

# Topology representing network enables highly accurate classification of protein images taken by cryo electron-microscope without masking

Toshihiko Ogura,<sup>a</sup> Kenji Iwasaki,<sup>b</sup> and Chikara Sato<sup>a,\*</sup>

<sup>a</sup> Neuroscience Research Institute and Biological Information Research Center (BIRC), National Institute of Advanced Industrial Science and Technology (AIST), Tsukuba, Ibaraki 305-8568, Japan

<sup>b</sup> Research Center for Ultra-High Voltage Electron Microscopy, Osaka University, 7-1, Mihogaoka, Ibaraki, Osaka 567-0047, Japan

Received 4 February 2003, and in revised form 28 July 2003

## Abstract

In single-particle analysis, a three-dimensional (3-D) structure of a protein is constructed using electron microscopy (EM). As these images are very noisy in general, the primary process of this 3-D reconstruction is the classification of images according to their Euler angles, the images in each classified group then being averaged to reduce the noise level. In our newly developed strategy of classification, we introduce a topology representing network (TRN) method. It is a modified method of a growing neural gas network (GNG). In this system, a network structure is automatically determined in response to the images input through a growing process. After learning without a masking procedure, the GNG creates clear averages of the inputs as unit coordinates in multi-dimensional space, which are then utilized for classification. In the process, connections are automatically created between highly related units and their positions are shifted where the inputs are distributed in multi-dimensional space. Consequently, several separated groups of connected units are formed. Although the interrelationship of units in this space are not easily understood, we succeeded in solving this problem by converting the unit positions into two-dimensional (2-D) space, and by further optimizing the unit positions with the simulated annealing (SA) method. In the optimized 2-D map, visualization of the connections of units provided rich information about clustering. As demonstrated here, this method is clearly superior to both the multi-variate statistical analysis (MSA) and the self-organizing map (SOM) as a classification method and provides a first reliable classification method which can be used without masking for very noisy images.

© 2003 Elsevier Inc. All rights reserved.

**Keywords:** Single-particle analysis; Topology representing network; Growing neural gas network; Cryo-electron microscopy; Image classification

## 1. Introduction

In single-particle analysis, a three-dimensional (3-D)<sup>1</sup> structure is constructed using electron microscopy (EM). This method is advantageous because it does not require a crystal (Frank, 2002; van Heel et al., 2000). Therefore,

single-particle analysis has been applied to membrane proteins whose crystals are difficult to obtain (Radermacher et al., 1994; Sato et al., 2001; Serysheva et al., 1995). Recently, the resolution of such analysis has reached a level better than 10 Å, even for asymmetric molecules (Matadeen et al., 1999; van Heel et al., 2000). In general, EM images of protein are very noisy and, therefore, the primary process of single-particle analysis is the classification of images according to their Euler angles, the images in each classified group then being averaged to reduce the noise level (Frank et al., 1978; van Heel and Frank, 1981). Thus, the method of classification employed is essential for single-particle analysis.

The methods used to classify EM images can be mainly categorized into two approaches: statistical and

\* Corresponding author. Fax: +81-29-861-6478.

E-mail address: [ti-sato@aist.go.jp](mailto:ti-sato@aist.go.jp) (C. Sato).

<sup>1</sup> Abbreviations used: 2-D, two-dimensional; 3-D, three-dimensional; Cryo-EM, cryo-electron microscopy; TRN, topology representing network; GNG, growing neural gas network; SOM, self-organizing map; SA, simulated annealing; MSA, multi-variate statistical analysis; HAC, hierarchical ascendant classification; SD, standard deviation; MRA, multi-reference alignment.

neural network strategies. Multivariate statistical analysis (MSA), in which a particle feature is extracted by reducing variables of the images, is one of the most widely used methods (Frank et al., 1982; van Heel and Frank, 1981). Other statistical methods are the hierarchical ascendant classification (HAC) (van Heel, 1984), the hybridized *k*-means to ascendant classification approach (Frank et al., 1988) and the fuzzy *c*-mean method (Carazo et al., 1990). The accuracies of these methods are decreased by the noise. To reduce the influence of the noise, these methods generally require manual masking, which is adopted for single particle images in most cases. However, in cryo-EM, it is hard to mask the protein image because the protein contrast is very low. In neural network methods, Kohonen's self-organizing map (SOM) is well known as a powerful method for classifying input data by using a two-dimensional (2-D) neuronal sheet (Kohonen, 1982). It has been widely utilized in various fields, including pattern classification (Kanaya et al., 2001; Marco et al., 1998). It has also been successfully applied in the classification of EM images (Marabini and Carazo, 1994; Pascual-Montano et al., 2001; Radermacher et al., 2001). The advantage of this approach is its robustness against noise, and thus the SOM can be applied without masking. However, when the input data have a complex topological structure which must be classified, the SOM is usually not able to set all the input receiving neuronal units at suitable positions (Martinetz and Schulten, 1994). Since a digitized micrograph has a monochrome density at each pixel, the image can be represented in the form of a multi-dimensional vector. In most cases, the distribution of single-particle projections in multi-dimensional space is highly complex because the protein molecule has a complicated structure and/or is freely rotated in a thin buffer layer. In the present paper, the SOM was revealed to produce inadequate unit images, which are the mixtures of the protein projections with different Euler angles. This problem arises mainly because an extremely complex distribution in multi-dimensional space is imposed to fit onto a very simple 2-D latticed, neuronal sheet. Consequently it is hard to set all the neurons at adequate positions in such a classification system.

In contrast, the topology representing network (TRN) (Martinetz and Schulten, 1994; Martinetz et al., 1993), is known to set all the neurons in a 2-D or 3-D complex distribution in response to the input data. Recently, the TRN has been used to combine a high-resolution 3-D structure acquired by X-ray crystallography with volumetric data of protein at lower resolution (Wriggers et al., 1998, 1999). The TRN constructs new nodes, i.e., connections between units which reflect the distribution of the input data. The growing neural gas network (GNG) is one of the TRN algorithms (Fritzke, 1994, 1995), the network structure of which is auto-

matically constructed by the growing process in response to inputs. We found that our newly developed procedure modified from GNG achieves high-performance classification of EM images.

## 2. Materials and methods

### 2.1. Purification of sodium channels and electron microscopy

The sodium channel is a glycosylated membrane protein with a molecular mass of 300 kDa. The extraction of voltage-sensitive sodium channels from the electric organ of *Electrophorus electricus* eels and their purification has been described previously (Sato et al., 1998, 2001). Apoferritin, a soluble protein with a molecular mass of 450 kDa, was kindly provided by Dr. Ichiro Yamashita (Advanced Technology Research Laboratory, Matsushita Electric Ind., Kyoto, Japan). Sodium channel and apoferritin images were recorded from unstained cryo samples using a JEM3000SFF and a JEM3000EFC electron microscope, respectively, at an acceleration voltage of 300 kV (Fujiyoshi, 1998). The micrograph was digitized with a Scitex Leafscan 45 scanner at a pixel size of 2.83 Å at the specimen level, and the applied underfocus values ranged from 3.7 to 7.6 μm for sodium channels, and from 3.0 to 5.4 μm for apoferritins.

### 2.2. Image processing of the learning data

A library of 11,000 images of sodium channels was constructed as previously described (Sato et al., 2001) and apoferritin images were interactively selected from whole cryo-EM images to create a library of 520 images. The images of each protein were aligned rotationally and translationally (van Heel et al., 2000) with the projections from its 3-D model and utilized as inputs. The size of model projections and cryo images were 40 × 40 and 61 × 61 pixels, respectively. Each image was masked by a circle equal in diameter to the side length of the image square. The average of the pixel intensities in each image was adjusted to 128, which is the median value of 8-bit densities.

### 2.3. Algorithms and construction of the growing neural gas network

The growing neural gas network (GNG) is a topology representing network (TRN) (Fritzke, 1994, 1995), in which the adaptation of the synaptic vectors is adopted as earlier proposed by Kohonen (1982). The most important difference of the GNG from the SOM is the process by which a unit-network system is grown, which includes flexible connections by nodes between units.

In the GNG, learning starts from two units. Each unit has its own initial vector composed of a matrix which is the averaged image of all the inputs as interpreted by its pixel densities. Therefore, the unit vector has the same dimensions as the input image,  $40 \times 40$

pixels (= 1600 dimensions) and  $61 \times 61$  pixels (= 3721 dimensions) for the model projection and the cryo-EM, respectively. In order to create the variations, each image of the different Gaussian noise, the parameter of which is set to a standard deviation of  $3\sigma$ , is added to

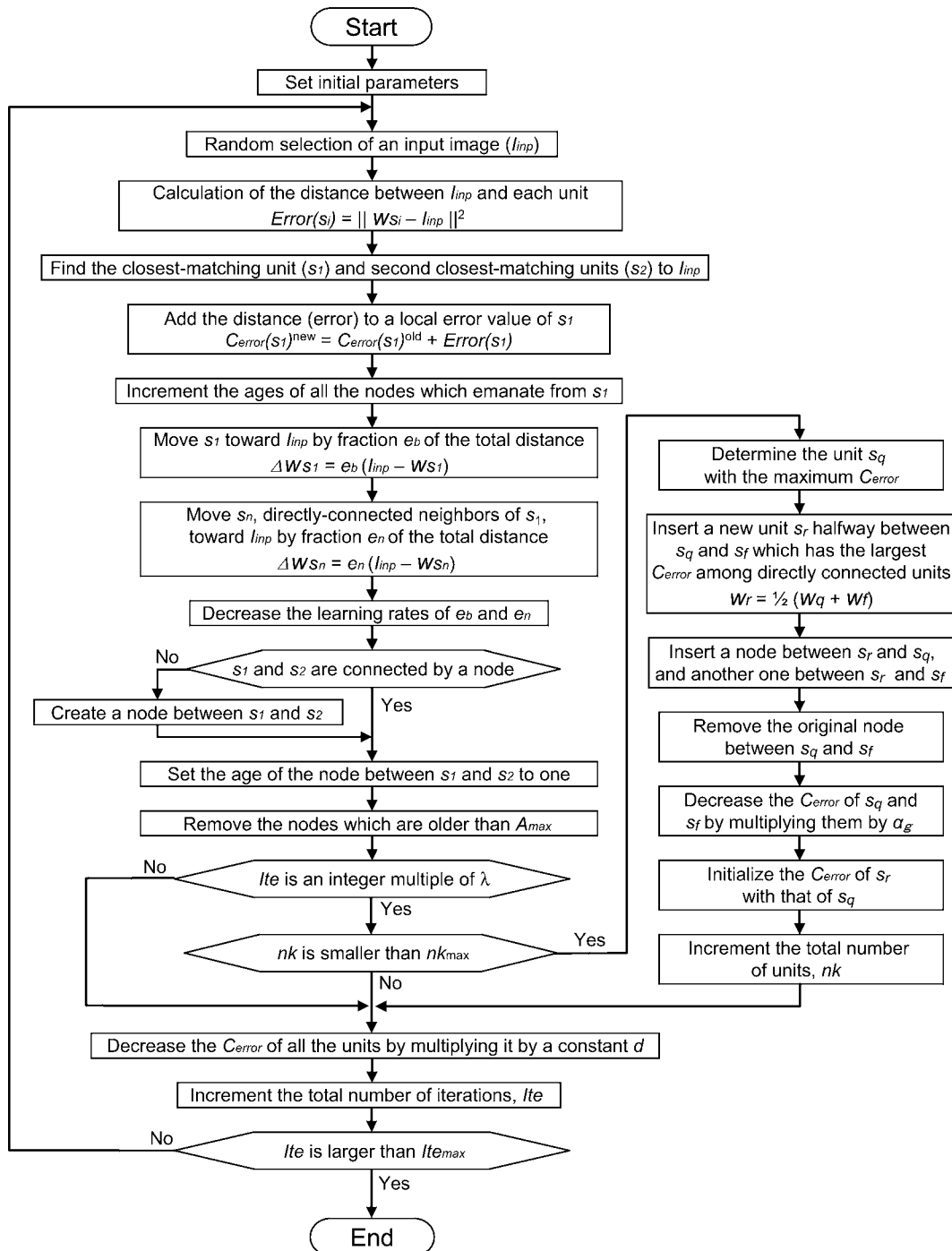


Fig. 1. Flowchart of the GNG algorithm. It comprises several stages: random selection of an input image, search for the unit image which best matches the input, and learning of the input image by the unit. Only the matched unit  $s_1$  and its neighbouring units, which are connected directly to  $s_1$ , learn the input image. The removal and creation of a unit connection by a node are programmed at every iteration, and the creation of a unit is programmed at every predefined iteration.

the initial vector. During learning, the vector is changed gradually by the input images and can be interpreted as a unit image which has the same dimensions as inputs. A unit is added at every  $\lambda$  learning iteration, and the learning is finished at iteration  $Ite_{\max}$ . The algorithm of the modified GNG is given as seen in Fig. 1 as follows:

0. Start from two units which have vectors of the averaged image with different Gaussian noises. Every unit possesses an error counter also, which is initially set at 0. The two units are connected by a node which has an age of 1.

1. Select an input image  $I_{\text{inp}}$  randomly from the input library.

2. Calculate the squared Euclidean distance,  $Error(s_i)$ , between  $I_{\text{inp}}$  and each unit  $s_i$ , which has a vector  $w_{s_i}$ .

$$Error(s_i) = \|w_{s_i} - I_{\text{inp}}\|^2. \quad (1)$$

3. Find the closest-matching unit  $s_1$  and the second closest-matching unit  $s_2$  by the squared distance,  $Error(s_i)$ , as follows. A set  $k$  consists of all the units.

$$s_1 = \arg \min_{s_i \in k} \{Error(s_i)\}, \quad (2)$$

$$s_2 = \arg \min_{s_i \in (k-s_1)} \{Error(s_i)\}. \quad (3)$$

4. Add  $Error(s_1)$  to the error counter,  $C_{\text{error}}$ , of unit  $s_1$ .

$$C_{\text{error}}(s_1)^{\text{new}} = C_{\text{error}}(s_1)^{\text{old}} + Error(s_1). \quad (4)$$

5. Increment the ages of all nodes which emanate from  $s_1$ .

6. Move  $s_1$  towards  $I_{\text{inp}}$  by fraction  $e_b$  of the total distance.

$$\Delta w_{s_1} = e_b(I_{\text{inp}} - w_{s_1}). \quad (5)$$

7. Move  $s_n$ , which are the neighbor units directly connected to  $s_1$ , toward  $I_{\text{inp}}$  by fraction  $e_n$  of the total distance.

$$\Delta w_{s_n} = e_n(I_{\text{inp}} - w_{s_n}). \quad (6)$$

For the first iteration, the learning rates,  $e_b$  and  $e_n$ , have initial values of  $e_{bs}$  and  $e_{ns}$ , respectively.

8. Decrease the learning rates,  $e_b$  and  $e_n$ , from the initial values,  $e_{bs}$  and  $e_{ns}$ , as follows.  $Ite$  and  $Ite_{\max}$  are the current and the maximum number of learning iterations, respectively.

$$e_b = e_{bs} \frac{Ite_{\max} - Ite}{Ite_{\max}}, \quad (7)$$

$$e_n = e_{ns} \frac{Ite_{\max} - Ite}{Ite_{\max}}. \quad (8)$$

This is in contrast to the original GNG method in which all the parameters are fixed (Fritzke, 1995). These steps enhance the convergence in spite of a huge amount of noise.

9. If  $s_1$  and  $s_2$  are connected by a node, set the age of this node to 1. If such a node does not exist, create a new node whose age is 1.

10. Remove the nodes which are older than  $A_{\max}$ .

11. If  $Ite$  is an integer multiple of a parameter  $\lambda$  on the condition that the total number of units,  $nk$ , is smaller than  $nk_{\max}$ , insert a new unit as follows:

- Determine the unit  $s_q$  with the maximum accumulated error,  $C_{\text{error}}$ .

$$s_q = \arg \max_{s_i \in k} \{C_{\text{error}}(s_i)\}. \quad (9)$$

- Insert a new unit  $s_r$  halfway between units  $s_q$  and  $s_f$  which has the largest error,  $C_{\text{error}}$ , amongst the directly connected neighbor units of  $s_q$ .  $w_q$ ,  $w_f$  and  $w_r$  are the vectors of units  $s_q$ ,  $s_f$  and  $s_r$ , respectively. The coordinate of the new unit  $s_r$  is calculated as follows:

$$w_r = \frac{1}{2}(w_q + w_f). \quad (10)$$

- Insert a new node between  $s_r$  and  $s_q$  and another one between  $s_r$  and  $s_f$ . Remove the original node between  $s_q$  and  $s_f$ .

- Decrease the error counters,  $C_{\text{error}}$ , of  $s_q$  and  $s_f$  by multiplying them by a constant  $\alpha_g$ . Initialize the error counter of  $s_r$  with the new error counter of  $s_q$ .

- Increment variable  $nk$  which is the total number of units in the system.

12. Decrease the error counters,  $C_{\text{error}}$ , of all the units by multiplying them by a constant  $d$ .

13. Increment the number of learning iterations,  $Ite$ .

14. If  $Ite$  is not yet  $Ite_{\max}$ , go back to step 1 and iterate.

In the present paper, the original GNG algorithm (Fritzke, 1995) has been modified as follows for use in the classification of protein images in EM. The learning rates,  $e_b$  and  $e_n$ , are decreased during the learning by the annealing method, as shown in Eqs. (7) and (8). Moreover, in the original algorithm, a unit which was not connected by a node was removed. However, such a unit was hardly ever produced in the present classification of the projections. Therefore, our algorithm does not include removal of such a unit.

#### 2.4. Parameter setting of the GNG

In the GNG, eight parameters ( $\lambda$ ,  $Ite_{\max}$ ,  $nk_{\max}$ ,  $A_{\max}$ ,  $e_{bs}$ ,  $e_{ns}$ ,  $\alpha_g$ , and  $d$ ) must be set. In these parameters, the initial learning rates,  $e_{bs}$  and  $e_{ns}$ , are especially important for classification.  $e_{bs}$  determines the amount of change in the unit image which is most similar to the input. Therefore, the parameter has to be adjusted depending on the signal-to-noise ratio of input images. When the ratio is low, the parameter must also be low. In our case of the cryo-EM,  $e_{bs}$  and  $e_{ns}$  were 0.01 and 0.0005, respectively. The maximum iteration constant,  $Ite_{\max}$ , depends on the number of input images in a library.  $Ite_{\max}$

determines the average number of iterated presentations of an image, which we found to be more than five in the case of the cryo-EM presented here, which was sufficient. If the library contains 1000 images, 5000 iterations (1000 images  $\times$  5) or more is suitable. The maximum unit number,  $nk_{\max}$ , is also determined by the total number of inputs and by the signal-to-noise ratio of the inputs. To achieve a good signal-to-noise ratio of unit images, the ratio of the total number of inputs to  $nk_{\max}$  should be more than 20 in the case of the cryo-EM. The number of iterations,  $\lambda$ , which determines the intervals between the creations of units, is calculated by dividing the maximum iteration constant  $Ite_{\max}$  by the maximum unit number  $nk_{\max}$ . It should be slightly smaller than  $Ite_{\max}/nk_{\max}$ . The node age,  $A_{\max}$ , above which a node is eliminated, is important for control of the density of connections between units. If  $A_{\max}$  is small, many nodes are deleted at the early stages and node density is decreased. The value of  $A_{\max}$  which results in an adequate density of nodes ranges from 30 to 50 in the cases presented here. The constants to decrease errors,  $\alpha_g$  and  $d$ , are fixed at 0.5 and 0.995, respectively, as shown by Fritzke (1995) as well as shown here. However, as the nature of the input data considerably varied in both studies, it was not necessary to change these two parameters.

### 2.5. The algorithm of the simulated annealing method

Simulated annealing (SA) is a powerful optimization algorithm that was exploited to anneal the physical process (Kirkpatrick et al., 1983). This method is utilized here to show the interrelatedness simply by rearranging the unit positions on a 2-D map which has a structure of a  $300 \times 300$  lattice grid. Accordingly, the SA algorithm is applied to minimize the node lengths by shifting the unit positions to the optimum, on condition that a certain minimum distance between each pair of units is maintained. First, the acquired GNG map in high-dimensional space is converted into a conventional 2-D connected map. In this step, the positions of the units are initialized randomly according to the 2-D Gaussian distribution. Accordingly, a new coordinate  $(x, y)$  of each unit is randomly extracted from the 2-D normal distribution, the parameter of which was set to a standard deviation (SD) of  $10\sigma$ . After the conversion, the units are reconnected as they were in the previous GNG map in high-dimensional space. The free energy of all the networks,  $E_{\text{all}}$ , is calculated as follows:

$$E_{\text{unit}} = \frac{1}{2} \sum_{i,j} \|U_{(i)} - U_{(j)}\|^2, \quad (11)$$

$$E'_{\text{node}} = \frac{1}{2} \sum_{i,j} E_{\text{node}}(i, j), \quad (12)$$

$$E_{\text{node}}(i, j) = \begin{cases} \|U_{(i)} - U_{(j)}\|^2 & \text{if units } i \text{ and } j \\ & \text{are connected by node,} \\ 0 & \text{otherwise,} \end{cases} \quad (13)$$

$$E_{\text{all}} = E'_{\text{node}} + \frac{(20)^4}{E_{\text{unit}}^2}. \quad (14)$$

Here,  $(20)^4/E_{\text{unit}}^2$  is the local sum of energy caused by the repulsions between units. In this term,  $E_{\text{unit}}$ , the sum of squared Euclidian distance, is calculated by measuring the distances between units as shown in Eq. (11), in which  $U_{(n)}$  is the position  $(x, y)$  of unit  $n$  on the 2-D map. In this map, both  $x$  and  $y$  are integers because  $U_{(n)}$  is always positioned on an intersection of the grid. The local sum of the energy of nodes,  $E'_{\text{node}}$ , becomes larger as the lengths of the nodes increase. Therefore, the total energy,  $E_{\text{all}}$ , grows smaller as the lengths of the nodes decrease, but it increases drastically if the distances between the units decrease too much.

From all the units on the grid, the position of a randomly selected unit,  $n$ , is shifted to a new position  $U_{(n)}^{\text{new}}$  by the addition of vector  $(\Delta x, \Delta y)$ , in which both  $\Delta x$  and  $\Delta y$  are the integers randomly extracted from the 2-D normal distribution whose parameter is set at a SD of  $3\sigma$ . The new total energy,  $E_{\text{all}}^{\text{new}}$ , is calculated by Eqs. (11)–(14). If the change results in a reduction of total energy, it is accepted unconditionally. If the change results in an increase in the total energy, the value of  $P(\Delta E)$  is calculated based on the following equations to judge whether to accept it or not:

$$\Delta E = E_{\text{all}} - E_{\text{all}}^{\text{new}}, \quad (15)$$

$$P(\Delta E) = \exp\left(\frac{\Delta E}{T}\right). \quad (16)$$

Here,  $\Delta E$  is the change in total energy and  $T$  is the current temperature. A random number between 0 and 1 is then generated and compared with  $P(\Delta E)$ . The change is accepted if the random number is less than  $P(\Delta E)$ .

For annealing, we start at a high temperature of  $T_0 = 10,000$ , then decrease the temperature exponentially as iterations by Eq. (17). The end of the annealing is fixed at 200,000 iterations, and the time constant  $\tau$  is set at 20,000.

$$T = T_0 \exp\left(\frac{-S_{\text{ite}}}{\tau}\right). \quad (17)$$

Here,  $S_{\text{ite}}$  is the current number of iterations in the annealing process.

### 2.6. The learning of the SOM

All results of the SOM were calculated according to Kohonen's algorithm (Kohonen, 1982; Marabini and

Carazo, 1994). The map of the SOM has a 2-D lattice structure in high-dimensional space, and both the number of units and the connections are unchanged by learning. Only the positions of units are changed by the learning rates,  $\alpha_s$  and  $\sigma_s$ , which are gradually decreased by Eqs. (18) and (19) during the learning (Marabini and Carazo, 1994).  $\alpha_{s0}$  and  $R_{\max}$  are the initial learning rate and the initial neighborhood radius, respectively.

$$\alpha_s = \alpha_{s0} \left( \frac{Ite_{\max} - Ite}{Ite_{\max}} \right), \quad (18)$$

$$\sigma_s = 1 + (R_{\max} - 1) \left( \frac{Ite_{\max} - Ite}{Ite_{\max}} \right). \quad (19)$$

Here,  $Ite$  and  $Ite_{\max}$  are the current and the maximum number of iterations, respectively.

### 2.7. Image analysis system

All the calculations in the GNG, the SOM, and the SA were performed with the image-processing toolbox of Matlab Version 6 (MathWorks) on a personal computer (Pentium 4: 2 GHz, 2 GB RAM) running Windows 2000. Every system was programmed using the Matlab script M-files. The calculation of the MSA combined with the HAC was performed with Imagic V (van Heel et al., 1996).

## 3. Results

### 3.1. Classification of the model data by the GNG, the MSA, and the SOM

The GNG has a variable structure during the learning because of its flexibility to connect or disconnect units by the growth algorithm (Fritzke, 1994, 1995), unlike the neuronal map of the SOM which possesses a pre-defined connection structure of the lattice. The flow diagram presented in Fig. 1 outlines the GNG learning procedure of EM images. The GNG has two units at the beginning, and the units increase at every specified number of iteration increments,  $\lambda$ , as described in Section 2. This method can position units flexibly to adapt to the distribution of the input images in multi-dimensional space. Each coordinate (unit vector) can be interpreted as a 2-D unit image. This image represents the local basic pattern of the surrounding images, which is similar to the class averages and is expected to be better than those by other methods. The classification ability of the modified growing neural gas network (GNG) was compared with those of the SOM and the MSA, using model projection images immersed in noise and further using cryo-EM of proteins. As a first step to this approach, we adopted a simple model data set comprised of four images which are the projections at different

Euler angles from a 3-D structure model of a sodium channel (Sato et al., 2001). In the single particle analysis of the cryo-EM, the multi-reference alignment (MRA) easily fails to align molecules, especially at the early stages of iterative analysis. That is due to the huge noise in the image and the immature references created by the reference-free method. Furthermore, some special 3-D structures of a molecule also enhance misalignment. Therefore, we considered that the ability to classify an image library which contains misalignments is critical for the classification method. In order to create misalignment in the library, each of the projections was shifted 5 pixels from the center towards one of the four corners of the square (Fig. 2A). Furthermore, 100 various noise images were added separately to each projection to construct a total of four hundred input images (Fig. 2B, 1–4 columns). The monochrome densities at pixels in each image can be represented in the form of a high-dimensional vector, and the inputs which originated from four projections are distributed in four areas of multi-dimensional space. Because the SOM has a squared neuronal map with four corners, it is assumed to easily perform clear classification. As a result, both the GNG and the SOM automatically created unit images, which are similar to class averages, as unit vectors by the classification algorithms. After learning, the SOM with  $3 \times 3$  units showed four clear projections at the corners of the map; however, five mixed images were also observed in other positions (Fig. 2C). This characteristic of the SOM is basically independent of the map size. Even a SOM with  $7 \times 7$  units still showed mixed images in the central area of the map (Fig. 2D), especially in the central unit. This image is clearly a mixture of all four projection images (Figs. 2D and E) because it is similar to the average of all the inputs (Fig. 2B, right end). In order to quantify the inputs classified by a unit, the cross-correlation between an input and each unit image is calculated to find the maximumly related inputs, i.e., member inputs of unit after learning. The number of member inputs which is maximumly related to each unit image is shown in Fig. 2C. Each unit at the corner has 100 member inputs, whereas all the others are 0. Thereby, it was possible to identify the mixed image by the extremely small number of the member inputs in the case of the model data set.

In contrast, the GNG with 9 units created fine averages in all the units (Fig. 3A) in response to the same inputs as those in Fig. 2. It classified projections without confusion. Moreover, the number of member inputs of any unit, which is calculated as in Fig. 2C, was similar to that of any other. The increase in units did not basically influence the result in the GNG, and clear unit images without mixture appeared all over the map (Fig. 3B). The units of the GNG are sequentially numbered upon creation; therefore, the numbers are not related to the connections. The connections among the units are

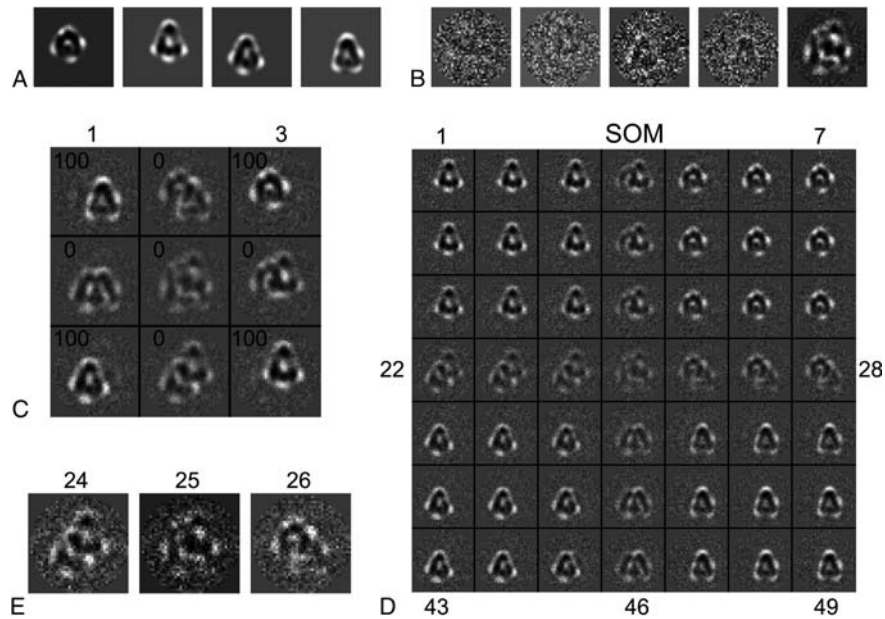


Fig. 2. Classification ability of the SOM using four simple projection images of a sodium channel which are immersed in noise. (A) Four projections from the sodium channel 3-D model (Sato et al., 2001): view close to the top, view close to the side, and views at oblique angles. Each image was shifted 5 pixels from the center towards one of the four corners in the square which has  $40 \times 40$  pixels, as shown. (B) Four examples of input data and an average of all the inputs. To create the input images, 100 various Gaussian noise images, whose parameter was set at a standard deviation (SD) of  $40\sigma$ , were added separately to each projection in (A). The total number of images to be learned was four hundred. The SD of the pixel intensities of the original projection image was almost 9, approximately 1/4 of that of the noise. The total average of the inputs is shown at the right end. (C) The unit images of a SOM with  $3 \times 3$  units after learning. Each unit image is displayed corresponding to its unit position on the square, and a clear projection image is observed at each corner of the map. Mixed images were also created in all the units except these four. The number of member inputs of each unit is shown at the top left of its image. The parameters,  $Iter_{max} = 27,000$ ,  $\alpha_{s_0} = 0.05$  and  $R_{max} = 4$ , were set as described in materials and methods. (D) The unit images of a SOM with  $7 \times 7$  units. Many mixed images were created near the center of the map as in (C). The parameters were set as in (C). (E) Magnified images of the central unit and the neighbouring units which are located to its immediate left and right in (D). The central unit obviously has a mixture of all the four projections as compared with the right-end image in (B).

modified drastically with the creation and removal of nodes by the growing algorithm. Therefore, to understand the grouping, the positions of the units must be reordered according to the connections (Fig. 3C).

Among the classification tasks, the MSA is the most frequently adopted method for single particle analysis and it usually utilized with masking. Here, the MSA combined with the HAC was performed with circular masking equal in diameter to the side length of the image square to compare its ability to classify the same image area used in the GNG and the SOM (Fig. 2B). The resulting 49 class averages are shown in Fig. 3D. Each average contains a higher level of noise than the unit image of GNG: One of the averages was blurred in the MSA (Fig. 3D, No. 48). These results demonstrate that the GNG is more useful than the SOM and the MSA to classify a simple data set.

### 3.2. Comparison of the classification accuracy between the GNG and the SOM

We compared the quality of unit images between the GNG (Fig. 3B) and the SOM (Fig. 2D) with 49 units in response to the same inputs (Fig. 2B). The cross-corre-

lation between a unit image and each of the four original projections (Fig. 2A) was calculated, and the maximum value of the four cross-correlations was adopted to evaluate the quality of the unit image. In the histogram of the SOM, the 49 maxima of 49 units have two peaks at 0.95 and 0.65 (Fig. 4A, upper row, black columns). The higher peak is attributed to the images which are similar to only one of the four original images, whereas the lower peak is attributed to the mixed images (Fig. 2D). On the other hand, the GNG shows a single peak at 0.75. The peak of the GNG is positioned lower than that of the SOM. The minimum value in the histogram of the GNG, however, is positioned higher than that of the SOM (Fig. 4A, upper row). Next, in order to determine the effect of the library size on the unit images, we increased the number of inputs into the GNG. The peak becomes sharp and shifts to a high position close to 1 due to the increase (Fig. 4A, red columns in the middle and lower rows). This demonstrates that unit images improve as the number of inputs increases. With 1200 inputs, the peak of the GNG becomes higher than that of the SOM (Fig. 4A, lower row). In contrast, the position of the peak is independent of the number of learning images in the SOM (Fig. 4A, black columns).

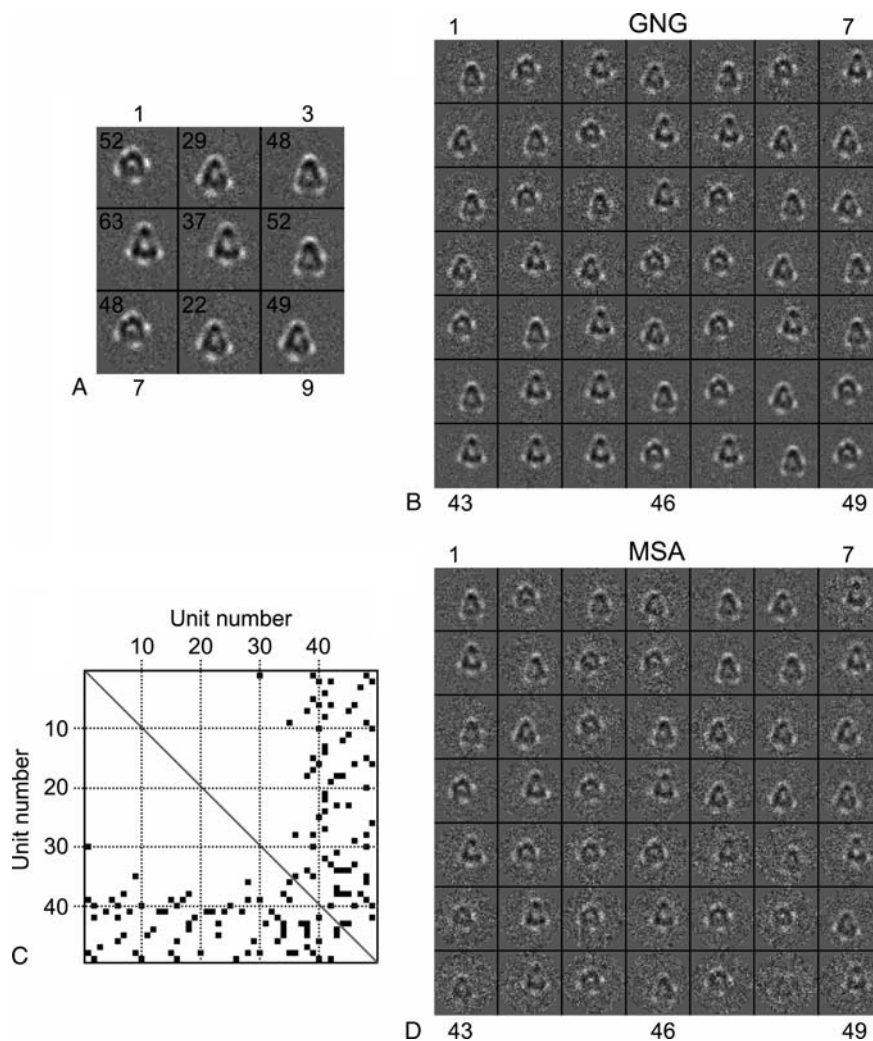


Fig. 3. Classification ability of the GNG and the MSA using the same inputs as in Fig. 2. (A) The unit images of a GNG with 9 units after learning. All the units in the GNG have clear projections. The number of member inputs of each unit is shown at the top left of its image. The total number of iterations,  $Ite_{max}$ , was 7000. A new unit was created in the network at every 500 iterations by the growth algorithm, whereas the maximum unit number  $nK_{max}$  was 9. The parameters,  $e_{hs} = 0.05$ ,  $e_{ns} = 0.005$ ,  $\alpha_g = 0.5$ ,  $d = 0.995$ , and  $A_{max} = 50$ , were set as described in Section 2. (B) The unit images of a GNG with 49 units after learning. The total number of iterations,  $Ite_{max}$ , was 27,000 and the maximum unit number  $nK_{max}$  was 49. The other parameters were the same as in (A). (C) Connections created in the GNG in (B) after the learning. Ordinate and abscissa axes show the unit numbers. A connection between units is shown by a black square. (D) The class averages created by the MSA combined with the HAC. The total number of eigenimages adopted in the analysis was 69, whereas the number of iterations in the HAC was 24. The images are classified using a circular mask of maximum size in the image square.

Therefore, it can be concluded that the GNG can create a better unit image than the SOM if there are a relatively large number of inputs.

If a unit has an average image without a mixture, only one of the four cross-correlations between its unit image and the four original images is high because the unit image is similar to only one of the original images. If a unit has a mixed projection image, more than one of the four correlations are high because its unit image is similar to some of the original images. In this case, the maximum is closer to the second maximum of the correlations. Therefore, the subtraction of the second-maximum from the maximum must be an index of the degree of confusion in a unit. In the

GNG, the histogram of the subtractions (Fig. 4B, red column) shows a peak similar to that of the maximum (Fig. 4A, red columns), reflecting the small values of the second-maxima. In contrast, the histogram is clearly divided into two peaks, the higher peak at 0.85 and the lower at 0.15 in the SOM (Fig. 4B, black column), which reflects non-mixed and mixed unit images, respectively. Again, the total shape of this histogram did not change as the number of the inputs increased. These results coincide with the fact that the SOM had both mixed and non-mixed projection images in the same map. In contrast, the GNG was shown to possess only one of the original projections in every unit.



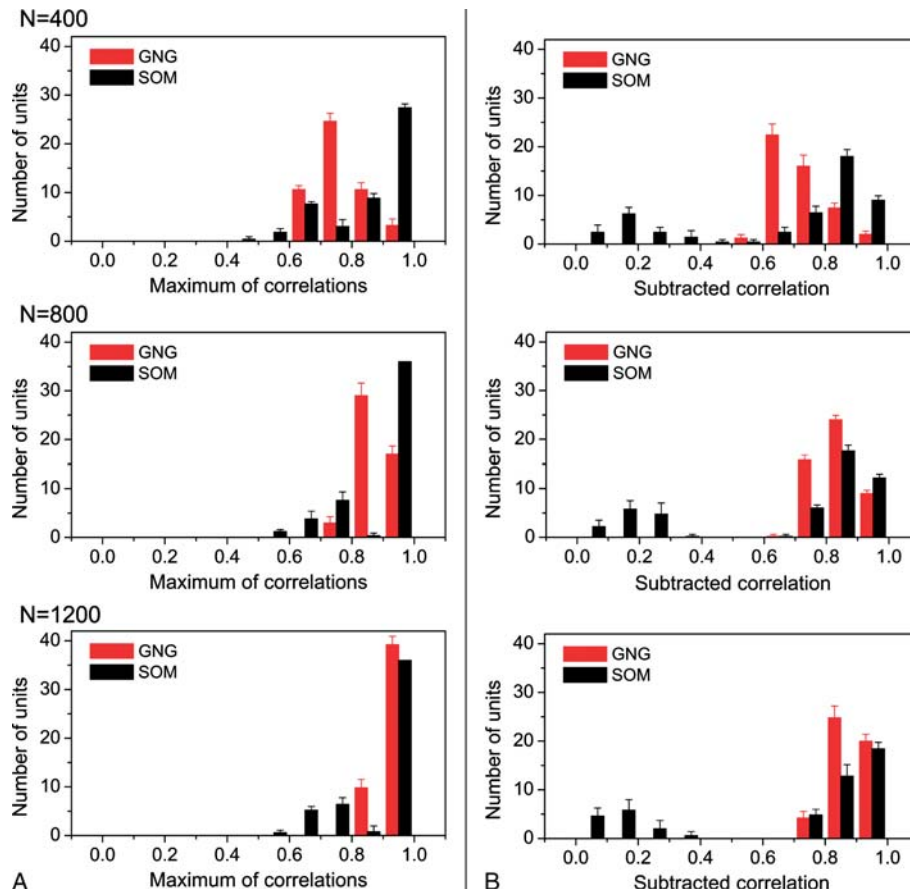


Fig. 4. Comparison of the SOM and the GNG as to classification ability. (A) A SOM with  $7 \times 7$  units (Fig. 2D) and a GNG with 49 units (Fig. 3B) after learning were compared by the cross-correlations between their unit images and the original four projection images (Fig. 2A) in order to determine their classification ability after learning. The abscissa axis shows the maximum value among the four cross-correlations between each unit image and the four original projections, whereas the ordinate axis shows the number of units. The SOM shows two peaks at 0.95 and 0.65 (black bar); basically neither peak was shifted by an increase in the total number of learning images,  $N$ , from 400 (upper) to 800 (middle), and further to 1200 (lower). In contrast, the GNG showed only one peak, which drastically shifted to a higher position as the number of inputs increased (red bar). Generally, each unit image of GNG has higher cross-correlations than the corresponding one of SOM after the learning of many inputs. (B) Detection of mixed unit images by the cross-correlations. The abscissa axis shows the subtraction of the second-maximum from the maximum of the cross-correlations calculated as in (A) for each unit. The SOM clearly shows two peaks (black bar) which can be divided by a value of 0.5. The lower peak at 0.15 and the higher peak at 0.85 reflect mixed and non-mixed unit images, respectively. In contrast, the GNG shows only one peak (red bar) as in (A), reflecting that its unit image is similar to only one of the four originals. The histogram and the error bar in (A,B) show the average and standard deviation (SD) of five independent learnings, respectively.

### 3.3. Optimum rearrangement of the GNG units in a 2-D space to show the connected relations

Each unit has a clear projection image after learning in the GNG. However, the map in Fig. 3B is simply ordered in the sequence of the unit creations. Therefore, we cannot understand the relationships of the images simply from the sequential numbers of units (Fig. 3B). To further classify the unit images based on the unit connections (Fig. 3C), the units in multi-dimensional space were converted into units in a 2-D space with node connections. Accordingly, the units were rearranged by shifting their positions to make the node lengths minima in a 2-D space. This type of problem is generally known as an optimum arrangement problem of the positions to

create minimum connections, and some of the problems can be solved by using a simulated annealing (SA) method (Kirkpatrick et al., 1983; Lambert and Hittle, 2000). Therefore, an SA method was designed here to minimize the node lengths on condition that a certain distance between units was maintained.

We optimized the GNG with 49 units (Fig. 3B) using the SA method, and the units were completely divided into four groups (Fig. 5), which corresponded to the four original images (Fig. 2A). Each unit in a group included the same original projection image, and node connections were created only within a group. No node connection was created between the groups. The optimized map shows relations of units clearly and, furthermore, it classifies unit images as connected groups.

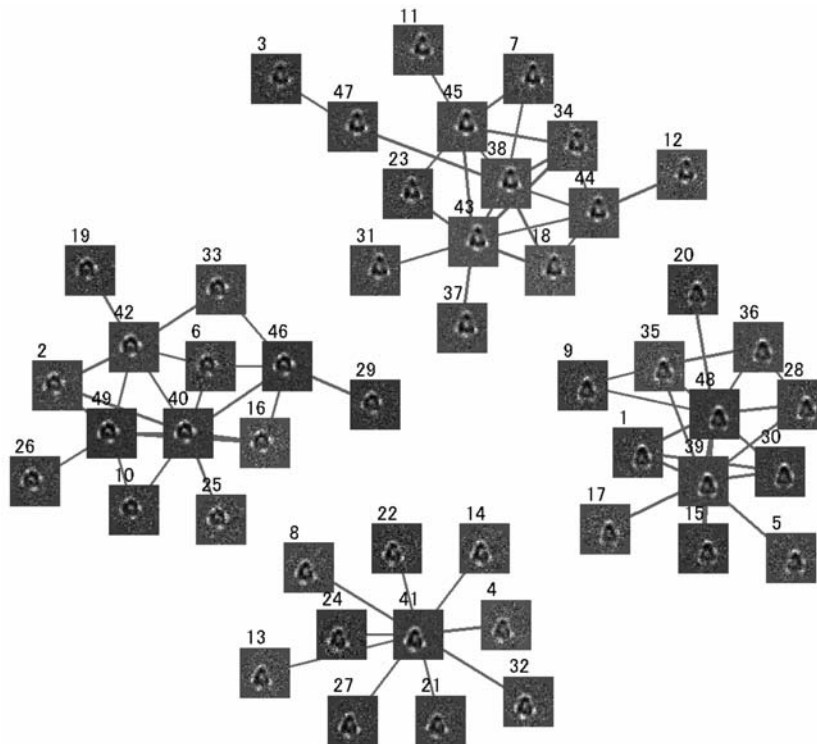


Fig. 5. Optimum rearrangement of the GNG units on a 2-D map by the simulated annealing method (SA). After learning, the GNG units in multi-dimensional space (Fig. 3B) were converted into 2-D space at random positions, and the positions were optimized in view of the minimization of the total of node lengths on the condition that a certain distance between units was maintained as described in Section 2. Each unit number corresponds to that in Fig. 3B. In the optimized map, the units were clearly divided into four groups corresponding to the original four projections (Fig. 2A). The maximum number of iterations was 200,000. The initial temperature  $T_0$  was 10,000, and the temperature  $T$  decreased exponentially as the iterations increased.

#### 3.4. Classification of a model image library including contamination

Generally, high-resolution single-particle analysis requires a huge number of particle images excised from electron-micrographs. Therefore, the data set of the particle images includes impurities, i.e., ice or contaminated protein particles in most cases. Some of them might be caused by incomplete purification of the protein, especially when the target belongs to some kinds of membrane proteins which are difficult to purify. The contaminated images are preferably discarded from the library in advance since their presence reduces the resolution of the final 3-D model. If they can be discarded by the present classification algorithm, it is very meaningful.

In order to construct a model of such an inhomogeneous library, a projection from the 3-D model of a sodium channel was rotated clockwise through  $360^\circ$  by  $45^\circ$  increments to create eight images, and further round and square white images were created artificially (Fig. 6A, upper row). To each image we added 100 various noise images synthesized artificially in order to create 100 images immersed in noise; thereby, a total of 1000 images were prepared (Fig. 6A, lower row). Both

the GNG and the SOM were basically able to classify the sodium channel from round and square images to some extent. In the case of the SOM, the bottom left unit has a clear artificial round image (Fig. 6B, No. 43); however, the neighbouring unit on the right has a confused round image with the sodium channel projection (Fig. 6B, No. 44). Similarly, the two units above the bottom left show mixed images of artificial square and round images (Fig. 6B, Nos. 29 and 36). Furthermore, the unit located at the bottom of the third column shows a typical mixed image of the projection and circle. Likewise, in image No. 21, there is a mixture of sodium channel projections in opposite directions. This and other examples of unit images created by the SOM are shown in Fig. 6C. In contrast, every unit of the GNG has a clear sodium channel projection or artificially created image without mutual confusion (Fig. 6D). Furthermore, the projections rotated at different angles are precisely classified and averaged in this figure. In the optimized connection map of the GNG, the units with artificial images are divided into two groups which can be distinguished from each other (Fig. 6E, right). Moreover, the sodium channel projections are classified according to their angles of rotation, and the units with the same projection are mostly in the same group in

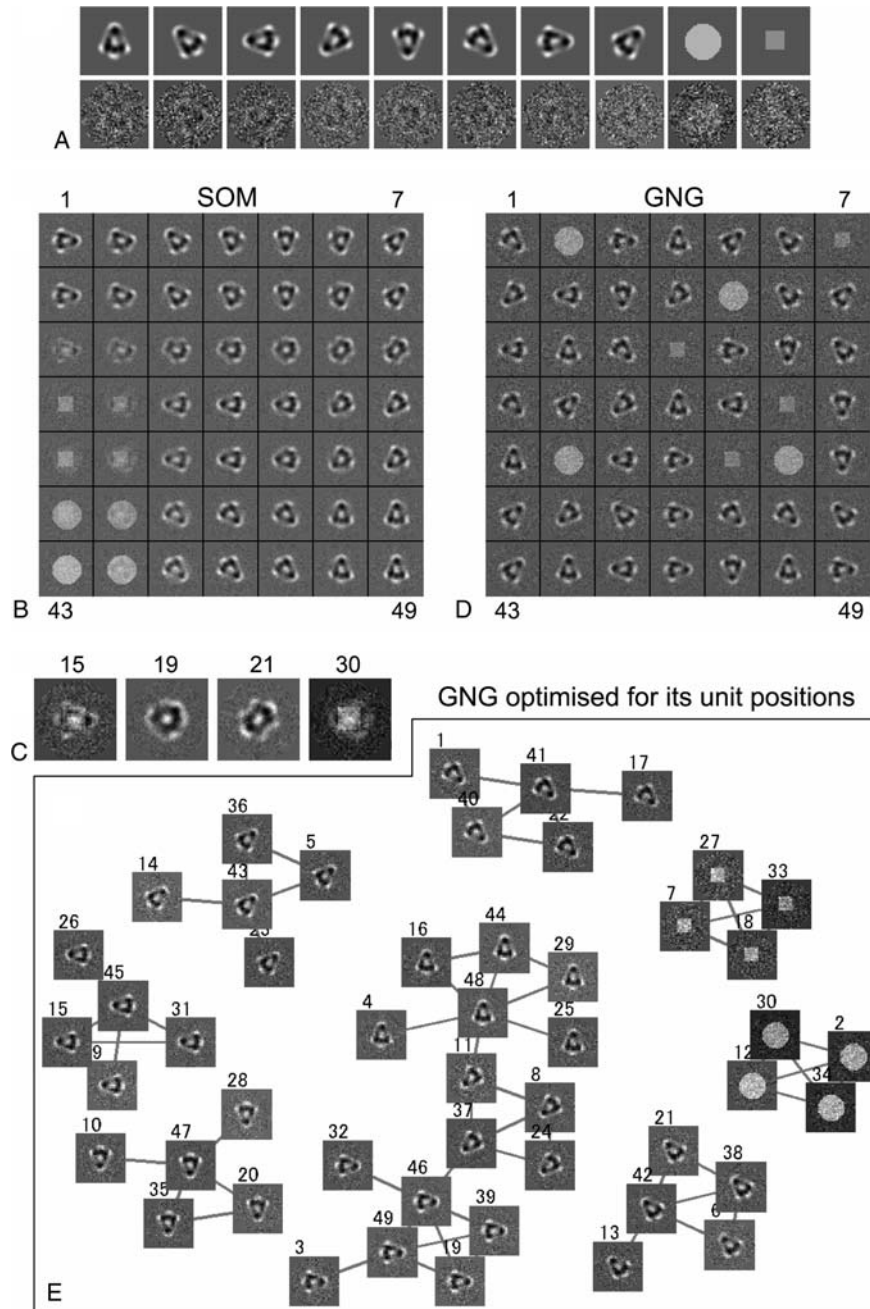


Fig. 6. Comparison of the GNG and the SOM as to classification ability using model projections and artificially created contaminants. (A) A projection image from a 3-D model of the sodium channel was rotated clockwise through  $360^\circ$  by  $45^\circ$  increments to provide eight original images, and a white circle and square images were also created artificially (upper row). To each image was added 100 various noise images, as in Fig. 2B, to provide a library of 1000 images in total. Ten examples are shown (lower row) where each image was created from the upper original image. (B) Unit images of a SOM with  $7 \times 7$  units after learning. Many units have the mixtures of images with different Euler angles. Moreover, a number of units contain a mixture of the projection and the artificial image. The parameters adopted here are the same as those in Fig. 2D. (C) Examples of mixed images extracted from (B). The unit number is shown above. (D) Unit images of a GNG with 49 units after learning. Every unit shows a clear projection, artificial square or round image without confusion. The parameters adopted here are the same as those in Fig. 3B. (E) The optimized connection map of the GNG by the SA. The unit positions in (D) are optimized as in Fig. 5. The units with artificial circular or square images are located in the right of the map and are clearly separated from those with sodium channel projections. The projections themselves are further classified by the connections depending on their Euler angles. The parameters adopted in the SA are the same as those in Fig. 5.

which units are connected with each other. In the SOM, it is difficult to discard these contaminated artificial images completely when the structure of the target

protein is unknown. That is because several projection images in its units are influenced by the contaminants. In contrast, the GNG clearly distinguishes the

projection image from the artificial images, and we can easily discard a group of contaminants based on the connections of the units.

### 3.5. Classification of membrane protein images taken by cryo-EM according to their Euler angles

Next, we compare the classification accuracy of the SOM and the GNG using 11,000 particle images of sodium channel protein taken by cryo-EM (Sato et al., 2001). The SOM map shows clear averages of projections at the four corners (Fig. 7A). However, the units

located near the center contain unclear projections which are mixtures of various Euler angle projections. Among them, some of the unit images also possess doughnut-like outlines (Fig. 7A, Nos. 3, 4, and 9) indicating confusion in the projections with closed Euler angles. In contrast, the units of the GNG have clear projections without confusion throughout the map, and the four massive corners near the molecular bottom are clearly observed by the auto-averaging function presented here (Fig. 7B, Nos. 2, 18, and 40). As a result, a certain kind of projections, which constitute large portion of the inputs, create a large number of similar unit

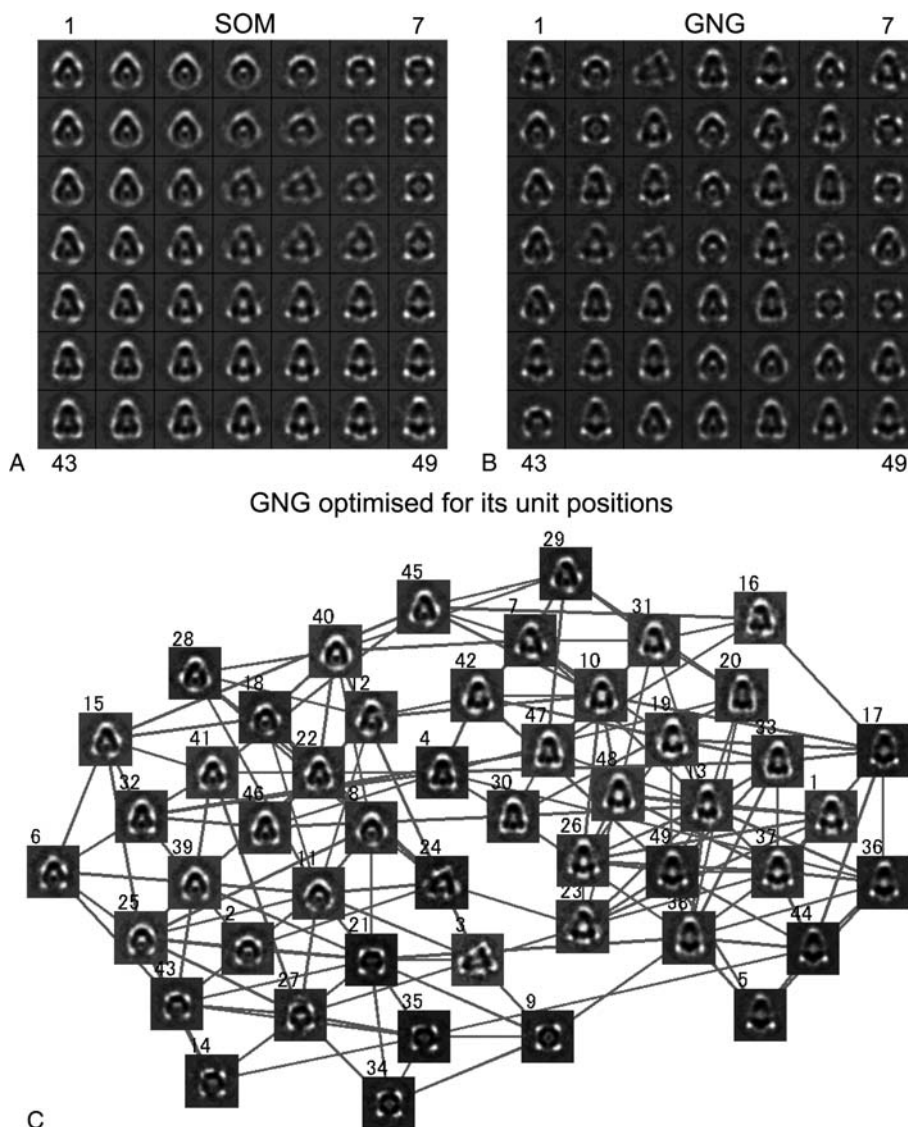


Fig. 7. Accuracy of the SOM and the GNG in classification of sodium channel images taken by cryo-EM. (A) Unit images of a SOM with  $7 \times 7$  units after learning. Several units possess clear images; however, a number of units show unclear mixed projections. The mixed images sometimes have round doughnut-like outlines which reflects the confusion which takes place between similar but not identically oriented projections. The parameters of the SOM were as follows:  $It_{\max} = 50,000$ ,  $\alpha_0 = 0.05$ , and  $R_{\max} = 4$ . (B) Unit images of a GNG with 49 units after learning. Every unit shows a clear projection. The parameters of the GNG were as follows:  $It_{\max} = 50,000$ ,  $\lambda = 1000$ ,  $e_{bs} = 0.01$ ,  $e_{ns} = 0.0005$ ,  $\alpha_g = 0.5$ ,  $d = 0.995$ , and  $A_{\max} = 30$ . Images located at Nos. 3 and 24 reflect the misalignments which had been included in the image library. (C) The optimized map of (B) by application of the SA. The actual top view of the sodium channel is located at the bottom center and gradually changes to the side view as the unit shifts in a clockwise direction. The parameters adopted in the SA were the same as those in Fig. 5.

images after learning. In the optimized connection map, the real top view of the sodium channel is located at the bottom center (Fig. 7C, No. 34), and the unit image is gradually changed to its side view as the unit is shifted in a clockwise direction by the rotation. In the map, all the units are connected and there is no isolated unit. This is in good accordance with the random orientation of a sodium channel particle in the frozen buffer layer (Sato et al., 2001). Again, by using the cryo-EM images as inputs, the optimized connection map of the GNG is shown to provide more abundant and beneficial information than the SOM.

### 3.6. Exclusion of contaminated protein using GNG from the cryo-EM image library

Impurities in an image library caused by imperfect protein purification is a significant problem in single particle analysis. Therefore, we mixed the 520 apoferritin images taken by cryo-EM (Fig. 8A) as impurities into the library of the sodium channel. In the SOM, the apoferritin projection is located at the bottom right on the map, and, again, its closest units have interpolations of the apoferritin into the sodium channel projection as shown in Fig. 8B (Nos. 42 and 48). The strange round bottom of the sodium channel molecule is also observed again in Fig. 8B (Nos. 46 and 47). The numbers of member inputs of these four units (Nos. 42, 46, 47, and 48) are 104, 264, 224, and 94, respectively (Fig. 8B, red numbers). These numbers are not so small. Therefore, only from the number of member inputs, it is hard to remove all the contaminated apoferritin images from the

unit images. The class averages calculated by the MSA combined with the HAC with the same mask are basically similar to the model projections; the image itself, however, is not so clear (Fig. 8C). This suggests that each class average contains wrong constituents with different Euler angles. In the GNG, clear apoferritin images were generated in two units (Fig. 8D, Nos. 3 and 46), and only one unit had a mixed image between apoferritin and the sodium channel (Fig. 8D, No. 28). In the optimized map of the GNG, the two units with apoferritin images protruded from a sodium channel cluster at the upper right (Fig. 9), whereas the unit with the mixed image was located on the boundary between apoferritin and the sodium channel on the map (Fig. 9, No. 28). Because the units with apoferritin images are located peripherally, it is easy to obviate the contaminated apoferritin from the sodium channel averages. Even by using a contaminated library, the GNG can classify sodium channel projections accurately based on their Euler angles.

## 4. Discussion

In single-particle analysis, the classification accuracy of particle images is one of the most important factors for high resolution in a 3-D reconstruction. Until now, the MSA in combination with the HAC has been widely used (Frank et al., 1982; van Heel, 1984; van Heel and Frank, 1981). On the other hand, we have developed a highly accurate automatic pickup method using a three-layer neural network (Ogura and Sato, 2001, 2003)

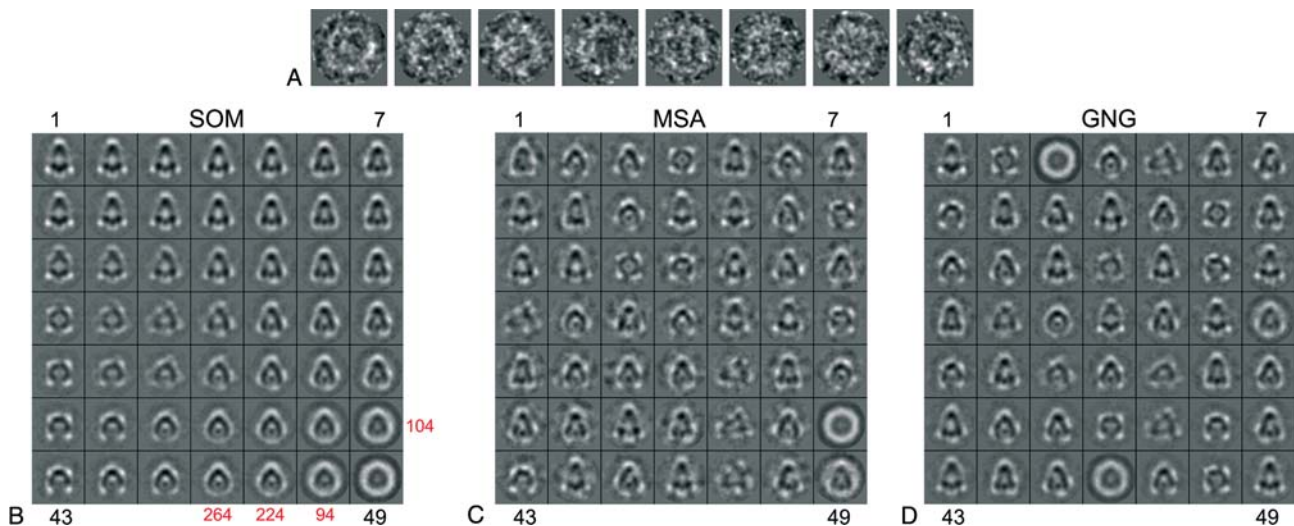


Fig. 8. Accuracy of the GNG, the SOM, and the MSA to classify a mixed library composed of sodium channel and apoferritin images taken by cryo-EM. (A) Eight examples of apoferritin images taken by cryo-EM. To create a mixed library, 520 apoferritin images and 11,000 sodium channel images were prepared. (B) Unit images of a SOM with  $7 \times 7$  units after the learning of the mixed library. A projection image of apoferritin is created in the bottom right unit on the map; however, the surrounding units have mixed images between the sodium channel and apoferritin. The parameters adopted here were the same as those in Fig. 7A. The number of member inputs of each unit is shown in red at the side for several confused unit images. (C) The class averages created by the MSA combined with the HAC. The inputs are classified with the circular mask as in Fig. 3D. (D) The unit images of a GNG with 49 units after the learning. Most of the unit images are clear without confusion.

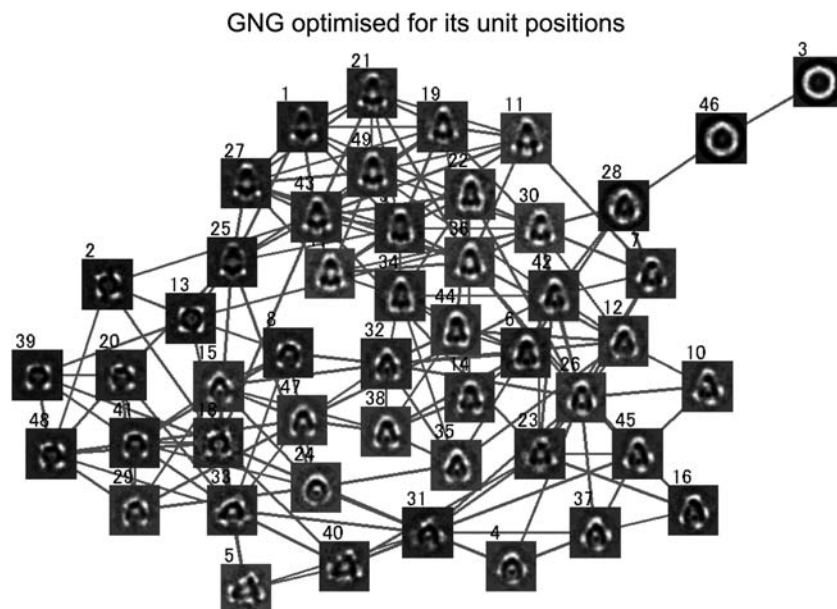


Fig. 9. The optimized map of Fig. 8D by applying the SA. Two apoferritin images protrude in the upper right from a large cluster of sodium channel images in the optimized map and form a group of two. However, a mixed image (No. 28) is observed between the sodium channel and apoferritin groups, which can be easily distinguished by its intermediate position on the map.

which enables the pickup of an enormous number of particles. Therefore, the classification methods are now required to process more than 100,000 images. In general, it is not easy to classify such an enormous number of images by the MSA or its hybridized method using a normal computer system. That is because the calculation is too heavy and requires a huge system memory. Furthermore, to reduce the effects of noise, these methods generally require masking which includes the possibility of obtaining artificial effects.

To date, the SOM has been a powerful method in the classification algorithms of EM because of its resistance to noise and its processing speed (Marabini and Carazo, 1994). However, the SOM is known to cause topological mismatches in response to a variety of learning data (Martinetz and Schulten, 1994). Generally, the number of dimensions included in an image is the same as the number of its constituent pixels, implying that an EM image contains high-dimensional data. Moreover, a solubilized single particle of protein tends to rotate freely in a thin buffer layer. In such a case, the particle images taken by the cryo-EM frequently show complex distribution in high-dimensional space. However, the SOM has a simple connection structure of a 2-D lattice or hexagonal cluster. Due to this disparity between the complexities of the inputs and the simple structure of the SOM, mixed images are produced in mismatched positions as revealed in the present paper. This problem has been attributed to the fixed connected relations of units in the SOM. By the learning data consisting of four groups (Fig. 2B), which correspond to four original projections (Fig. 2A) in high-dimensional space, unit

images similar to the original projections were created in the units at the corners of the map (Fig. 10A). Thereby, during the learning, each of these four units with the matched images responded most frequently to the

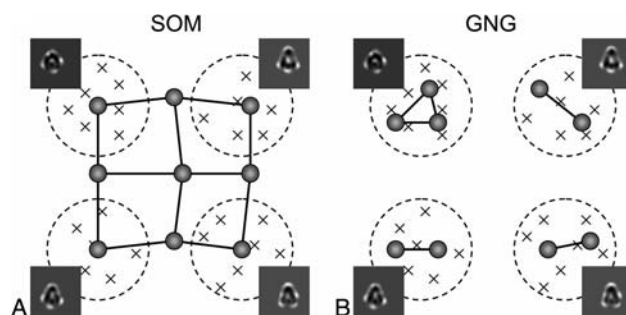


Fig. 10. Schematic representation of the classification mechanisms used by the SOM and the GNG. A cross represents the position of an input image in Figs. 2 and 3, and a dotted circle represents the area where a group of the inputs is distributed in multi-dimensional space. Each projection indicates the original image (Fig. 2A) of the group. A filled circle and a black line represent a neuronal unit and connection node, respectively. (A) A schematic demonstrating recognition by a SOM in Fig. 2C. The SOM has a lattice structure of neural units, and its connections by the nodes are not changed. Hence, it is hard to set all the units at the appropriate positions when a distribution of inputs is not suitable for this network structure. In the case of four groups of inputs as in Fig. 2, the central unit is far apart from the four areas whose inputs are distributed, consequently creating a mixed image of all four original projections. (B) A schematic demonstrating recognition by a GNG in Fig. 3A. The GNG produces units at the positions where inputs are crowded in a multi-dimensional space, and the nodes are formed to connect the units with unit images which are similar to each other. Therefore, the GNG is able to set all the units at the appropriate positions and classify them properly.



inputs, while the average of all the four original images was created in the center on the square, far away from the four corners. Therefore, in learning, the units located near the center were able to develop unit images which were less similar to the inputs and became gradually unresponsive to the inputs, thus losing further movement. The number of member inputs of such unit was found to be 0, suggesting the possibility of eliminating the confused images of the units by the number of the member inputs. However, using the EM library of sodium channel and apoferritin, it was found to be hard to eliminate mixed images by using the number of member inputs. Moreover, we consider that it is better to avoid a method which creates mixed images. When there are many mixed images, the non-necessary units with such images consume unnecessary calculation time.

To avoid producing an unresponsive unit, the kernel-based SOM has been developed recently (Pascual-Montano et al., 2001; Van Hulle, 1999), in which each unit has a receptive field which accepts various inputs. The algorithms of the kernel-based SOM are categorized into several types based on the regulation of the receptive fields and the learning processes. In this system, each unit can have equal probability to receive inputs by decreasing the number of the inactivated units. However, equal probability does not necessarily lead to a clear unit image, especially in the classification of images. When each unit has a broad receptive field that overlaps those of its neighbours, the units in kernel-based SOM tend to form mixed averages because the receptive field expedites activations of various image inputs at equal probabilities. Moreover, as with the conventional SOM, this system also adopts fixed connected relations in the 2-D lattice or hexagonal map. Therefore, the fundamental problem of a gap between the simple connection structure and the highly complex distribution of inputs remains and, therefore, would not be solved by this method.

In contrast, the GNG has a flexible structure of units and node connections (Fritzke, 1994, 1995). Therefore, by the growing process, the GNG is able to set its units at suitable positions based on the coordinates of the inputs in multi-dimensional space (Fig. 10B). By means of this strategy, the GNG creates clear projections in all the units after learning. However, the units are not numbered in a certain order determined by the connections because the connections are rearranged by learning. We succeeded in solving this problem by the optimization using the SA method. In an optimized connection map of the GNG, the relationships of the units can be more clearly understood than those of the SOM. If similar unit images arise in a connected group, these images can be merged according to the Euclidean distances between the units in multi-dimensional space. In GNG, each unit image is not a simple average, but an average of the member images weighted differently.

Accordingly, the unit image can be used almost as a class average. Therefore, the GNG can be utilized to reconstruct a 3-D structure in combination with the MRA. For the MRA, high quality references can also be created automatically in the picking up task of particles by the neural network method (Ogura and Sato, 2003) and it can be utilized to align precisely. Moreover, the GNG in combination with the SA can be used not only to classify images according to its Euler angles but to eliminate the contaminated images from the image library.

In the classification by the SOM, the results presented here clearly show that the user has to eliminate the confused unit images manually. The exclusion of the mixed images would be difficult to perform if the target protein structure is unknown and/or contaminated by other proteins. In addition, the calculation speed of the GNG is almost two times faster than that of the SOM because the GNG increases units starting from only two, making the initial calculations very light. Moreover, the number of units can be set in the GNG arbitrarily. The MSA combined with the HAC using a circular mask of maximum size in the image is shown to not be as effective as the GNG. Therefore, we conclude that the GNG is a much more powerful method for classifying images than the MSA and the SOM.

The robustness of the GNG against noise opens up new fields for single-particle analysis, such as protein complexes in which the components are repeatedly connected to and disconnected from each other. Such protein complexes are sometimes very important in biological functions, e.g., signal transductions, protein folding, and RNA transcriptions. In such a case, the diameter of the complex is changed according to the clustering, and the use of masking to reduce the effect of noise is very difficult. Therefore, the GNG classification system without masking must be exclusively effective for carrying out a dynamic analysis of a protein complex at various clustering states with a huge number of images, which are collected in combination with the automatic pickup program (Ogura and Sato, 2001, 2003).

## Acknowledgments

The authors express their cordial thanks to Dr. Yoshinori Fujiyoshi (Kyoto University, Japan) for the cryo-microscopy and illuminating suggestions and to Dr. Andreas Engel (MIH, Biozentrum, Basel, Switzerland) for his invaluable discussions. Dr. Ichiro Yamashita (ATRL, Matsushita Electric Ind. Co., Japan) kindly provided the apoferritin. This work was supported by a grant from the Japan New Energy and Industrial Technology Development Organization (NEDO) and by grants from AIST.

## References

- Carazo, J.M., Rivera, F.F., Zapata, E.L., Radermacher, M., Frank, J., 1990. Fuzzy sets-based classification of electron-microscopy images of biological macromolecules with an application to ribosomal particles. *J. Microsc.* 157, 187–203.
- Frank, J., 2002. Single-particle imaging of macromolecules by cryo-electron microscopy. *Annu. Rev. Biophys. Biomol. Struct.* 31, 303–319.
- Frank, J., Bretaudiere, J.P., Carazo, J.M., Verschoor, A., Wagenknecht, T., 1988. Classification of images of biomolecular assemblies: a study of ribosomes and ribosomal subunits of *Escherichia coli*. *J. Microsc.* 150, 99–115.
- Frank, J., Goldfarb, W., Eisenberg, D., Baker, T.S., 1978. Reconstruction of glutamine synthetase using computer averaging. *Ultramicroscopy* 3, 283–290.
- Frank, J., Verschoor, A., Boublik, M., 1982. Multivariate statistical analysis of ribosome electron micrographs. L and R lateral views of the 40 S subunit from HeLa cells. *J. Mol. Biol.* 161, 107–133.
- Fritzke, B., 1994. Growing cell structures—a self-organizing network for unsupervised and supervised learning. *Neural Networks* 7, 1441–1460.
- Fritzke, B., 1995. A growing neural gas network learns topologies. In: Tesauro, G., Touretzky, D.S., Lenn, T.K. (Eds.), *Advances in Neural Information Processing Systems 7*. MIT Press, Cambridge, pp. 625–632.
- Fujiyoshi, Y., 1998. The structural study of membrane proteins by electron crystallography. *Adv. Biophys.* 35, 25–80.
- Kanaya, S., Kinouchi, M., Abe, T., Kudo, Y., Yamada, Y., Nishi, T., Mori, H., Ikemura, T., 2001. Analysis of codon usage diversity of bacterial genes with a self-organizing map (SOM): characterization of horizontally transferred genes with emphasis on the *E. coli* O157 genome. *Gene* 276, 89–99.
- Kirkpatrick, S., Gelatt, C.D., Vecchi, M.P., 1983. Optimization by simulated annealing. *Science* 220, 671–680.
- Kohonen, T., 1982. Self-organized formation of topologically correct feature maps. *Biol. Cybern.* 43, 59–69.
- Lambert, T.W., Hittle, D.C., 2000. Optimization of autonomous village electrification systems by simulated annealing. *Solar Energy* 68, 121–132.
- Marabini, R., Carazo, J.M., 1994. Pattern recognition and classification of images of biological macromolecules using artificial neural networks. *Biophys. J.* 66, 1804–1814.
- Marco, S., Ortega, A., Pardo, A., Samitier, J., 1998. Gas identification with tin oxide sensor array and self-organizing maps: adaptive correction of sensor drifts. *IEEE Trans. Instrum. Measur.* 47, 316–321.
- Martinetz, T., Schulten, K., 1994. Topology representing networks. *Neural Networks* 7, 507–522.
- Martinetz, T.M., Berkovich, S.G., Schulten, K.J., 1993. Neural-gas network for vector quantization and its application to time-series prediction. *IEEE Trans. Neural Networks* 4, 558–569.
- Matadeen, R., Patwardhan, A., Gowen, B., Orlova, E.V., Pape, T., Cuff, M., Mueller, F., Brimacombe, R., van Heel, M., 1999. The *Escherichia coli* large ribosomal subunit at 7.5 Å resolution. *Structure* 7, 1575–1583.
- Ogura, T., Sato, C., 2001. An automatic particle pickup method using a neural network applicable to low-contrast electron micrographs. *J. Struct. Biol.* 136, 227–238.
- Ogura, T., Sato, C., 2003. Automatic particle pickup method using a neural network has high accuracy by applying an initial weight derived from eigenimages: a new reference free method for single-particle analysis. *J. Struct. Biol.* (in press).
- Pascual-Montano, A., Donate, L.E., Valle, M., Barcena, M., Pascual-Marqui, R.D., Carazo, J.M., 2001. A novel neural network technique for analysis and classification of EM single-particle images. *J. Struct. Biol.* 133, 233–245.
- Radermacher, M., Rao, V., Grassucci, R., Frank, J., Timmerman, A.P., Fleischer, S., Wagenknecht, T., 1994. Cryo-electron microscopy and three-dimensional reconstruction of the calcium release channel/ryanodine receptor from skeletal muscle. *J. Cell Biol.* 127, 411–423.
- Radermacher, M., Ruiz, T., Wiczorek, H., Gruber, G., 2001. The structure of the V<sub>1</sub>-ATPase determined by three-dimensional electron microscopy of single particles. *J. Struct. Biol.* 135, 26–37.
- Sato, C., Sato, M., Iwasaki, A., Doi, T., Engel, A., 1998. The sodium channel has four domains surrounding a central pore. *J. Struct. Biol.* 121, 314–325.
- Sato, C., Ueno, Y., Asai, K., Takahashi, K., Sato, M., Engel, A., Fujiyoshi, Y., 2001. The voltage-sensitive sodium channel is a bell-shaped molecule with several cavities. *Nature* 409, 1047–1051.
- Serysheva, I.I., Orlova, E.V., Chiu, W., Sherman, M.B., Hamilton, S.L., van Heel, M., 1995. Electron cryomicroscopy and angular reconstitution used to visualize the skeletal muscle calcium release channel. *Nat. Struct. Biol.* 2, 18–24.
- van Heel, M., 1984. Multivariate statistical classification of noisy images (randomly oriented biological macromolecules). *Ultramicroscopy* 13, 165–183.
- van Heel, M., Frank, J., 1981. Use of multivariate statistics in analysing the images of biological macromolecules. *Ultramicroscopy* 6, 187–194.
- van Heel, M., Gowen, B., Matadeen, R., Orlova, E.V., Finn, R., Pape, T., Cohen, D., Stark, H., Schmidt, R., Schatz, M., Patwardhan, A., 2000. Single-particle electron cryo-microscopy: towards atomic resolution. *Quart. Rev. Biophys.* 33, 307–369.
- van Heel, M., Harauz, G., Orlova, E.V., Schmidt, R., Schatz, M., 1996. A new generation of the IMAGIC image processing system. *J. Struct. Biol.* 116, 17–24.
- Van Hulle, M.M., 1999. Faithful representations with topographic maps. *Neural Networks* 12, 803–823.
- Wriggers, W., Milligan, R.A., McCammon, J.A., 1999. Situs: a package for docking crystal structures into low-resolution maps from electron microscopy. *J. Struct. Biol.* 125, 185–195.
- Wriggers, W., Milligan, R.A., Schulten, K., McCammon, J.A., 1998. Self-organizing neural networks bridge the biomolecular resolution gap. *J. Mol. Biol.* 284, 1247–1254.