

SWARM OPTIMIZATION AS A CONSENSUS TECHNIQUE FOR ELECTRON MICROSCOPY INITIAL VOLUME

C. O. S. SORZANO, J. L. VILAS, A. JIMENEZ-MORENO,
J. MOTA, T. MAJTNER, D. MALUENDA, M. MARTINEZ,
R. SANCHEZ, J. SEGURA, J. OTON, R. MELERO, L. DEL CANO,
P. CONESA, J. GOMEZ-BLANCO, Y. RANCEL, R. MARABINI,
J. M. CARAZO, J. VARGAS, AND R. MARABINI

ABSTRACT. Single Particle Analysis by Electron Microscopy aims at producing a three-dimensional model of a biological macromolecule using projection images acquired with an electron microscope. The task boils down to solving the inverse problem of estimating the three-dimensional structure from thousands of two-dimensional projections of it. The reconstruction process is iterative and needs an initial guess of the structure that can be obtained by several methods from the acquired data itself. The algorithm presented in this article eliminates the need of manually choosing one of these low quality volumes. Instead, it considers the whole population of initial volumes along with the acquired data and allows the whole population to evolve according to the dynamics given by swarm optimization. We show that this strategy successfully finds good initial estimates without the need of user intervention.

1. INTRODUCTION

The recent 2017 Chemistry Nobel Prize to Single Particle Analysis by Electron Microscopy has shown the maturity of Electron Microscopy as a structural technique for elucidating the three-dimensional (3D) structure of biological macromolecules [12,37]. The objective is to provide key information to better understand the biological mechanisms behind the physiological functions of these molecules. Electron microscopy acquires thousands of two-dimensional (2D) projection images from the macromolecule under study (Fig. 1, top). The goal is to produce a 3D model compatible with these measurements (Fig. 1, bottom, left), where each image is assigned to a projection direction with respect to this 3D model. As an intermediate step, similar projection images (supposed to come from a similar projection direction) are averaged in order to increase the Signal-to-Noise Ratio (Fig. 1, middle). These averages help to better understand the experimentally acquired images, help to remove contaminations and incorrectly selected particles, and are normally used in the construction of an initial model (Fig. 1, bottom, right).

2010 *Mathematics Subject Classification.* 68U10, 94A08, 90C26.

Key words and phrases. Electron microscopy, single particle analysis, initial volume, image processing, optimization, swarm consensus.

The initial volume is used to assign the projection directions to the experimental projections. Subsequently, these projections, along with their corresponding orientations, are used to update the reference volume. This process is iterated until convergence resulting in a high resolution structure like the one shown in the bottom left of Fig. 1. The correct choice of the initial volume is crucial because most refinement algorithms are greedy and they go in the direction of the closest local minimum [28]. Actually, there is well-described problem in the field called “Einstein-from-noise” [14] by which it is shown that pure noise projections (no macromolecular structure in them), aligned with a reference volume, result in the same reference volume. An alternative consequence of this problem is that by using an increasing number of pure noise particles, we may artificially increase the resolution of the reconstruction [15]. Severe consequences have been derived from the use of an incorrect initial volume [14, 17, 19, 33, 35] whose results are a totally incorrect 3D reconstruction of the macromolecule being studied.

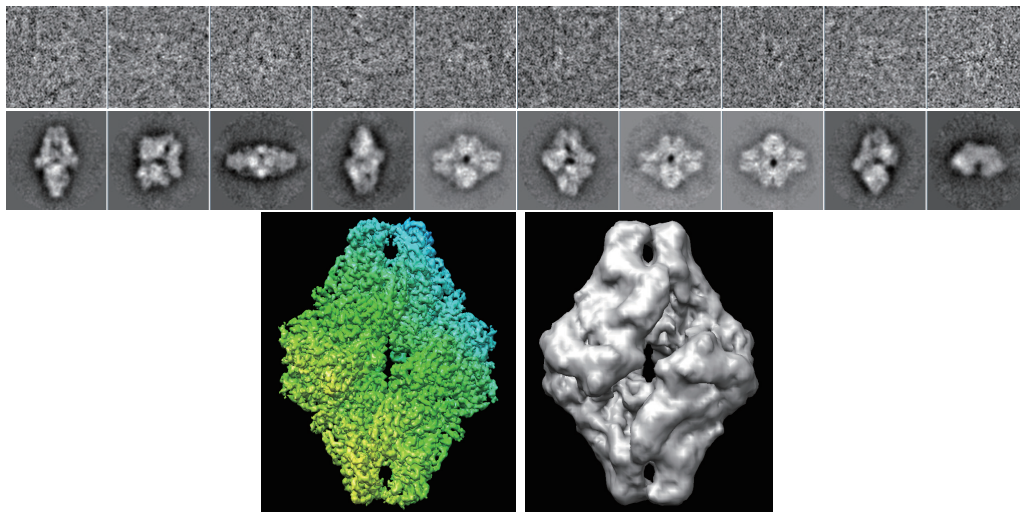


FIGURE 1. Top: Representative projections of the macromolecule under study. Middle: Representative 2D class averages of the projections. Bottom: 3D reconstruction from the dataset whose representative projections are shown on top (left) and its initial volume as calculated by swarm optimization (right).

Many algorithms have recently been proposed to computationally solve the initial volume problem [?, 5, 6, 10, 11, 21, 23–26, 30, 36, 38]. Most methods rely either on random angular assignments to the class averages of raw experimental projections; or on identifying common lines in the Fourier space. In any case, both families are rather error prone and algorithms normally produce a bunch of initial volume

candidates. The user is still left with the decision to choose a correct 3D structure to start the iterative process. This task is not easy due to the low resolution of the initial volumes, the fact that the true, underlying structure is unknown, and the fact that initial volume algorithms produce a non-negligible proportion of incorrect structures. Although there are some tools to help the user in this decision (like comparing the class averages with reprojections from the initial volume), there is a significant chance of starting from an incorrect structure.

In this article we introduce a method based on particle swarm optimization that relieves the user from having to take this decision. The method starts by estimating a collection of initial volumes using the standard approaches on class averages. Then, the whole group is considered as a population of candidates that evolve according to the swarm optimizer dynamics. We have empirically observed that as long as there is one volume in the group that can converge to the right structure, the whole population will also converge. In a way, the method can be seen as constructing a consensus between the different initial volume candidates. The method uses random subsets from the whole set of experimental projections to evolve the population. In this way, it is half-way between the initial volume algorithm (which uses the class averages) and the full refinement algorithms (which use all experimental projections).

2. METHODS

Particle swarm optimization. Let us consider the unconstrained optimization problem

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{arg\,min}} f(\mathbf{x})$$

Particle Swarm Optimization [9] (PSO) is an optimization technique that does not require the gradient of the objective function, f . Although it is not guaranteed to converge to the global optimum, it normally provides a much better exploration of the search space than local optimizers. PSO was originally designed to mimic the behavior of a flock of birds looking for food. At iteration k , each bird in the flock is flying and it has its own position, $\mathbf{x}_i^{(k)}$, and speed, $\mathbf{v}_i^{(k)}$. It evaluates the amount of food at its position $f_i^{(k)} = f(\mathbf{x}_i^{(k)})$ and knows its own best position along time, f_i^{best} (that occurred at \mathbf{x}_i^{best}), and the best position ever found by any of the rest of birds, \mathbf{x}^{best} . Then, it updates its velocity and position according to the iterative equations

$$(2.1) \quad \begin{aligned} \mathbf{v}_i^{(k+1)} &= \mathbf{v}_i^{(k)} + c_1 u_1 (\mathbf{x}_i^{best} - \mathbf{x}_i^{(k)}) + c_2 u_2 (\mathbf{x}^{best} - \mathbf{x}_i^{(k)}) \\ \mathbf{x}_i^{(k+1)} &= \mathbf{x}_i^{(k)} + \mathbf{v}_i^{(k+1)} \end{aligned}$$

c_1 and c_2 are acceleration constants typically fixed at $c_1 = c_2 = 2$ [9], and u_1 and u_2 are random variables uniformly distributed in the (0,1) range (for each iteration and each volume these random numbers are different).

In our case the $\mathbf{x}_i^{(k)}$ vectors are the initial volume candidates (let us assume there are I of them). Initially, they are started with the population provided by the initial volume algorithms and the initial speed vectors are set to $\mathbf{0}$. Due to the volumetric

interpretation of our bird positions, it is recommended that the input candidate volumes are aligned to a single reference so that they all lie in the same position, although this is not a strong requirement of the method. Internally, after updating each volume, it is always aligned to the average of all volumes so that volumes do not drift from a central position and they all stay in similar orientations.

In this article, we use a modified version of the original PSO in which the updated volume undergoes some transformation $R_i^{(k)}$ that considers the next candidate volume as well as the experimentally acquired data.

$$(2.2) \quad \mathbf{x}_i^{(k+1)} = R_i^{(k+1)} \left\{ \mathbf{x}_i^{(k)} + \mathbf{v}_i^{(k+1)} \right\}$$

where $\mathbf{x}_i^{(k)}$ is the current estimate of the initial volume, and $\mathbf{v}_i^{(k+1)}$ has been defined in Eq. (2.1). Inspired by Stochastic Gradient Descent [22], the operator $R_i^{(k+1)}$ is a reconstruction operator that is different for each volume and each iteration, which tends to prevent overfitting. The full dataset contains N_{img} images (typically in the order of several tens of thousands). Full refinement algorithms make use of all of them to refine the current estimate of the 3D model. However, for an initial volume we do not need all of them. Actually, it might be better if different candidates update using a different subset of the data. Even better, if this subset of data changes over time so that a particular candidate cannot overfit a particular subset. In this way, the operator $R_i^{(k+1)}$ takes, for every candidate and every iteration, a subset of the experimental images of size N_R (in our examples we choose $N_R = 500$). It performs an alignment and reconstruction using the algorithm introduced in [30] strongly based on the concept of statistical significance. Then, it denoises the reconstructed volume using again the concept of statistical significance and as shown in the postprocessing section of [32]. Finally, it aligns the denoised volume with respect to the population average.

To calculate the goal function $f(\mathbf{x})$, we randomly select N_E images from the experimental dataset (in our examples below we set $N_E = 100$). We perform a 3D angular assignment to determine their projection direction and generate reprojections of the volume \mathbf{x} along the same projection direction. $f(\mathbf{x})$ is the average of the cross-correlation of the experimental images with the reprojections. Note that the images used for the evaluation of the volumes are different from the N_R images used for the reconstruction by $R_i^{(k+1)}$. This goal function is evaluated for each volume in the swarm, and the \mathbf{x}_i^{best} and \mathbf{x}^{best} volumes are updated after each iteration. The algorithm is stopped after a fixed number of iterations.

For clarification of the whole process, Algorithm 1 summarizes the steps of our proposed approach. Note that the evaluation of the volumes involves choosing a random subset of the experimental images, aligning them and averaging the cross-correlation of the best alignment. Similarly, the refinement of the proposed volume involves choosing a random subset of the experimental images, aligning them and performing a 3D reconstruction. For the reconstruction and alignment we use the algorithm described in [30].

Algorithm 1: Swarm optimization of initial volumes

Data: A list of N_v initial volumes, $\mathbf{x}_i^{(0)}$ ($i = 1, 2, \dots, N_v$).
Data: A list of N_{img} images, \mathbf{y}_j ($j = 1, 2, \dots, N_{img}$).
Data: Number of iterations K
Result: List of evolved volumes, $\mathbf{x}_i^{(K)}$
Result: Average of the evolved volumes \mathbf{x}
// Keep the best volume in each trajectory
Set $\mathbf{x}_i^{best} = \mathbf{x}_i^{(0)}$
// Keep the best volume in the population
Evaluate $f(\mathbf{x}_i^{(0)})$ and choose \mathbf{x}^{best}
// Average all volumes

$$\mathbf{x} = \frac{1}{N_v} \sum_{i=1}^{N_v} \mathbf{x}_i^{(0)}$$
// Let the volumes evolve
for $k = 0, 1, \dots, K - 1$ **do**
 // Update velocity

$$\mathbf{v}_i^{(k+1)} = \mathbf{v}_i^{(k)} + c_1 u_1 (\mathbf{x}_i^{best} - \mathbf{x}_i^{(k)}) + c_2 u_2 (\mathbf{x}^{best} - \mathbf{x}_i^{(k)})$$
 // Propose new volume

$$\tilde{\mathbf{x}}_i^{(k+1)} = \mathbf{x}_i^{(k)} + \mathbf{v}_i^{(k+1)}$$
 // Refine new proposed volumes with
 // random subset of images

$$\hat{\mathbf{x}}_i^{(k+1)} = R_i^{(k+1)} \left\{ \tilde{\mathbf{x}}_i^{(k+1)} \right\}$$
 // Denoise and align volume with respect to the average

$$\mathbf{x}_i^{(k+1)} = A \left\{ D \left\{ \hat{\mathbf{x}}_i^{(k+1)} \right\}, \mathbf{x} \right\}$$
 // Update best volumes in each trajectory
 // and in the population
 Evaluate all $\mathbf{x}_i^{(k+1)}$ and choose $\mathbf{x}^{best}, \mathbf{x}_i^{best}$
 // Average all volumes

$$\mathbf{x} = \frac{1}{N_v} \sum_{i=1}^{N_v} \mathbf{x}_i^{(k+1)}$$
end

Relationship to other optimization problems. Our algorithm is related to two different optimization problems. On one hand we have the problem related to the update in Eq. (2.1), which we will refer to as the swarm reconstruction. On the other hand, we have the problem related to the update in Eq. (2.2) that we will call the stochastic reconstruction.

Swarm reconstruction. In this section we will look for a functional whose gradient descent iteration resembles the swarm iterative step. Let us consider the 3D reconstruction problem

$$(2.3) \quad \mathbf{x}^* = \arg \min_{\mathbf{x}} \|\mathbf{y} - P\mathbf{x}\|^2 + \frac{1}{2} c_1 \|\mathbf{x}_i^{best} - \mathbf{x}\|^2 + \frac{1}{2} c_2 \|\mathbf{x}^{best} - \mathbf{x}\|^2$$

where \mathbf{y} is a vector with the experimental measurements, \mathbf{x} is the volume to be reconstructed, $\mathbf{x}_i^{(best)}$ and $\mathbf{x}^{(best)}$ are two fixed volumes (for the moment, let us assume that they are known at this point), c_1 and c_2 are constants that define the importance of the regularization terms, and P is a projection matrix that transforms the volume into a set of projections images (in Electron Microscopy, this projection matrix may also include the effect of the Contrast Transfer Function). The associated gradient descent (Landweber) iterative algorithm would be [31]:

$$(2.4) \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mu_k P^T (\mathbf{y} - P\mathbf{x}^{(k)}) + \mu_k c_1 (\mathbf{x}_i^{best} - \mathbf{x}^{(k)}) + \mu_k c_2 (\mathbf{x}^{best} - \mathbf{x}^{(k)})$$

where μ_k is the relaxation factor and controls the convergence speed.

If we compare Eq. (2.1) with Eq. (2.4), we see that they are very similar, although they have some notable differences:

- The pull towards $\mathbf{x}_i^{(best)}$ and $\mathbf{x}^{(best)}$ is performed stochastically in the particle swarm update (Eq. (2.1)) as if μ_k were a (0,1)-uniform random variable.
- The velocity term is self-updated in the particle swarm update (Eq. (2.1)) while there is not such a self-reference in the gradient descent iteration. However, this self update was proposed by [?] as a way of accelerating the convergence of gradient descent by adding momentum. The idea was that in order to minimize an error function $E(\mathbf{x})$, instead of the standard gradient descent iteration

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mu_k \nabla E(\mathbf{x}^{(k)})$$

we could consider the “momentum” of the trajectory (as if it were a heavy ball falling down the function landscape)

$$\begin{aligned} \mathbf{v}^{(k+1)} &= \beta_k \mathbf{v}^{(k)} + \nabla E(\mathbf{x}^{(k)}) \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \mu_k \mathbf{v}^{(k+1)} \end{aligned}$$

where β_k is some suitable sequence of relaxation factors.

- The velocity term in the particle swarm update (Eq. (2.1)) does not directly incorporate any information from the experimental dataset. The experimental information indirectly comes through the evaluation of the new candidate $\mathbf{x}^{(k+1)}$. However, Eq. (2.4) explicitly incorporates the data \mathbf{y} in a reconstruction iteration ($P^T(\mathbf{y} - P\mathbf{x}^{(k)})$).

Stochastic reconstruction. Let us consider now the J individual experimental images, \mathbf{y}_j . The alignment and reconstruction problem can be formulated as the problem of finding a volume, the angular orientation, and in-plane shifts of each image such that

$$(2.5) \quad \mathbf{x}^*, P_j^* = \arg \min_{\mathbf{x}, P_j} \frac{1}{J} \sum_j \|\mathbf{y}_j - P_j \mathbf{x}\|^2$$

The stochastic gradient descent [3] introduced the idea that to minimize an error function of the form

$$(2.6) \quad \mathbf{x}^* = \arg \min_{\mathbf{x}} E(\mathbf{x}) = \arg \min_{\mathbf{x}} \frac{1}{J} \sum_j E_j(\mathbf{x})$$

one could work with subsets of these individual error functions, $E_j(\mathbf{x}) = \|\mathbf{y}_j - P_j \mathbf{x}\|^2$, and approximate the gradient descent update

$$(2.7) \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mu_k \nabla E(\mathbf{x}^{(k)})$$

by an approximation to this gradient

$$(2.8) \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mu_k \frac{1}{|J_k|} \sum_{j \in J_k} \nabla E_j(\mathbf{x}^{(k)})$$

where J_k is a random subset of experimental images (at each iteration, k , this subset is different), and $|J_k|$ is its cardinality. Stochastic gradient descent (SGD) has proved to be very useful in optimizing functions with a large number of parameters (as in deep learning). Although, in theory, it is not guaranteed to converge to a global minimum, in practice, it typically converges to it, or at least to a “good-enough” minimum, usually much better than their deterministic counterparts.

The gradient descent optimization is usually addressed in Electron Microscopy by an iterative algorithm that decomposes the original problem into two subproblems. Given the current estimate of the 3D reconstruction, $\mathbf{x}^{(k)}$, we look for the angular orientation and in-plane shifts of each experimental image

$$(2.9) \quad P_j^{(k+1)*} = \arg \min_{P_j} \|\mathbf{y}_j - P_j \mathbf{x}^{(k)}\|^2$$

Then, we use these alignment parameters to refine the current volume estimate

$$(2.10) \quad \mathbf{x}^* = \arg \min_{\mathbf{x}} \sum_j \|\mathbf{y}_j - P_j^{(k+1)*} \mathbf{x}\|^2$$

Our algorithm uses this idea of random subsets (SGD) to try to avoid getting trapped into a local minimum. This issue is particularly important in the problem at hand, initial volume estimation, because it is well known the difficulty to get away from local minima in Electron Microscopy [14, 17, 19, 33, 35] and the quality of the final reconstruction extremely depends on the quality of its initial volume. However, we use the state-of-the-art algorithms to solve the alignment [30] and reconstruction [1] subproblems (Eqs. (2.9) and (2.10)).

Features of our initial volume algorithm. With this digression, we see that our initial volume algorithm alternates between a swarm reconstruction iteration and a stochastic reconstruction iteration. The idea is to exploit the properties of these two algorithms to tend to escape from local minima. At the same time we tried to accelerate the convergence of the swarm reconstruction by allowing the algorithm to access to an update that directly uses the experimental images (instead of an indirect use through the correlation coefficient which is a much weaker link to the underlying 3D reconstruction problem).

The random nature of the optimization (random update of the particle velocities as well as the random subsets in the stochastic reconstruction) helps to avoid overfitting a particular set of images. Additionally, the fact that we have as many particles in the swarm as input volumes in the original set of initial volumes help the algorithm to explore promising areas of the landscape of solutions (let us remind that this set of initial volumes were proposed by the current state-of-art algorithms

for the initial volume). We have observed that the algorithm converges to a good starting volume as long as there is at least one sufficiently good initial volume in the original input set. In practice, we have not seen the algorithm to fail producing a good initial volume, despite the fact that many (sometimes most) of the input volumes are not correct. In this way, we relieve the user from having to choose a good volume from the initial dataset.

The fact that this algorithm is working with random subsets of several hundreds of experimental images (instead of a few dozens of class averages) also helps to avoid the overfitting typically encountered by the standard initial volume algorithms.

3. RESULTS

In order to validate the usefulness of our method, we have tested it with two macromolecular structures under very different imaging conditions: 1) our first example is the β -galactosidase under cryo-EM conditions [2], this molecule is D_2 symmetric; 2) the PriL-CTD of the α -polymerase protein in negative staining [20] acquired with a Random Conical Tilt, this protein is asymmetric. Both cases pose important challenges to the algorithm. The first example has very low signal-to-noise ratio (SNR) and contrast, making the alignment of the experimental images difficult. The second example has higher SNR, but its asymmetry makes the landscape of solutions more complicated, especially with the lack of an initial volume. Random Conical Tilt is an experimental methodology able to estimate an initial volume exploiting the geometrical relationship between two tilted views of the same microscopy field [27]. However, the resulting volumes lacks a wide region of the volume content in Fourier space. We show that our method avoids the need to perform a Random Conical Tilt reconstruction and that the method is capable of producing a suitable initial volume without any *a priori* geometrical knowledge and without this missing region in Fourier space.

β -galactosidase. This dataset held the resolution record in 2015 with 2.2Å [2]. Fig. 1 shows some experimental images as well as some of the two-dimensional classes obtained from it using CL2D [?]. We collected 3,460 projection images from 15 electron micrographs (the specific details of the experimental acquisition can be seen at [4]). We corrected the phase flip introduced by the Contrast Transfer Function (CTF) at the level of micrograph before extracting the particles. These images were classified into 32 classes and these classes were input to the RANSAC initial volume algorithm [36]. We specifically manipulated the RANSAC parameters to make it fail as to put the swarm consensus algorithm in an unfavourable situation. Fig. 2 (top) shows the initial volumes produced by these algorithms. Note that none of them is a correct initial volume (in practice the proportion of incorrect volumes produced by RANSAC and EMAN [34] initial volume, two of the most popular algorithms to initiate the 3D reconstruction process may oscillate, depending on the dataset, between 20 and 90%). This population of 20 volumes (10 from RANSAC and 10 from EMAN) was fed into the swarm consensus algorithm presented in this paper. Fig. 2 (bottom) shows the final volumes produced by swarm consensus after 10 iterations. We updated the volumes each time with 500 experimental images, and evaluated their quality with another 100 experimental images. During the

reconstruction we enforced D2 symmetry for this molecule. The whole process took 4h in a laptop with 8 cores and 16GB of RAM. It can be seen that all of the original volumes, regardless of their original quality, converged to the correct structure. We performed a 3D reconstruction starting from the average of the evolved volumes and we were able to reach high resolution. Given the large number of reconstructions to perform (20) we did not reconstructed to high resolution all 20 faulty initial volumes. However, for a few of them we did and, as expected, the 3D refinement algorithm was not capable of escaping from this local minimum.

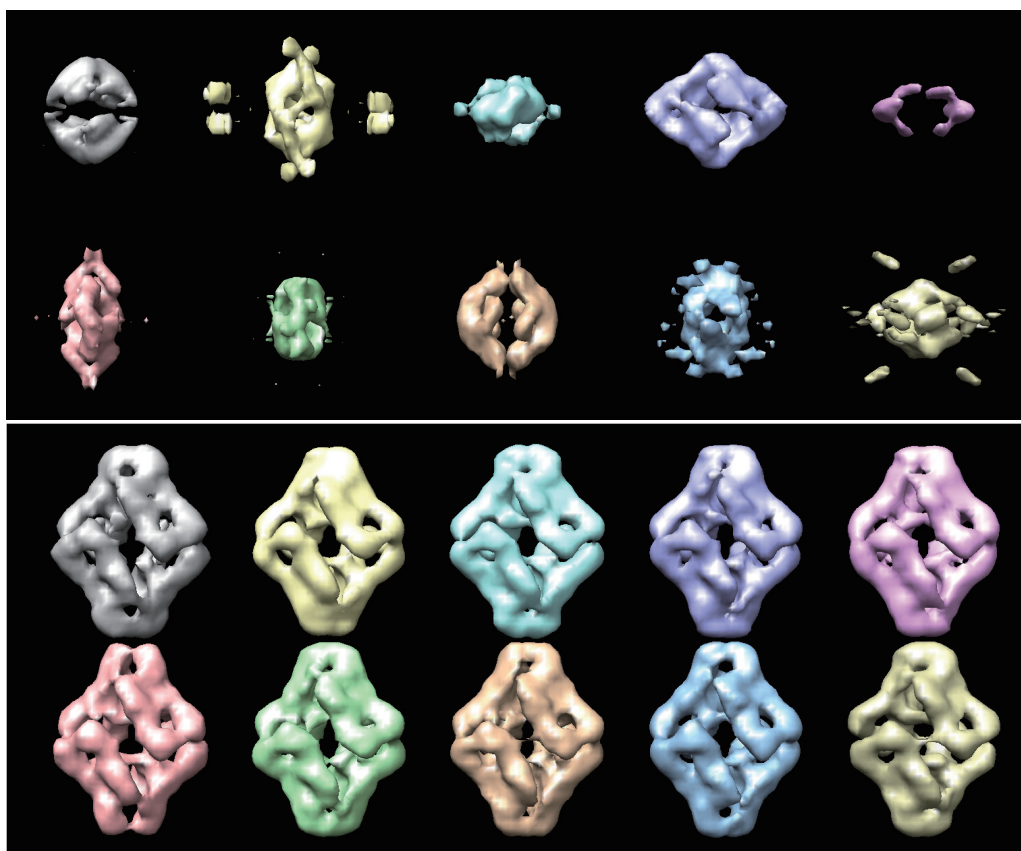


FIGURE 2. Top: Initial set of volumes of the β -galactosidase. Bottom: Final set of volumes produced by consensus swarm.

PriL-CTD of the α -polymerase. This dataset was used as an example to illustrate the Random Conical Tilt procedure [27] and the experimental acquisition details are described at [20]. In [27] it was shown that with as few as 13 projection pairs, an initial volume could be reconstructed (although it lacked large areas in Fourier space). The fact that the protein is asymmetric poses an important challenge to the angular and 3D reconstruction optimizer because the number of degrees of freedom is not simplified by symmetry. Fig. 3 shows some of the experimental projections as well as some of the 2D classes calculated from them. We collected 8,526 projection images from 12 electron micrographs (6 tilted and 6 untilted). We

corrected the phase flip introduced by the CTF at the level of micrograph before extracting the particles. They were classified into 32 classes using CL2D [?] and these classes were given to RANSAC obtaining the set of initial volumes shown in Fig. 4 (top). Note that none of these structures is correct although some of them show the two main domains of the central part of the protein. We processed this set of 20 initial volumes (10 from RANSAC and 10 from EMAN) with 10 iterations of the algorithm described in this paper. Obtaining the results shown in Fig. 4 (bottom), where all of the structures are valid initial volumes. As in the previous example, we updated the volumes with 500 experimental images, and evaluated their quality with another 100 experimental images. The whole process took 8h in a laptop with 8 cores and 16GB of RAM.

In Fig. 5 we show the evolution of the average cross correlation of reprojections of the volumes in the swarm with the random subsets of 100 experimental images (this is the objective function of the algorithm). We can see that the swarm stabilizes after relatively few iterations (about 8-9 iterations). For the β -galactosidase, the stabilization is even earlier (2-3 iterations, data not shown) thanks to the simplified landscape of solutions implied by the D2 symmetry of the molecule. This can be compared with the thousands of iterations reported in the Stochastic Gradient implementation in CryoSparc [21].

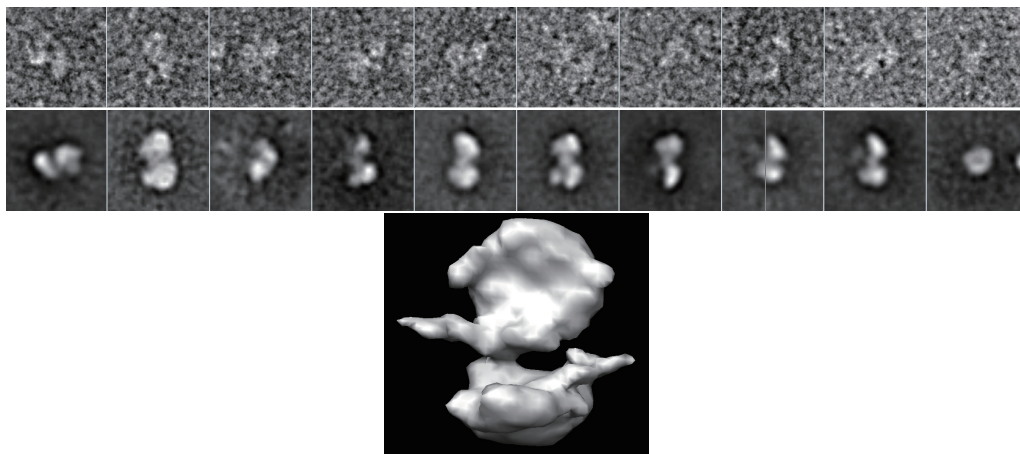


FIGURE 3. Top: Representative experimental images of PriL-CTD of the α -polymerase obtained by negative staining. Middle: Examples of 2D classes of PriL-CTD of the α -polymerase. Bottom: Random Conical Tilt volume (note the elongation along the missing cone direction).

CONCLUSIONS

In this article we have presented a new algorithm to combine the information from a set of initial volume candidates. With this approach, we relieve the user from having to choose from this set and we, thus, decrease the probability of producing an incorrect structure by an incorrect choice. Our algorithm is especially designed to

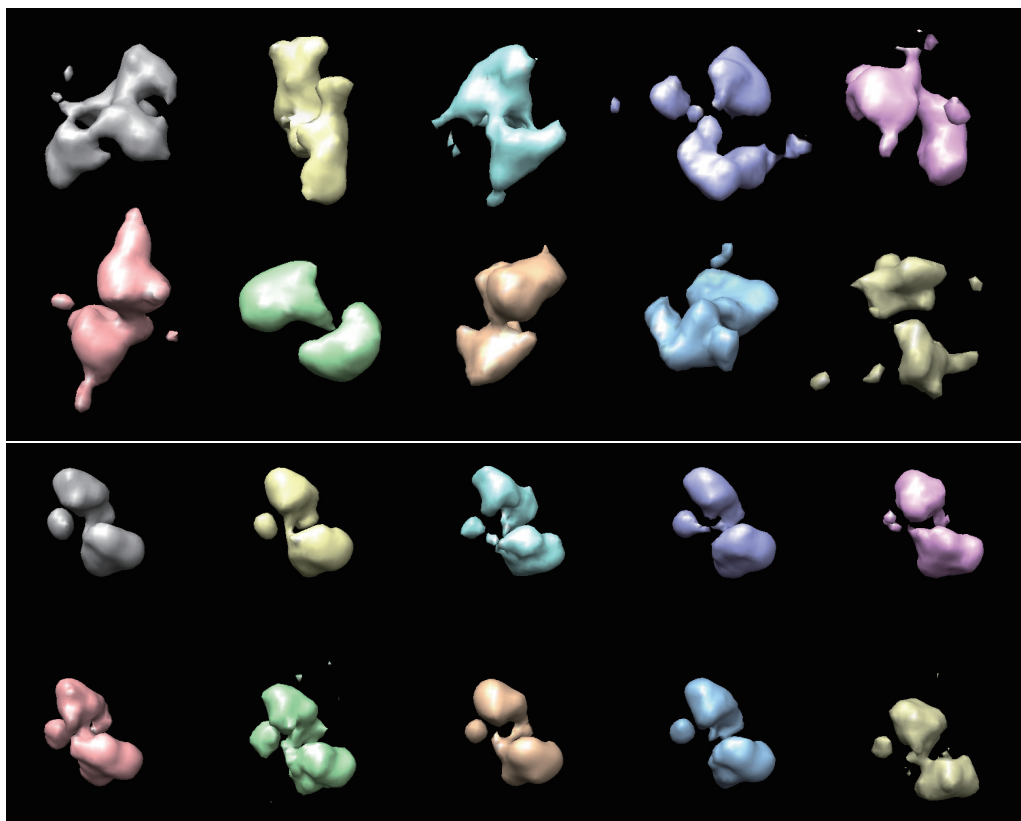


FIGURE 4. Top: Initial set of volumes of PriL-CTD of the α -polymerase. Bottom: Final set of volumes produced by consensus swarm.

avoid getting trapped into local minima, although there is no proof of convergence to the global minimum. The algorithm borrows ideas from particle swarm optimization (which, in its turn, is related to the momentum gradient descent) and stochastic gradient descent. We alternate between both types of iterations trying to avoid getting trapped into local minima as well as trying to accelerate the convergence to a reasonable initial volume to start the angular refinement of all the experimental images.

Two features are particularly interesting from the algorithm: 1) it explores a wide area of the landscape of solutions by taking an input set of initial volumes (these volumes have been proposed by sensible algorithms, although prone to commit errors); 2) it randomizes the gradient descent updates in two different ways (by random momentum coefficients as done by particle swarm optimization and by random subsets of the experimental images as done by the stochastic gradient descent). Both actions result in a rather robust algorithm which we have shown to converge to correct structures under rather challenging conditions.

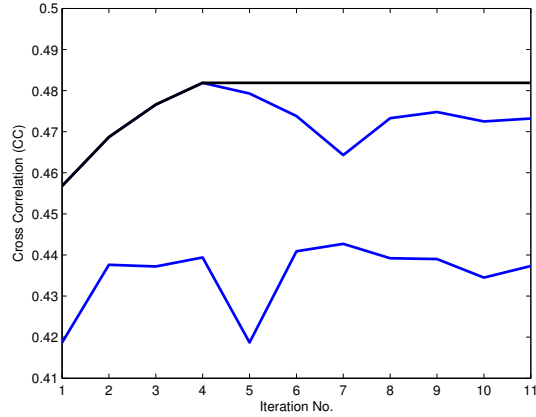


FIGURE 5. Evolution over iterations of the average cross correlation between reprojections of the swarm of volumes with a random subset of 100 experimental images. The blue lines represent the minimum and maximum average cross correlation in the swarm. The black line represents the best historical value (this is the average cross correlation for \mathbf{x}^{best}).

The algorithm has been implemented in the Xmipp package [7, 29] and it is freely available through the Scipion image processing framework [8] under the name **swarm consensus**.

ACKNOWLEDGEMENT

The authors would like to acknowledge economical support from: The Spanish Ministry of Economy and Competitiveness through Grants BIO2016-76400-R(AEI/FEDER, UE), the “Comunidad Autónoma de Madrid through Grant: S2017/BMD-3817, Instituto de Salud Carlos III, PT13/0001/0009, PT17/0009/0010, European Union (EU) and Horizon 2020 through grants CORBEL (INFRADEV-1-2014-1, Proposal: 654248), West-Life (EINFRA-2015-1, Proposal: 675858), Elixir - EXCELERATE (INFRADEV-3-2015, Proposal: 676559), iNEXT (INFRAIA-1-2014-2015, Proposal: 653706), EOSCpilot (INFRADEV-04-2016, Proposal: 739563), INSTRUCT - ULTRA (INFRADEV-03-2016-2017, Proposal: 731005). The authors acknowledge the support and the use of resources of Instruct-ERIC.”

REFERENCES

- [1] V. Abrishami, J. R. Bilbao-Castro, J. Vargas, R. Marabini, J. Carazo and C. O. S. Sorzano, *A fast iterative convolution weighting approach for gridding-based direct Fourier three-dimensional reconstruction with correction for the contrast transfer function*, *Ultramicroscopy* **157** (2015), 79–87.
- [2] A. Bartesaghi, A. Merk, S. Banerjee, D. Matthies, X. Wu, J. L. S. Milne and S. Subramaniam, *2.2 Å resolution cryo-EM structure of β -galactosidase in complex with a cell-permeant inhibitor*, *Science* 2015.
- [3] L. Bottou, *Large-scale machine learning with stochastic gradient descent*, in: Proc. COMPSTAT'2010, Springer, 2010, pp. 177–186.
- [4] S. Chen, G. McMullan, A. R. Faruqi, G. N. Murshudov, J. M. Short, S. H. W. Scheres and R. Henderson, *High-resolution noise substitution to measure overfitting and validate resolution in 3D structure determination by single particle electron cryomicroscopy.*, *Ultramicroscopy* **135** (2013), 24–35.
- [5] R. R. Coifman, Y. Shkolnisky, F. Sigworth and A. Singer, *Cryo-EM structure determination through eigenvectors of sparse matrices*, Dept. Computer Science, Univ. Yale, 2007.
- [6] R. R. Coifman, Y. Shkolnisky, F. Sigworth and A. Singer, *Reference free structure determination through digenvectors of center of mass operators*, *Appl. Comput. Harmon. Anal.* **28** (2010), 296–312.
- [7] J. M. de la Rosa-Trevín, J. Otón, R. Marabini, A. Zaldívar, J. Vargas, J. M. Carazo and C. O. S. Sorzano, *Xmipp 3.0: an improved software suite for image processing in electron microscopy*, *J. Structural Biology* **184** (2013), 321–328.
- [8] J. M. de la Rosa-Trevín, A. Quintana, L. Del Cano, A. Zaldívar, I. Foche, J. Gutiérrez, J. Gómez-Blanco, J. Burguet-Castell, J. Cuenca-Alba, V. Abrishami, J. Vargas, J. Otón, G. Sharov, J. L. Vilas, J. Navas, P. Conesa, M. Kazemi, R. Marabini, C. O. S. Sorzano and J. M. Carazo, *Scipion: A software framework toward integration, reproducibility and validation in 3D electron microscopy*, *J. Structural Biology* **195** (2016), 93–99.
- [9] R. C. Eberhart and Y. Shi, *Particle swarm optimization: developments, applications and resources*, in: Proc. 2001 Congress on Evolutionary Computation, vol. 1. IEEE, 2001, pp. 81–86.
- [10] D. Elmlund and H. Elmlund, *SIMPLE: Software for ab initio reconstruction of heterogeneous single-particles*, *J. Structural Biology* **180** (2012), 420–427.
- [11] D. Elmlund, R. Davis and H. Elmlund, *Ab initio structure determination from electron microscopic images of single molecules coexisting in different functional states*, **18** (2010), 777–786.
- [12] J. Frank, *Advances in the field of single-particle cryo-electron microscopy over the last decade*, *Nature protocols* **12** (2017), 209–212.
- [13] P. Geladi and B. R. Kowalski, *Partial least-squares regression: a tutorial*, *Analytica Chimica Acta* **185** (1986), 1–17.
- [14] R. Henderson, *Avoiding the pitfalls of single particle cryo-electron microscopy: Einstein from noise*, *Proc. Natl. Acad. Sci. USA* **110** (2013), 18037–18041.
- [15] B. Heymann, *Validation of 3DEM Reconstructions: The phantom in the noise*, *AIMS Biophysics* **2** (2015), 21–35.
- [16] M. J. Keiser, B. L. Roth, B. N. Armbruster, P. Ernsberger, J. J. Irwin and B. K. Shoichet, *Relating protein pharmacology by ligand chemistry*, *Nat Biotechnol* **24** (2007), 197–206.
- [17] Y. Mao, L. R. Castillo-Menendez and J. G. Sodroski, *Reply to Subramaniam, van Heel, and Henderson: Validity of the cryo-electron microscopy structures of the HIV-1 envelope glycoprotein complex*, *Proc. Natl. Acad. Sci. USA* **110** (2013), E4178–E4182.
- [18] C. F. Reboul, F. Bonnet, D. Elmlund and H. Elmlund, *A stochastic Hill climbing approach for simultaneous 2D alignment and clustering of cryogenic electron microscopy images.*, *Structure* **24** (2016), 988–996.
- [19] Y. Mao L. Wang, C. Gu and A. Herschhorn, A. Désormeaux, A. and A. Finzi, S. H. Xiang and J. G. Sodroski, *Molecular architecture of the uncleaved HIV-1 envelope glycoprotein trimer*, *Proc. Natl. Acad. Sci. USA* **110** (2013), 12438–12443.

- [20] R. Núñez-Ramírez, S. Klinge, L. Sauguet, R. Melero, M. A. Recuero-Checa, M. Kilkenny, R. L. Perera, B. García-Alvarez, R. J. Hall, E. Nogales, L. Pellegrini and O. Llorca, *Flexible tethering of primase and DNA Pol α in the eukaryotic primosome*, *Nucleic acids research* **39** (2011), 8187–8199.
- [21] A. Punjani, M. A. Brubaker and D. J. Fleet, *Building proteins in a day: Efficient 3D molecular structure estimation with electron cryomicroscopy*, *IEEE Trans. Pattern Anal. Machine Intelligence* **39** (2017), 706–718.
- [22] A. Punjani, J. L. Rubinstein, D. J. Fleet and M. A. Brubaker, *cryoSPARC: algorithms for rapid unsupervised cryo-EM structure determination*, *Nature methods* **14** (2017), 290–296.
- [23] Y. Shkolnisky and A. Singer, *Viewing direction estimation in cryo-EM using synchronization*, *SIAM J. imaging sci.* **5** (2012), 1088–1110.
- [24] A. Singer, R. R. Coifman, F. Sigworth, D. Chester and Y. Shkolnisky, *Detecting consistent common lines in cryo-EM by voting*, *J. Structural Biology* **169** (2010), 312–322.
- [25] A. Singer and Y. Shkolnisky, *Three-dimensional structure determination from common lines in cryo-EM by eigenvectors and semidefinite programming*, *SIAM J. Imaging Sci.* **4** (2011), 543–572.
- [26] A. Singer, Z. Zhao, Y. Shkolnisky and R. Hadani, *Viewing angle classification of cryo-electron microscopy images using eigenvectors*, *SIAM J. Imaging Sci.* **4** (2011), 723–759.
- [27] C. O. S. Sorzano, M. Alcorlo, J. M. de la Rosa-Trevín, R. Melero, I. Foche, A. Zaldívar-Peraza, L. del Cano, J. Vargas, V. Abrishami, J. Otón, R. Marabini and J. M. Carazo, *Cryo-EM and the elucidation of new macromolecular structures: random conical tilt revisited*, *Sci. Rep.* **5** (2015), 14290.
- [28] C. O. S. Sorzano, R. Marabini, A. Pascual-Montano, S. H. W. Scheres and J. M. Carazo, *Optimization problems in electron microscopy of single particles*, *Annals of Operations Research* **148** (2006), 133–165.
- [29] C. O. S. Sorzano, R. Marabini, J. Velázquez-Muriel, J. R. Bilbao-Castro, S. H. W. Scheres, J. M. Carazo and A. Pascual-Montano, *XMIPP: A new generation of an open-source image processing package for electron microscopy*, *J. Structural Biology* **148** (2004), 194–204.
- [30] C. O. S. Sorzano, J. Vargas, J. M. de la Rosa-Trevín, J. Otón, A. L. Álvarez-Cabrera, V. Abrishami, E. Sesmero, R. Marabini and J. M. Carazo, *A statistical approach to the initial volume problem in single particle analysis by electron microscopy*, *J. Structural Biology* **189** (2015), 213–219.
- [31] C. O. S. Sorzano, J. Vargas, J. Otón, J. L. Vilas, M. Kazemi, R. Melero, L. del Caño, J. Cuenca, P. Conesa, J. Gómez-Blanco, R. Marabini, J. M. Carazo, *A survey of the use of iterative reconstruction algorithms in electron microscopy*, *BioMed Research Intl.* **2017** (2017), 2017, 6482567.
- [32] C. O. S. Sorzano, J. Vargas, J. M. del Rosa-Trevín, A. Jiménez-Moreno, R. Melero, M. Martínez, P. Conesa, J. L. Vilas, R. Marabini and J. M. Carazo, *High-resolution reconstruction of single particles by electron microscopy*, *J. Structural Biology* (2018), submitted.
- [33] S. Subramaniam, *Structure of trimeric HIV-1 envelope glycoproteins*, *Proc. Natl. Acad. Sci. USA* **110** (2013), E4172–E4174.
- [34] G. Tang, L. Peng, P. R. Baldwin, D. S. Mann, W. Jiang, I. Rees and S. J. Ludtke, *EMAN2: an extensible image processing suite for electron microscopy*, *J. Structural Biology* **157** (2007), 38–46.
- [35] M. van Heel, *Finding trimeric HIV-1 envelope glycoproteins in random noise*, *Proc. Natl. Acad. Sci. USA* **110** (2013), E4175–E4177.
- [36] J. Vargas, A. L. Álvarez-Cabrera, R. Marabini, J. M. Carazo and C. O. S. Sorzano, *Efficient initial volume determination from electron microscopy images of single particles*, *Bioinformatics* **30** (2014), 2891–2898.
- [37] K. R. Vinothkumar and R. Henderson, *Single particle electron cryomicroscopy: trends, Issues and future perspective*, *Quarterly Reviews of Biophysics* **49** (2016), e13: 1-25.
- [38] L. Wang, A. Singer and Z. Wen, *Orientation determination of cryo-EM images using least unsquared deviations*, *SIAM J. imaging Sci.* **6** (2013), 2450–2483.

*Manuscript received December 29 2017
revised April 22 201*

C.O.S. SORZANO, J.L. VILAS, A. JIMENEZ-MORENO, J. MOTA, T. MAJTNER, D. MALUENDA,
M. MARTINEZ, R. SANCHEZ, J. SEGURA, J. OTON, R. MELERO, L. DEL CANO, P. CONESA, J.
GOMEZ-BLANCO, Y. RANCEL, R. MARABINI, J.M. CARAZO
National Center of Biotechnology (CSIC), c/Darwin, 3, Campus Univ. Autonoma de Madrid,
28049 Cantoblanco, Madrid, Spain

E-mail address: `coss@cnb.csic.es`

J. VARGAS

Dept. Anatomy and Cell Biology, McGill Univ., Montreal, Canada

R. MARABINI

Dept. Computer Science, Univ. Autonoma of Madrid, Madrid, Spain