# Energy Minimization of Protein Tertiary Structure by Parallel Simulated Annealing using Genetic Crossover

Tomoyuki Hiroyasu[†], Mitsunori Miki[†], Shinya Ogura[††],
Keiko Aoi[††], Takeshi Yoshida[††], Yuko Okamoto [‡] Jack Dongarra [‡‡]

[†] Department of Engineering, Doshisha University@
[††] Graduate School of Engineering, Doshisha University@
[‡] Department of Theoretical Studies, Institute for Molecular Science@
[‡‡] Department of Computer Science, University of Tennessee@

From our recent research, it has been determined that Parallel Simulated Annealing using Genetic Crossover (PSA/GAc) has a high searching capability on an energy minimization of a small protein called Met-enkephalin. PSA/GAc performs genetic crossover (one of the operations of Genetic Algorithm (GA)) among the Parallel SA (PSA) to exchange their information. In this paper, PSA/GAc is applied to the energy minimization of C-peptide and Parathyroid Hormone Fragment(1-34). Also among the parameters of PSA/GAc, crossover interval and total number of searching steps, which are supposed to have greater influence on the searching capability, are modified to some values in order to examine and study their influences. The results show that PSA/GAc provides lower energy of the target proteins than Parallel SA. Furthermore, higher frequency of crossover and longer searching steps are proven to derive lower energies. From the results we conclude that PSA/GAc is also effective on the predictions of larger proteins.

*keywords*: Prediction of Protein Tertiary Structure, Simulated Annealing, Genetic Algorithm, Optimization

## 1  Introduction

Protein is a substance directly connected to the biological phenomena and its biological functions are said to be derived from its tertiary structure. Thus, clarification of its tertiary structure leads to an explanation of the biological phenomena process. Tertiary structure of the protein is believed to correspond to the conformation with the lowest residual potential energy. Therefore, as one of the methods to predict the tertiary structure of proteins, minimization of the energy function has been studied. Simulated Annealing (SA) has often been employed as the optimization method to predict the tertiary structure of proteins by the minimization of the energy[4] .

However, the energy function that determines the tertiary structure is complicated with large numbers of global and local minima. Therefore in order to acquire a fast convergence to a global optimal, a hybrid method that combines Sequential SA (SSA) with a mechanism to improve the searching capability has been developed. Thus, as a combination method of SA with an operator of Genetic Algorithm (GA), Parallel Simulated Annealing using Genetic Crossover (PSA/GAc) is what we propose[3] . GA is another optimization method that imitates an evolution of living things. GA is said to be superior to SA at searching for the optimum in a wider area of the searching domain.

Our previous research shows that PSA/GAc has higher searching capability than SSA on the energy minimization of a small protein called Met-enkephalin[3] . In this paper, PSA/GAc is applied to the energy minimization of the proteins larger than Met-enkephalin; C-peptide and Parathyroid Hormone Fragment (1-34). Also among the parameters of PSA/GAc, crossover interval and total number of searching steps, which are supposed to have greater influence on the searching capability, are modified to some values to examine and study their influences. The results of PSA/GAc are compared to those of Parallel SA (PSA).

## 2   SA and GA

### 2.1   Simulated Annealing

The term annealing derives from the physical process of heating and then slowly cooling a substance to obtain a strong crystalline structure. Simulated Annealing (SA) is the optimization method to stochastically simulate this physical process of annealing on the computers. In SA, the simulation proceeds by randomly generating a solution and then determining its acceptance in certain probability. A temperature parameter is used to determine this probability. In the basic algorithm of SA, three operations bear an important role; generate, accept, and cool.

The generation operation modifies a current solution $x$ and generates a next solution $x'$ using a probability distribution $G(x, x')$.

The accept operation is a judgment to decide whether to accept the modification or not. The acceptance of the modification is determined from a difference $\Delta E(= E' - E)$ of a current energy $E = f(x)$ and a modified energy $E' = f(x')$, and the temperature parameter $T$. Metropolis et al.[6] introduced a simple algorithm to provide an efficient simulation. That is, if $\Delta E \leq 0$, the modification is accepted. In case $\Delta E > 0$, the modification is accepted at certain probability. This algorithm can be defined as:

$$P_{ACCEPT} = \begin{cases} 1 & if \quad \Delta E \leq 0 \\ exp(-\dfrac{\Delta E}{T}) & otherwise \end{cases} \quad (1)$$

The temperature $T$ is the important parameter to control the acceptance of the modified solution. At the beginning of the simulation, both the temperature and the acceptance levels are high. As the simulation proceeds and the temperature decreases, solutions that have the bigger fitness values.

Cooling is the operation to generate the temperature of the next state $T_{k+1}$ from the temperature of the current state $T_k$.

The simulation begins with the initial solution $x$ with energy of $E$ and the initial temperature $T$. The solution is then randomly modified to $x'$ with energy of $E'$. The acceptance of the modification is calculated from the difference of energy, $\Delta E(= E' - E)$, and the temperature $T_k$.

If accepted, the solution $x'$ becomes the starting point of the next step. These operations are repeated long enough at each temperature for the system to reach a steady state, or equilibrium. When reaching the steady state at $T_k$, the temperature is cooled to $T_{k+1}$ and the simulation is repeated until reaching steady state again. Users define the terminal criterion, such as the number of evaluations. The simulation is concluded when the temperature becomes low enough and a terminal criterion is met.

### 2.2   Genetic Algorithm

Genetic Algorithm (GA) is the optimization algorithm that imitates the evolution of living creatures. In nature, inadaptable creatures to an environment meet extinction, and only adapted creatures can survive and reproduce. A repetition of this natural selection spreads the superior genes to conspecifics and then the species prospers. GA models this process of nature on computers.

GA can be applied to several types of optimization problems by encoding design variables of individuals. Searching for the solution proceeds by performing the three genetic operations on the individuals; selection, crossover, and mutation, which play an important role in GA.

Selection is an operation that imitates the survival of the fittest in nature. The individuals are selected for the next generation according to their fitness.

Crossover is an operation that imitates the reproduction of living creatures. The crossover exchanges the information of the chromosomes among individuals.

Mutation is an operation that imitates the failure that occurs when copying the information of DNA. Mutating the individuals in a proper probability maintains the diversity of the population.

## 3   Parallel SA using Genetic Crossover

Parallel Simulated Annealing using Genetic Crossover (PSA/GAc) is the optimization method to exchange the informations of solutions among SAs running in parallel by the genetic

crossover[2] . In this algorithm, the total number of SAs running in parallel is defined as population size and each SA is defined as individual[2] .

Procedure of PSA/GAc is stated below:

**step 1** SAs start running in parallel with initial solutions.

**step 2** When number of annealing reaches a certain number $d$, two solutions (individuals) from the parallel SAs are randomly paired.

**step 3** For each pair of individuals, the crossover operation is performed and two new individuals are generated. The crossover operation used in PSA/GAc is one-point crossover between the design variables.

**step 4** Among the two parents and the two offsprings, the two individuals with the better fitness are selected as new searching points.

**step 5** A certain number $d$ of annealing is performed from the new searching points.

**step 6** Each pair of individuals performs step 3 to step 5.

**step 7** Step 2 to step 6 are repeated until the terminal criterion is met.

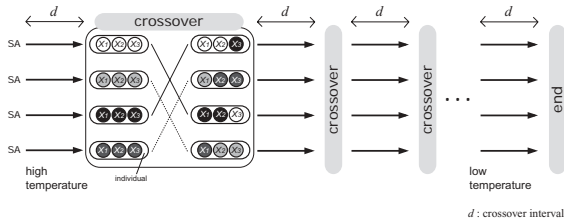A concept of PSA/GAc is illustrated in Figure 1.



Figure 1: Concept of Parallel Simulated Annealing using Genetic Crossover

Figure 2 describes the concept of steps 3 and 4. The figure describes the case of three design variables where the randomly selected boundary between the first and the second variables becomes a crossover point. After evaluating the individuals, the individuals with higher evaluations (parent 2 and offspring 2) are chosen as the next searching points.
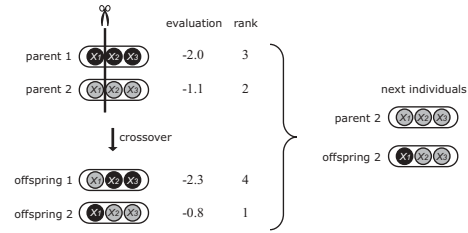


Figure 2: Concept of crossover and selection

In case the optimal values of some design variables are obtained, the crossover operation works to transfer the values to other individuals. Thus, the convergence of the annealing is precipitated.

The former experiments of applying PSA/GAc on some continuous minimization problems show that PSA/GAc has a better convergence property than SSA, and on the problems difficult for GA to solve, PSA/GAc has also been confirmed to be effective[2] .

# 4 Energy Minimization of Proteins by PSA/GAc

In this research, PSA/GAc is applied to the energy minimization of larger proteins to verify its performance and also to study some of its parameters.

## 4.1 Target Proteins

The target proteins of this research are C-peptide and Parathyroid Hormone Fragment(1-34). C-peptide consists of 13 amino residues, has 13 pairs of dihedral angles $(\phi, \psi)$ in the backbone, and 38 dihedral angles in the side chains. PTH(1-34) consists of 34 amino residues, has 34 pairs of dihedral angles in the backbone, and 110 dihedral angles in the side chains.

The energy minimization of these proteins is gas-phase simulation based on the energy parameters of ECEPP/2[7, 8, 12] . The dihedral angles of backbone and side chains are applied as design variables. Values of the dihedral angles are in the range of [-180°, 180°]. Each dihedral angle is generated and given the accept criterion sequentially, and then the temperature is cooled. In this paper, this series of operations are defined as a Monte Carlo Sweep (MCsweep). Thus, 64

Metropolis criterions are performed in one MC-sweep in the case of C-peptide and 178 in the case of PTH(1-34).

Okamoto's experiments, using the energy function that adopts ECEPP/2 energy parameters, show that C-peptide has the lowest energy conformation when eight residues (4-11) form an $\alpha$-helix. The energy of this lowest energy conformation was -42.0 $kcal/mol$[9] . Also, the lowest energy conformation of C-peptide obtained in this experiment corresponds with the structure deduced from the X-ray crystallography experiment[1] .

It is known by the NMR experiment that PTH(1-34) has two $\alpha$-helices[5] . The lowest energy conformation of PTH(1-34) deduced from Okamoto's experiment also had two $\alpha$-helices and its energy is -210.0 $kcal/mol$[11] .

## 4.2 Energy Minimization of C-peptide

In the energy minimization of C-peptide, Okamoto made 20 runs of SSA with each run consisting of 10,000 MCsweeps[9] . In order to equalize the total number of MCsweeps with Okamoto's experiment, the parameters of PSA/GAc were set to 24 individuals × 4,165 MCsweeps and two runs were made. All the experiments were started from random initial conformations. Next, states of the dihedral angles are generated stochastically from the neighborhood range using uniform distribution. Neighborhood range $[max, min]$ is determined as:

$$\begin{cases} max = 180° - \dfrac{180° \times 0.7 \times \#sweep}{Total \ \# \ sweeps} \\ min = -max \end{cases} \quad (2)$$

Other parameters are shown in Table 1. The initial and the final temperature are the same values as Okamoto's[10] .

The result of the experiment signified that the lowest energy conformation was obtained when the crossover interval was set to 32. Table 2 lists the main-chain dihedral angles of the conformation with the lowest energy, which was -53.9 $kcal/mol$. The conformation is considered as $\alpha$-helical when the dihedral angles $(\phi, \psi)$ fall in the range $(-60 \pm 45°, -50 \pm 45°)$[9] , therefore,

the lowest energy conformation exhibits $\alpha$-helices in nine residues. The $\alpha$-helical residues are indicated with **a** in Table 2. Figure 3 shows the lowest energy conformation of C-peptide obtained in the experiment. The lowest energy conformation derived from the experiment had a lower energy than that of Okamoto and the $\alpha$-helices were in accord with the known conformation. Consequently, PSA/GAc has high convergence property on the energy minimization of C-peptide.

Table 1: Parameters of PSA/GAc

| Parameter | Value |
| --- | --- |
| Initial Temperature | 2.0 (1000K) |
| Final Temperature | 0.1 (50K) |
| Crossover Interval | 8, 16, 32, 64 |

Table 2: Amino acid sequence of C-peptide and the dihedral angles of the lowest energy conformation

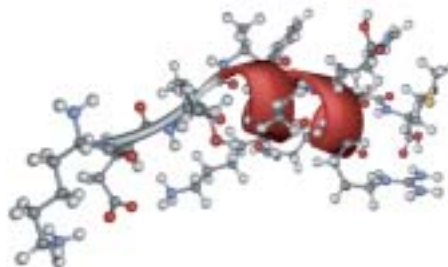| Sequence | Ly+ | Gl- | Thr | Ala | 5<br>Ala |
| --- | --- | --- | --- | --- | --- |
| | 23 | -79 | -76 | -67 | -66 |
| | -66 | 102 | 86 | -27 | -32 |
| | - | - | - | a | a |
| | Ala | Ly+ | Phe | Glu | 10<br>Ar+ |
| | -79 | -62 | -65 | -65 | -63 |
| | -40 | -43 | -41 | -41 | -41 |
| | a | a | a | a | a |
| | 13<br>Gln | Hi+ | Met | | |
| | -72 | -69 | -82 | | |
| | -30 | -40 | 105 | | |
| | a | a | - | | |

$Energy = -53.9 \ kcal/mol$



Figure 3: Lowest energy conformation of C-peptide obtained by PSA/GAc

## 4.3 Energy Minimization of PTH(1-34)

In the energy minimization of PTH(1-34), Okamoto made 20 runs of SSA with each run consisting of 10,000 MCsweeps[11] . In order to equalize the total number of MCsweeps with Okamoto's experiment, the parameters of PSA/GAc was set to 24 individuals $\times$ 4,165 MCsweeps, and two runs were made. All the experiments were started from random initial conformations. Other parameters are shown in Table 1.

The result of the experiment signified that the lowest energy conformation was obtained when the crossover interval was set to eight. Table 3 lists the main-chain dihedral angles of the conformation with lowest energy, which was -236.6 $kcal/mol$. The lowest energy conformation exhibits $\alpha$-helices in residues 2-8 and 15-18. Figure 4 shows the lowest energy conformation of PTH(1-34) obtained in the experiment.
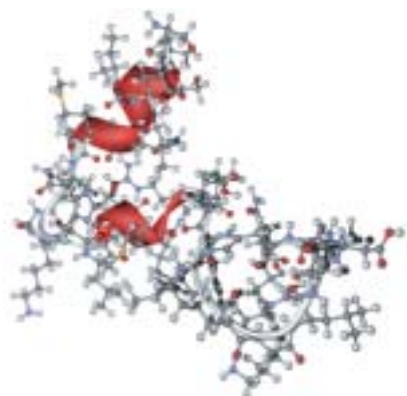


Figure 4: Lowest energy conformation of PTH(1-34) obtained by PSA/GAc

As stated in section 4.1, PTH(1-34) is known to have two $\alpha$-helices from both NMR and Okamoto's experiment. The lowest energy conformation obtained by PSA/GAc also exhibited two $\alpha$-helices and retained lower energy than the Okamoto's. Consequently, PSA/GAc has high convergence property on the energy minimization of PTH(1-34).

Table 3: Amino acid sequence of PTH(1-34) and the dihedral angles of the lowest energy conformation

| Sequence | Ser | Val | Ser | Glu | 5 Ile |
|---|---|---|---|---|---|
| | 93 | -65 | -68 | -76 | -71 |
| | 157 | -26 | -36 | -37 | -35 |
| | - | a | a | a | a |

| | Gln | Leu | Met | His | 10 Asn |
|---|---|---|---|---|---|
| | -64 | -70 | -55 | -64 | -112 |
| | -41 | -44 | -43 | -96 | 145 |
| | a | a | a | - | - |

| | Leu | Gly | Lys | His | 15 Leu |
|---|---|---|---|---|---|
| | -105 | 76 | -95 | -90 | -73 |
| | 108 | -90 | 165 | -1 | -22 |
| | - | - | - | - | a |

| | Asn | Ser | Met | Glu | 20 Ar+ |
|---|---|---|---|---|---|
| | -85 | -85 | -74 | -88 | -142 |
| | -20 | -37 | -30 | 64 | -60 |
| | a | a | a | - | - |

| | Val | Glu | Trp | Leu | 25 Ar+ |
|---|---|---|---|---|---|
| | -134 | -85 | -64 | -65 | -155 |
| | -58 | 73 | 131 | 108 | 112 |
| | - | - | - | - | - |

| | Lys | Lys | Leu | Gln | 30 Asp |
|---|---|---|---|---|---|
| | -113 | -105 | -148 | -64 | -65 |
| | 63 | -32 | 148 | -55 | 145 |
| | - | - | - | a | - |

| | Val | His | Asn | Phe | 34 |
|---|---|---|---|---|---|
| | 48 | -55 | -60 | -84 | |
| | 71 | -46 | 153 | 118 | |
| | - | a | - | - | |

$Energy = -236.6 \ kcal/mol$

## 4.4 Study on the Number of MCsweeps and Crossover

### 4.4.1 Abstract of the Experiment

PSA/GAc performs genetic crossover, one of the operations of GA, to transfer the solutions among SAs running in parallel. In case the optimal values of some design variables are obtained, the crossover operation works to transfer the values to other individuals. Thus, the convergence of the annealing can be precipitated. Therefore the frequency of the crossover operation gives significant influence on the performance of PSA/GAc.
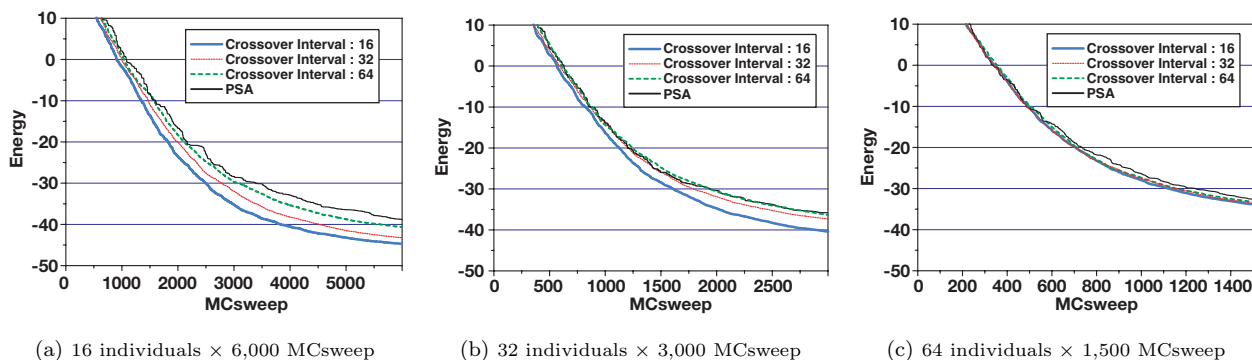
(a) 16 individuals × 6,000 MCsweep  (b) 32 individuals × 3,000 MCsweep  (c) 64 individuals × 1,500 MCsweep

Figure 5: Energy history of C-peptide. Each figure presents the result of PSA and PSA/GAc with crossover intervals ranging from 16 to 64.



(a) 16 individuals × 6,000 MCsweep  (b) 32 individuals × 3,000 MCsweep  (c) 64 individuals × 1,500 MCsweep
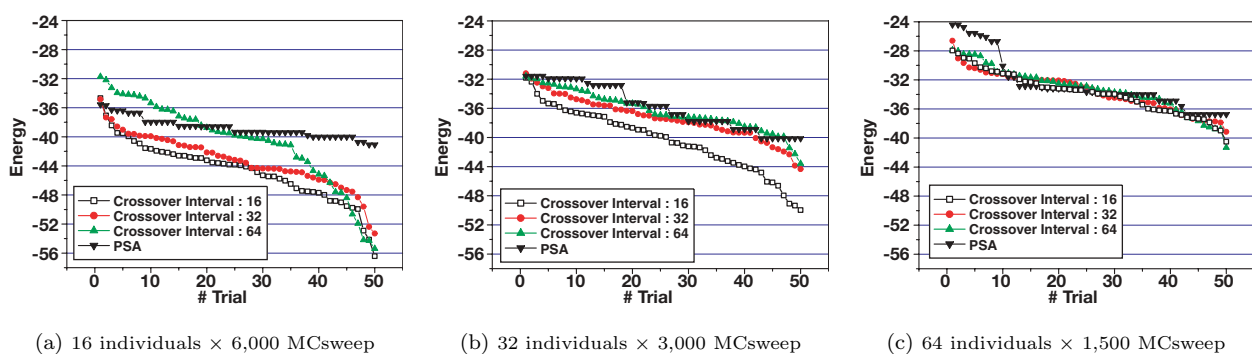
Figure 6: Lowest energies of C-peptide obtained in each run are plotted in the descending order. Each figure presents the result of PSA and PSA/GAc with crossover intervals ranging from 16 to 64.

Thus, in this section, the influence of the total number of MCsweeps and the frequency of the crossover on the convergence property of PSA/GAc are studied. In the experiment, parameters of PSA/GAc are set to 16 individuals × 6,000 MCsweeps, 32 individuals × 3,000 MCsweeps, and 64 individuals × 1,500 MCsweeps with crossover intervals of 16, 32, and 64. Among PSA/GAcs, Parallel SA (PSA) is also applied. PSA is the method to fetch the best solution from plural SSAs running in parallel.

Total number of evaluations in every run of the experiment is approximately united.

### 4.4.2  Result of C-peptide

The results of the energy minimization of C-peptide by PSA/GAc with parameters of 16 individuals × 6,000 MCsweep, 32 individuals × 3,000

MCsweep, and 64 individuals × 1,500 MCsweep are summarized in Figure 5(a), Figure 5(b), and Figure 5(c), respectively which show the best of the individuals as the number of MCsweeps proceeds. The figures show the average of 50 runs. Each figure also shows the result with crossover intervals of 16, 32, and 64. The horizontal and the vertical axes each represent the number of proceeded MCsweeps and the energy of protein ($kcal/mol$).

Figure 6(a), Figure 6(b), and Figure 6(c) present the lowest energies obtained in each run of PSA and PSA/GAc with crossover intervals of 16, 32, and 64, respectively. Each energy is plotted in the descending order. The horizontal and the vertical axes each represent the number of runs and the energy of proteins($kcal/mol$).
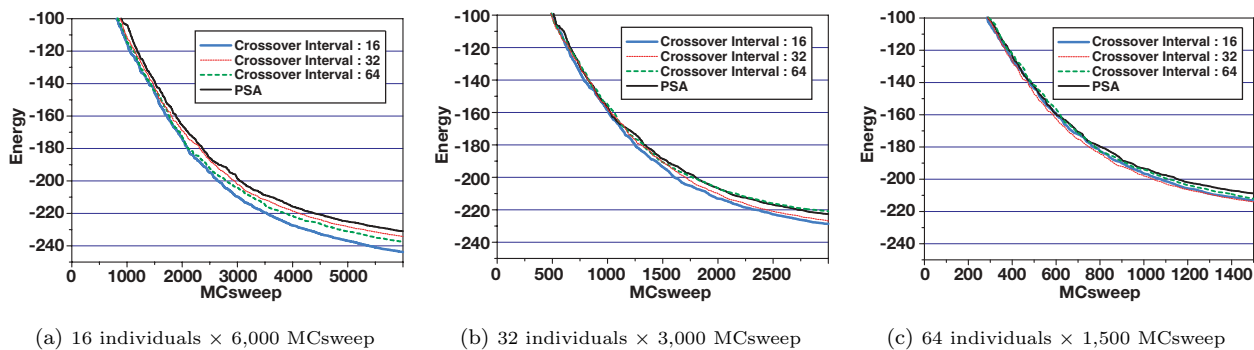
6

(a) 16 individuals × 6,000 MCsweep     (b) 32 individuals × 3,000 MCsweep     (c) 64 individuals × 1,500 MCsweep

Figure 7: Energy history of PTH(1-34). Each figure presents the result of PSA and PSA/GAc with crossover intervals ranging from 16 to 64.



(a) 16 individuals × 6,000 MCsweep     (b) 32 individuals × 3,000 MCsweep     (c) 64 individuals × 1,500 MCsweep
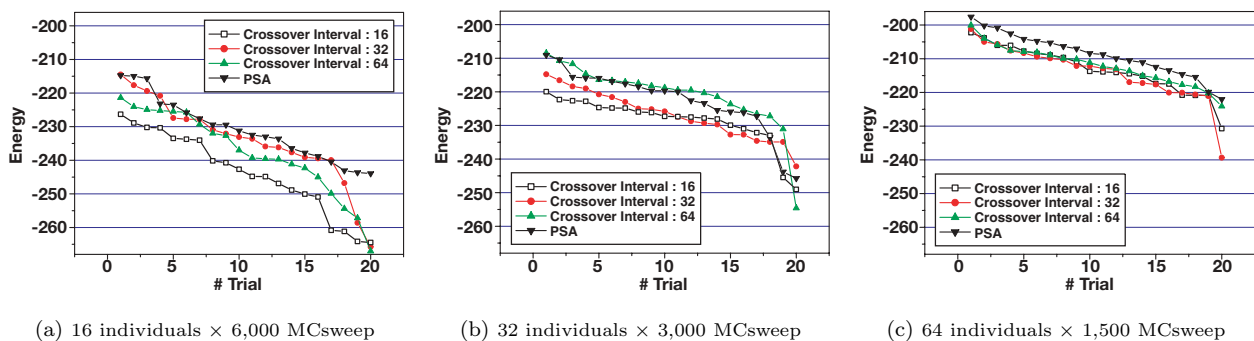
Figure 8: Lowest energies of PTH(1-34) obtained in each run are plotted in descending order. Each figure presents the result of PSA and PSA/GAc with crossover intervals ranging from 16 to 64.

### 4.4.3   Result of PTH(1-34)

The result for the energy minimization of PTH(1-34) by PSA/GAc with parameters of 16 individuals × 6,000 MCsweep, 32 individuals × 3,000 MCsweep, and 64 individuals × 1,500 MCsweep are summarized in Figure 7(a), Figure 7(b), and Figure 7(c), respectively which show the best of the individuals as the number of MCsweeps proceeds. The figures show the average of 20 runs. Each figure also shows the result with crossover intervals of 16, 32, and 64. The horizontal and the vertical axes each represent the number of proceeded MCsweeps and the energy of proteins ($kcal/mol$).

Figure 8(a), Figure 8(b), and Figure 8(c) present the lowest energies obtained in each run of PSA and PSA/GAc with crossover interval of 16, 32, and 64. Each energy is plotted in the de-

scending order. The horizontal and the vertical axes each represent the number of runs and the energy of proteins ($kcal/mol$).

### 4.4.4   Result of Met-enkephalin

Results for the energy minimization of small protein consisted of five residues called Met-enkephalin are shown below for comparison.

Figure 9(a) presents the transition of energy obtained by PSA/GAc with the parameters of 16 individuals × 6,000 MCsweeps, respectively. The transition is an average of 50 runs. Figure 9(b) and Figure 9(c) each present the lowest energy obtained in 50 runs of PSA/GAc with the parameters of 16 individuals × 6,400 MCsweeps and 64 individuals × 1,500 MCsweeps. The energies are plotted in the descending order.

(a) Transition of energies obtained by PSA/GAc with 16 individuals × 6,000 MCsweeps

(b) The lowest energies of each runs obtained by PSA/GAc with 16 individuals × 6,000 MCsweeps

(c) The lowest energies of each runs obtained by PSA/GAc with 64 individuals × 1,500 MCsweeps
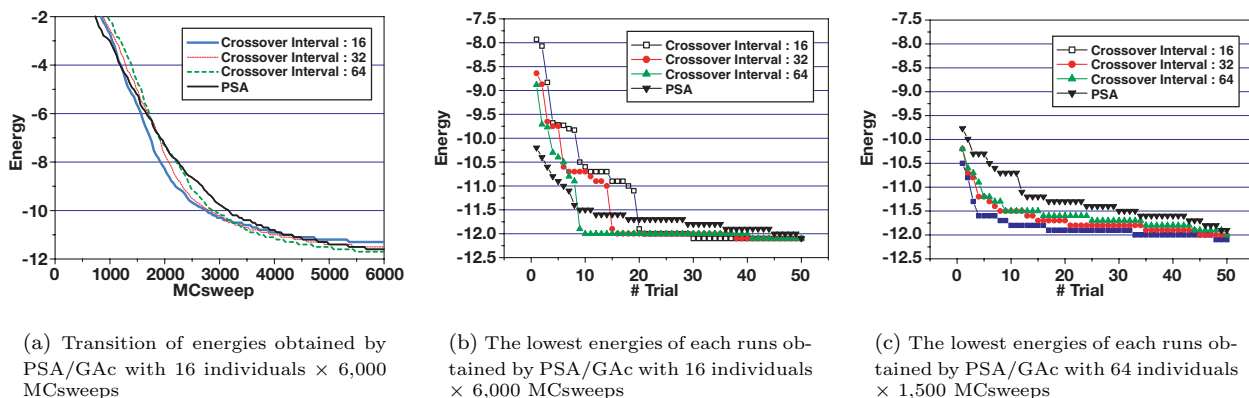
Figure 9: Transition of the energy and the lowest energies obtained in the energy minimization of Met-enkephalin.

### 4.4.5   Discussion of the Results

The following observations can be made from the results of both C-peptide and PTH(1-34).

Firstly, we made a comparison of PSA/GAc with the large and small total number of MCsweeps. In each case PSA/GAc with the larger number of MCsweeps acquired better convergence. This is because the total number of MCsweeps is linked with the cooling schedule; thus, PSA/GAc with longer MCsweeps could search in both global and local areas of the searching domain sufficiently. On the other hand, Figure 9(c) reveals that PSA/GAc with 64 individuals × 1,500 MCsweeps has the best convergence property in the energy minimization of Met-enkephalin. Because Met-enkephalin is a small protein with 19 dihedral angles, it can be considered to be simple as an optimization problem. Therefore, the small number of MCsweeps and the short crossover interval are sufficient for the search on the energy minimization of Met-enkephalin. On the other hand, for the energy minimization of larger proteins such as C-peptide and PHT(1-34), larger number of MCsweeps is required to obtain enough convergence.

Secondly, we focused on the effect of the crossover operation. By comparison the results of PSA/GAc with those of PSA that does not have the cross over operation, it can be said that PSA/GAc derived better solutions. Thus, crossover operation is working effectively on the

energy minimization of these proteins.

Comparison of the crossover interval also revealed that the convergence property of PSA/GAc is improved by the crossover operation. The smaller the crossover interval is (the higher frequency of the crossover operation is) in either results the faster the convergence becomes. When the total number of MCsweeps is large, this trend is prominent.

In the case of the energy minimization of Met-enkephalin, it is obvious from Figure 9(a) that the crossover interval has almost no influence on the convergence property of PSA/GAc in spite of the total number of MCsweeps. Figure 9(b) indicates that the performance of PSA/GAc with a crossover interval of 16 is inferior to PSA/GAc with other crossover intervals. That is, when the total number of MCsweeps for the energy minimization of small proteins is unnecessarily large, setting a short crossover interval yields a decrement in the diversity of individuals, and therefore, the performance decreases. From the discussion, it can be speculated that PSA/GAc and the crossover operation are especially effective for the energy minimization of large proteins that require the large number of MCsweeps.

## 5   Conclusions

In this paper, Parallel Simulated Annealing using Genetic Crossover (PSA/GAc) has confirmed its effectiveness on the energy minimization of the

small protein called Met-enkephalin and applied to the energy minimization of larger proteins; C-peptide and PTH(1-34).

Firstly, to compare the performance of PSA/GAc with Okamoto's Sequential SA (SSA), PSA/GAc was applied to the energy minimization of C-peptide and PTH(1-34). The result demonstrated that PSA/GAc has a better convergence property than SSA on the energy minimization of both proteins.

Secondly, PSA/GAc and Parallel SA (PSA) were applied to the energy minimization of above proteins to make a comparison of the performances. (Number of individuals) × (total number of MCsweeps) was set to constant (total number of evaluation is approximately equal) to study the influence of total numbers of MCsweeps and the crossover operation. The result revealed that PSA/GAc obtains lower energy than PSA. In addition, the result also indicates that the shorter the crossover interval is, the more frequently the lower energies were obtained.

Studying the above results, we can say that the crossover operation of PSA/GAs is working effectively on the energy minimization of proteins. Additionally, it can be concluded that PSA/GAc is also effective on the energy minimization of larger proteins.

# References

1) Ulrich H. E. Hansmann and Yuko Okamoto. Tertiary Structure Prediction of C-Peptide of Ribonuclease A by Multicanonical Algorithm. *J. Phys. Chem. B*, 102(4):653–656, 1998.

2) Tomoyuki Hiroyasu, Mitsunori Miki, and Maki Ogura. Parallel simulated annealing using genetic crossover. *Proceedings of the 44th Institute of Systems, Control and Information Engineers*, pages 113–114, 2000.

3) Tomoyuki HIROYASU, Mitsunori MIKI, Maki OGURA, and Yuko OKAMOTO. Examination of Parallel Simulated Annealing using Genetic Crossover. *Journal of Information Processing Society of Japan*, 43(SIG7(TOM6)):70–79, 2002.

4) Hikaru Kawai, Takeshi Kikuchi, and Yuko Okamoto. A prediction of tertiary structures of peptide by the Monte Carlo simulated annealing method. *Protein Engineering*, 3(2):85–94, 1989.

5) W. Klaus, T. Dieckmann, V. Wray, D. Schomburg, E. Wingender, and H. Mayer. *Biochemistry*, 30:6936–6942, 1991.

6) N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. Equation of State Calculations by Fast Computing Machines. *The Jurnal of Chemical Physics*, 21(6):1087–1092, 1953.

7) F.A. Momany, R.F. McGuire, A.W. Burgess, and H.A. Scheraga. *J. Phys. Chem.*, 79:2361–2381, 1975.

8) G. Nemethy, M.S. Pottle, and H.A. Scheraga. *J. Phys. Chem.*, 87:1883–1887, 1983.

9) Yuko Okamoto, Masataka Fukugita, Takashi Nakazawa, and Hikaru Kawai. $\alpha$-Helix folding by Monte Carlo simulated annealing in isolated C-peptide of ribonuclease A. *Protein Engineering*, 4(6):639–647, 1991.

10) Yuko Okamoto, Takeshi Kikuchi, and Hikaru Kawai. Prediction of Low-Energy Structures of Met-Enkephalin by Monte Carlo Simulated Annealing. *CHEMISTRY LETTERS*, pages 1275–1278, 1992.

11) Yuko Okamoto, Takeshi Kikuchi, Takashi Nakazawa, and Hikaru Kawai. $\alpha$-Helix structure of parathyroid hormone fragment (1-34) predicted by Monte Carlo simulated annealing. *INTERNATIONAL JOURNAL OF PEPTIDE & PROTEIN RESEARCH*, 42:300–303, 1993.

12) M.J. Sippl, G. Nemethy, and H.A. Scheraga. *J. Phys. Chem.*, 88:6231–6233, 1984.