2015. CTFFIND4: Fast and accurate defocus estimation from electron micrographs. *J. Struct. Biol.* 192 (2), 216–221.
- Sorzano, C., Vargas, J., de la Rosa-Trevín, J., Jiménez, A., Maluenda, D., Melero, R., Martínez, M., Ramírez-Aportela, E., Conesa, P., Vilas, J., et al., 2018. A new algorithm for high-resolution reconstruction of single particles by electron microscopy. *Journal of structural biology* 204 (2), 329–337.
- Sorzano, C. O., de la Rosa Trevín, J., Otón, J., Vega, J., Cuenca, J., Zaldívar-Peraza, A., Gómez-Blanco, J., Vargas, J., Quintana, A., Marabini, R., et al., 2013. Semi-automatic, high-throughput, high-resolution protocol for three-dimensional reconstruction of single particles in electron microscopy. In: *Nanoimaging*. Springer, pp. 171–193.
- Svergun, D., Barberato, C., Koch, M. H., 1995. Crysolv-a program to evaluate x-ray solution scattering of biological macromolecules from atomic coordinates. *Journal of applied crystallography* 28 (6), 768–773.
- Vilas, J. L., Gómez-Blanco, J., Conesa, P., Melero, R., de la Rosa-Trevín, J. M., Otón, J., Cuenca, J., Marabini, R., Carazo, J. M., Vargas, J., et al., 2018.

Monores: automatic and accurate estimation of local resolution for electron microscopy maps. *Structure* 26 (2), 337–344.

Wagner, T., Merino, F., Stabrin, M., Moriya, T., Antoni, C., Apelbaum, A., Hagel, P., Sitsel, O., Raisch, T., Prumbaum, D., et al., 2019. Sphire-cryolo is a fast and accurate fully automated particle picker for cryo-em. *Communications Biology* 2 (1), 218.

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

Zivanov, J., Nakane, T., Forsberg, B. O., Kimanis, D., Hagen, W. J., Lindahl, E., Scheres, S. H., 2018. New tools for automated high-resolution cryo-em structure determination in relion-3. *Elife* 7, e42166.