

DFS generated pathways in GA crossover for protein structure prediction

Author

Hoque, Md Tamjidul, Chetty, Madhu, Lewis, Andrew, Sattar, Abdul, Avery, Vicky M

Published

2010

Journal Title

Neurocomputing

DOI

[10.1016/j.neucom.2010.02.021](https://doi.org/10.1016/j.neucom.2010.02.021)

Rights statement

© 2010 Elsevier B.V. This is the author-manuscript version of this paper. Reproduced in accordance with the copyright policy of the publisher. Please refer to the journal's website for access to the definitive, published version.

Downloaded from

<http://hdl.handle.net/10072/35904>

Griffith Research Online

<https://research-repository.griffith.edu.au>

DFS generated pathways in GA crossover for protein structure prediction

Md Tamjidul Hoque^{a, b, *}, Madhu Chetty^c, Andrew Lewis^b, Abdul Sattar^b and Vicky M Avery^a

^aDiscovery Biology, Esikitis Institute for Cell & Molecular Therapies and ^bInstitute for Integrated and Intelligent Systems (IIIS), Griffith University, Nathan QLD 4111, Australia.

^cGippsland School of Information Technology (GSIT), Monash University, Churchill VIC 3842, Australia.

Elsevier use only: Received date here; revised date here; accepted date here

Abstract

Genetic Algorithms (GAs), as nondeterministic conformational search techniques, are promising for solving protein structure prediction (PSP) problems. The crossover operator of a GA can underpin the formation of potential conformations by exchanging and sharing potential sub-conformations. However, as the optimum PSP conformation is usually compact, the crossover operation may result in many invalid conformations (by having non-self-avoiding-walk). Although a crossover-based converging conformation suffers from limited pathways, combining it with depth-first search (DFS) can partially reveal potential pathways and make an invalid crossover valid and successful. Random conformations are frequently applied for maintaining diversity as well as for initialization in many GA applications. Random-move-only-based conformation generator has exponential time complexity in generating random conformations, whereas the DFS based random conformation generator has linear time complexity and performs relatively faster. We have done extensive experiments using popular 2D as well as useful 3D models to justify our hypothesis empirically.

Keywords: Depth-first search; protein structure prediction; genetic algorithm; HP lattice; FCC model;

1. Introduction

The protein structure prediction (PSP) is a problem of determining the native state of a protein from its primary structure and is of great importance because three dimensionally folded structures determine the biological function [1] and hence proves very useful in applications such as drug design [2]. Particular folded structures are essential for the functioning of living cells as well as for providing body structure. High-resolution protein modeling is possible, provided a homologue of the target protein exists [3]. The application of the high-resolution model becomes less effective without a homologous template. However, the homologous template is unable to explain “how and why” a protein adopts a specific structure. Thus, low resolution model based *ab initio* [4] (meaning ‘from

the origin’) or the *de novo* becomes essential. In an *ab initio* approach, the building of a 3D conformation (structure) is essentially based on the properties of amino acids, since protein is a three dimensionally folded molecule composed of amino acids [5] linked together (called the primary structure) in a particular order specified by the DNA sequence of a gene [6]. In this paper, our efforts are to investigate the *ab initio* protein structure prediction problem.

Lattice protein models introduced by Dill [7] are widely used for investigating the underlying principles of protein folding [8]. Protein conformation as a *self-avoiding walk* in the lattice model has been proven [9, 10] to be *NP-complete* for the 2D square and 3D cube HP models. Therefore, a deterministic algorithm for folding prediction is not feasible. Reasonably a nondeterministic approach with robust strategies that can extract minimal energy conformations efficiently from these models becomes

necessary. Nevertheless, this is a very challenging task as there exists an astronomical number of possible conformations even for a very short sequence of amino acids [11, 12].

Due to its superior performance, Genetic Algorithm (GA) having crossover as one of its key operation [13], is often chosen as a vehicle for providing solutions to the PSP problems. Not only within GA itself, but also in many PSP solving algorithms, the core concepts of GAs and their components are often adapted for effectiveness [14-18]. While crossover can be very effective in joining two different potential sub-conformations, it can be repeatedly unsuccessful with the conformations (hence the sub-conformations) converging. This is because to the conformation being compact in nature, it leaves limited pathways available to a valid (i.e., self-avoiding-walk) conformation, thereby causing many potential conformations to be lost. This motivates us to apply partial pathways, based on depth first search (DFS) [19] to regain potential sub-conformations, leading to effective algorithms and superior conformations resulting better PSP solutions.

2. Background and preliminaries

In nature, a protein folds fast, requiring between a tenth of a millisecond and one second in general, whereas any algorithm on any modern computer is still unable to simulate this task in anything close to real-time folding [13, 20]. Current research confronted with immensely complex of the protein structure prediction problems, has lead to the manifestation of several important issues and approaches, which are yet to be investigated [13, 21, 22]:

First, the energy function, which is a combination of several factors that determines the free energy of a folded protein, is not fully understood. Therefore, existing formulations for energy functions do not suggest any obvious path to a solution for the PSP problem.

Second, conformational search algorithms are promising approaches for solving this hard optimization problem. However, the PSP problem still needs considerable research to find an effective algorithm. The aim of the search is to identify an optimum conformation within a very vast and convoluted search landscape.

Third, Cyrus Levinthal postulated, in what is popularly known as the Levinthal paradox, that proteins fold into their specific 3D conformations in a time-span far shorter than it would be possible for the molecule to actually search the entire conformational space (which is astronomically large) for the lowest energy state [23]. As proteins cannot, while folding, be sampling all possible conformations, folding pathways must therefore exist.

While focusing on the second issue [24-29], we are utilizing DFS strategies, developing novel search algorithms in a form to address the pathway hypothesis. Energy landscape of the protein folding pathways is very

convoluted, resembling a microscopic funnel energy landscape, where at any point of the surface the energy of conformation drops increasingly with decreasing search space. This reflects in structure prediction as a crossover failure because the converging congested conformations face more collisions producing invalid pathways, which fail to capture the microscopic pathways facing crests and troughs. Thus, it has been concluded that conformational searching is a major bottleneck in protein folding prediction with the observed folding rates have been found to be proportional to the number of microscopic folding routes [30]. The macroscopic routes can be captured by the crossover operation applied to the suboptimal conformations and then partially applying DFS can mimic the existing microscopic path guided by the converging sub-conformation. In contrast, a crossover operation alone can encounter more collisions [15] (while mating dissimilar converging conformations) before having a SAW conformation, and thus often can reject the potential sub-conformation as being unfit when paired with the available counterpart of the crossover portion (from a dissimilar conformation).

2.1. The HP lattice model

The simplified HP lattice model [31, 32] is based on *hydrophobicity* [33], dividing the amino acids into two different beads – *hydrophobic* (H) and *hydrophilic* (or *polar* (P)). The model allows HP protein sequences to be configured as self-avoiding walks (SAW) on the lattice path favoring an energy free state due to HH interaction. The energy of a given conformation is defined as the number of *topological neighbor* (TN) contacts between those Hs, which are not adjacent in the sequence. This contact between two neighboring H residues (or HH contact) is a TN and is assigned a value for the potential, termed *interaction potential* which is defined as -1 for the HP model. The value is chosen negative, since each protein's stable folded state is assumed corresponding to the global minimum free energy [34]. Further, the HP interaction and PP interaction potential value is assigned 0, which basically implies that there is no interaction between an H and a P of HP contact or between the Ps of PP contacts.

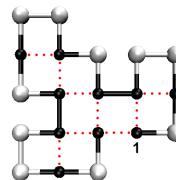


Fig. 1. HP conformation in the 2D HP square model shown by a solid line. 2D square lattice having fitness = - (TN Count) = -9. ● indicates a hydrophobic and ○ indicates a hydrophilic residue. The dotted line indicates a TN. Starting residue is indicated by a '1' in the figure.

To define PSP formally, assume for an amino-acid sequence $s = s_1, s_2, s_3, \dots, s_n$, a conformation c needs to be formed where $c^* \in C(s)$, $C(s)$ is the set of all valid (i.e., SAW) conformations of s , n is the total number of amino acids in the sequence and energy $E^* = E(C) = \min\{E(c) | c \in C\}$ [17]. If the number of TNs (for HH contact) in a conformation c is q then the value of $E(c)$ is defined as $E(c) = -1 \times q = -q$ and the *fitness function* is $F = -q$. The optimum conformation will have a maximum possible value of $|F|$. With respect to the configurations, the lattice model can be of many types, however the 2 different useful forms of configuration will be discussed in 2.1.1 and 2.1.2.

2.1.1. 2D square configuration

A 2D square lattice model (Figure 1) is popular being easy to implement within the research community [13, 25, 31, 32, 35-41], many search algorithms are developed using the 2D square HP lattice model. Hence any algorithm development in the same configuration allows validation and comparison of new techniques for protein structure prediction (PSP) [24-26, 28, 29, 42].

In a 2D HP square lattice model, a non-terminal and a terminal residue, each with 4 neighbours, can have a maximum of 2 TNs and 3 TNs, respectively and the model will have a parity problem. However, we use a 2D square HP lattice model to encourage interested readers to do further research and to compare and then we have extended the developed algorithm for the 3D face-centred-cube (FCC) HP configuration as described in 2.1.2.

2.1.2. 3D face-centered-cube (FCC) configuration

In a 3D FCC configuration, a residue can have 12 neighbors (see Figure 2(a) and 2(b)). As can be seen from Figure 2(b), the distributions of the neighbours follow the following arrangement: 6 are in the same plane, 3 are above in an upper plane and the remaining 3 are below in the lower plane. Any two adjacent residues are always a unit distance apart [38].

By placing the centre residue at the origin $(0, 0, 0)$, the 6 neighbours form a uniform hexagon with their x , y and z coordinates given as $(1, 0, 0)$, $(1/2, \sqrt{3}/2, 0)$, $(-1/2, \sqrt{3}/2, 0)$, $(-1, 0, 0)$, $(-1/2, -\sqrt{3}/2, 0)$ and $(1/2, -\sqrt{3}/2, 0)$. The coordinates of the upper layer, with 3 residues, are $(0, 1/2, \sqrt{3}/2)$, $(-\sqrt{3}/4, -1/4, \sqrt{3}/2)$ and $(\sqrt{3}/4, -1/4, \sqrt{3}/2)$. For the lower 3 residues the coordinates, therefore, will be $(0, -1/2, -\sqrt{3}/2)$, $(\sqrt{3}/4, 1/4, -\sqrt{3}/2)$ and $(-\sqrt{3}/4, 1/4, -\sqrt{3}/2)$.

The algorithms developed for 2D HP model are extended to the 3D FCC model in this paper. This is done to handle the immense complexity to solve PSP and the realistic sample generated by the 3D FCC model can be fed into a high-resolution model in a hierarchical manner [43]. The choice for the 3D FCC model is based on the following reasons [44]:

i) 3D FCC is the densest sphere packing configuration [45]. Thus, it can provide the most compact or densest protein core [38, 39] in PSP. However, the protein core may not necessarily be the most compact one. The 3D FCC configuration is also parity problem free [46].

ii) Following *(i)*, for a region of fixed volume of space, logically inferring, FCC can offer highest degree of freedom for placing a residue in a suitable neighboring position.

iii) Following *(ii)*, FCC can provide the most perfect discrete mapping of the real folded protein, thus can be a perfect model configuration at the top of the high-resolution protein structure prediction in a hierarchical manner.

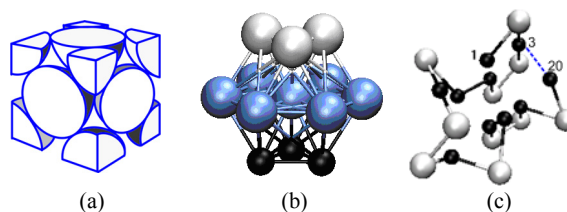


Fig. 2. (a) Face-centred-cube arrangement [47]. Each of the half spheres (or residues) at the centres of the six faces of the cubic chunk coincide with a full sphere at the centre of the cube. Each of the one-eighth spheres is placed at each of the eight corners of that cube. (b) Layers are separated by both colours and size, and dedicated connections are used to assist the visualisation of the concept. (c) A sample conformation, black and white spheres are indicating hydrophobic and hydrophilic residues respectively. There are total of 20 topological neighbours (TNs), a dotted line in between residue 3 and 20 indicates a sample TN.

We will continue using the 2D square HP model to explain the model relevant concepts in subsequent sections. Discussions related to 3D FCC will be re-appear along with the simulation results in section 3 on experiments.

2.2. Complexity of the lattice model

Even if we use a simplified lattice model and even if the sequence length is short, we have an inordinate number of valid (i.e., SAW) conformations [11, 12, 48]. For instance in a 2D HP model, for a sequence of n amino acids, the number of valid conformations is proportional to μ^n , where the connective constant or the effective coordinate number μ , is lattice dependent [12]. Prediction of the optimal conformation using the lattice model is also an *NP-complete* problem [9, 10]. To predict the backbone conformation of the folded protein from its amino acid sequence based on global interactions such as *hydrophobicity*, lattice models are used for approximation [31, 32, 35-37]. For *ab initio* prediction in *Critical Assessment of Structure Prediction* (CASP) [35-37], most

successful approaches followed the hierarchical paradigm where the lattice-based, backbone conformational sampling works very effectively at the top of the hierarchy. With further advancement toward all-atom or full modeling from the lattice, the energy functions include atom-based potentials from molecular mechanics packages such as CHARMM, AMBER, ECEPP and so on [49, 50]. Conformational search algorithms built on lattice models, which play a key role in solving PSP, are discussed next.

2.3. Rationale of low resolution model

Due to the involvement of immense computational complexity, a high-resolution model can only be applied only when a homologue of the target protein exists. Even though, it is applied, it lacks the ability to answer how and why a protein adopts its specific structure [3, 4] as again due to computational complexity it hardly covers a reasonable area of the search landscape. Thus, *ab initio* modeling, especially using a low resolution model is essential for a complete solution to the PSP problem including the investigation of the physicochemical principle of protein folding, as the recent data indicate that the fundamental physics underlying the folding process may be simpler than was previously thought [21]. Many changes in amino-acid sequence usually do not vary the overall topology of a protein which suggest that folding mechanisms depend more on the low resolution geometrical properties of the native state [21] and therefore the simplified model can be applied to understand the physical principles governing the folding processes [51].

Low resolution models can potentially be used to predict initial approximation of the protein structure, folding rate that depends on the height of the free-energy barrier, and the effects of mutations on the folding rate that depend on the region(s) of the protein ordered near the top of the barrier and so on, which in turn allows gaining insight into how adaptation and selection operate among large collections of sequences versus structures mapping. Also, an initial coarse sampling of the energy landscape makes the conformational search feasible as well as faster [4]. Thus low resolution models are found to be promising [21, 52-54].

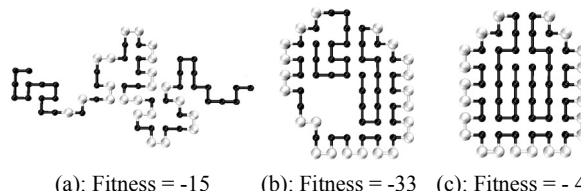
2.4. Nondeterministic conformational search algorithms

For solving *ab initio* PSP using the lattice model numerous nondeterministic approaches have been investigated: Monte Carlo (MC) simulation, Evolutionary MC (EMC) [14, 15], Simulated Annealing (SA), Tabu Search with Genetic Algorithm (GTB) [16], Ant Colony Optimisation [17], and Immune Algorithm (IA) based on Artificial Immune System (AIS) [55]. Due to their simplicity and search effectiveness, Genetic Algorithms (GAs) [13, 28, 34, 56-59] are the most attractive. They also provided superior performance over MC [57, 58]. The

concepts of GAs are also widely adapted within these algorithms. For instance, a new MC algorithm [14] adopted the population-based cut-and-paste (i.e., crossover) operation to achieve higher fitness. The evolutionary Monte Carlo (EMC) [15] algorithm incorporated the evolutionary features of genetic algorithms, such as a population which is updated by crossover and mutation operations. Jiang *et al.* applied the GA with a Tabu (GTB) search to solve PSP using lattice models [16]. Also, the conformational space annealing (CSA) [18, 60] algorithm is based on a concept similar to GA, where a “bank” in CSA is equivalent to the “population” in GA.

2.5. Focus of the paper

Given the widespread adaptation of GAs for PSP, the essence of GA, i.e. its crossover operation, can be made more effective by combining it with DFS which can have a positive impact on solving the PSP problem. In solving PSP using a conventional GA, where the optimum conformation is mostly physically compact (see Figure 3), a crossover-based converging conformation suffers from limited pathways and the algorithm, thus increasingly generates invalid conformations. Our hypothesis is that the combination of DFS with crossover can instead reveal potential pathways in solving PSP. Thus by using DFS, a repeatedly failing crossover having congested but potential sub-conformation can be allowed for a limited number of pathways as a possible candidate for crossover counterparts obtained, if there exists at least one path.



(a): Fitness = -15 (b): Fitness = -33 (c): Fitness = -42
 Fig. 3. As the search proceeds the conformation gets more compact: For a typical run, conformations at generation 1, 1434 and 5646 have been shown in (a), (b) and (c) respectively, showing the fitter conformation is relatively more compact.

2.6. Defining the GA operators for PSP problem

Here, we define the GA operators for the PSP problem based on the HP lattice model:

Crossover operation: For PSP, this aids the construction of global solutions by the cooperative combination of many local substructures [13]. We particularly followed the commonly-used crossover operation pioneered by Unger *et al.* [57], as illustrated in Figure 4, a single-point crossover. We follow this single-point crossover, since otherwise the converging conformation, being compact in nature, would generate more collisions or invalid conformations [15]. In

addition, the ability to rotate before joining within the crossover, in addition, provides a mutation-equivalent operation especially when *relative encoding* is followed [42]. For example, if we emulate the crossover in Figure 4 without the rotation, we can write using relative encoding that:

Crossover (a: ‘LFLRRRLRLLFLRFRLFL’, b: ‘RFFFRFRFLFLRLLFL’) \rightarrow would output, c: ‘LFLRRRLRLLFL*RLLFL’ without the rotation before joining. (Here, ‘*’ indicates an undefined move in relative encoding and here it indicates a non-SAW move.) But, with rotation, the conformation can have SAW, i.e. c: ‘LFLRRRLRLLFLRLLFL’.

Comparing c: ‘LFLRRRLRLLFL*RLLFL’ and c: ‘LFLRRRLRLLFLRLLFL’, it becomes clear that the ‘*’ is replaced by an ‘R’ after the rotation, which is genotypically a single-point mutation.

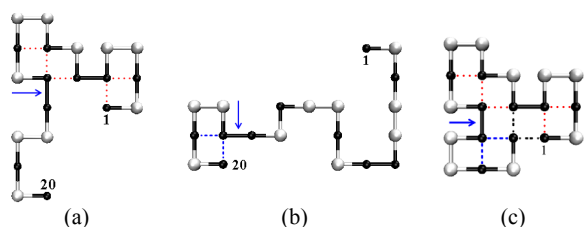


Fig 4. An example of the crossover operation [57]. Conformations are randomly cut and pasted with the cut point chosen randomly between residues 14 and 15. The first 14 residues of (a) are rotated first as needed (as allowed by the degree of freedom by the model configuration) and then joined with the last 6 residues of (b) to form (c), where fitness, $F = -9$. ‘ \rightarrow ’ indicates crossover positions.

Crossover failure: This implies that before joining two parts, all possible rotated positions at the joining point have been tried but failed to produce at least one valid conformation (i.e., a SAW).

Combination of crossover and DFS: For generating a conformation this implies that a DFS-generated random and partial path has been joined with the first half of the sub-conformation.

DFS after crossover failed: This implies that ‘combination of crossover and DFS’ has been performed after an occurrence of ‘crossover failure’.

Mutation operation: This involves pivot rotation (Figure 5) as basically pioneered by Unger *et al.* [57]. We employed single-point mutation to avoid more collisions.

Ordinary random conformation generation: This implies the generation of a SAW conformation based on random-move-only (RMO). In a 2D square lattice model *Left*, *Right* and *Forward* moves are permissible but *Backward* moves are prohibited. For a conformation, once a path search has failed after looking in the three possible degrees of the freedom the whole process restarts.

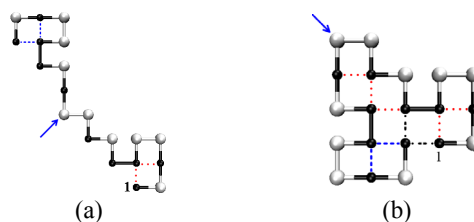


Fig 5. An example of the mutation operation [57]. Dotted lines indicate TN. Residue number 11 is chosen randomly as the pivot. For the move to apply, a 180° rotation (among a number of possible degree of freedom defined by the model configuration) alters (a) with $F = -4$ to (b) $F = -9$. ‘ \rightarrow ’ indicates the mutation residue.

Random conformation generation by DFS: This implies that we apply DFS to generate a SAW conformation. As the DFS proceeds, it stores the possible pathways using a stack-memory [19] and, upon total failure after trying all possible degrees of freedom on a particular location (i.e. lattice point), it can backtrack to restart from the stored options instead of restarting the creation of the whole conformation.

3. Experiments and results

We carried out experiments to empirically verify our hypothesis that combining DFS with crossover will be advantageous. The simple GA (SGA) applied for PSP is illustrated in Figure 6 and the crossover variations with the possible implementation have been shown in Figure 7.

Fig 6. Genetic Algorithm for solving PSP problem.[†]

1. Initialize fixed size current population (Pop_z) of randomly generated conformations.
2. Obtain new solution (S_{new}) from the current population by using **Crossover** and **Mutation** operations at the pre-specified rates (p_c and p_m respectively).
3. Assess quality or fitness, F , of S_{new} .
4. Promote the obtained S_{new} , and elite and untouched chromosomes, to the next generation and assign the new generation as the current population.
5. IF END-OF-SOLUTION is not reached THEN repeat from Step 2.

As shown in Figure 7, we have experimented with four variations of the crossover operation. *Crossover* (a) (see Figure 7(a)) represents a conventional crossover operation for PSP without DFS. *Crossover*(b) (see Figure 7(b)) applies DFS-based partial path generation with the sub-conformation immediately the sub-conformation fails to

[†]Terms in **bold** and *italic* are explained in section 2.6.

join with its counterpart sub-conformation after trying all possible degrees of freedom. *Crossover(d)* (see Figure 7(d)) is similar to *Crossover(b)* in operation but allows more time for a failed crossover to search for a suitable counterpart sub-conformation to match. *Crossover(c)* is a the most dissimilar variation of *Crossover(d)* where, instead of a sub-conformation looking for its counterpart sub-conformation in the population, *Crossover(c)* directly uses DFS to generate the rest of the path to complete the conformation. This alternative was investigated to determine an effective rate of DFS.

The default GA parameters for all experiments were set as population size (Pop_z) to 200, crossover rate (p_c) to 0.85 or 85%, mutation rate (p_m) to 5% and for elitism the elite rate was set to 5% [61, 62].

Fig 7. Crossover operation and variation details.[†]

(a)
<ol style="list-style-type: none"> 1. DO single-point <i>Crossover</i>. 2. IF '<i>Crossover failure</i>' = TRUE then 3. REPLACE one of the parents. 4. DO single-point <i>Crossover</i>. 5. END IF
(b)
<ol style="list-style-type: none"> 1. DO single-point <i>Crossover</i>. 2. IF '<i>Crossover failure</i>' = TRUE then 3. DO '<i>DFS after crossover failed</i>'. 4. END IF
(c)
<ol style="list-style-type: none"> 1. DO single-point '<i>Combination of crossover and DFS</i>'.
(d)
<ol style="list-style-type: none"> 1. DO apply option: (a). 2. IF no improvement for 5 consecutive generations, 3. DO apply option: (b). 4. END IF

The fold for longer PSP sequences generally has complex energy landscapes [32, 63-68], and hence these sequences would normally require longer time to converge. So, we chose these longer benchmark sequences (see Table 1) to highlight the true impact of this approach.

3.1. Experiments using 2D HP square model

A maximum of 2000 generations was allocated for each of the 10 iterations carried out per sequence and in each category of experiments. The benchmark PSP sequences use are shown in Table 1 for the 2D square HP lattice

[†]Terms in **bold** and *italic* are explained in section 2.6.

model [7], with the length ranging from 50 to 100 [69, 70]. The results are given in Table 2.

In Table 2, we include two other algorithms in their generic form, namely Unger's GA (UGA [57]) and Conformational Space Annealing (CSA) algorithm [18, 60] along with our proposed algorithm. It may be noted that UGA has been reported to have outperformed many MC variations [13, 57]. We emulated UGA in our experiment without changing the original parameter for cooling. The initial cooling temperature was set to 2 and was decreased by 0.99 every 200000 steps until the temperature reached 0.15.

Table 1

Benchmark protein sequences for 2D HP model

Len.	Sequences	Ref.
50	H2(PH)3PH4PH(P3H)2P4H(P3H)2PH4P(HP)3H2	[70]
60	P2H3PH8P3H10PHP3H12P4H6PH2PHP	[70]
64	H12(PH)2(P2H2)2P2HP2H2PPHP2H2P2(H2P2)2(HP)2 H12	[70]
85	4H4P12H6P12H3P12H3P12H3P1H2P2H2P2H2P1H1P 1H	[69]
100	3P2H2P4H2P3H1P2H1P2H1P4H8P6H2P6H9P1H1P2H 1P1H2P3H1P2H1P1H2P1H1P3H6P3H	[69]

'H' and 'P' in the sequence indicate hydrophobic and hydrophilic amino acids, respectively. *Len.* indicates length.

Table 2

Run results of 10 iterations on each PSP sequence (see Table 1 for the sequences). GA runs with four different crossover options (shown in Figure 7), have been compared

Len.	<i>X(a)</i>	<i>X(b)</i>	<i>X(c)</i>	<i>X(d)</i>	CSA	UGA
50	-17.3/ -20	-17.6 /- 20	-14.5/ -17	-18/-20	-17 / -19	-16.6 / -18
60	-29.2/ -32	-29.8/ -32	-27.8/ -31	-30.5/ -32	-30.4/- 32	-29/ -31
64	-29.1/ -31	-29.3/ -31	-25.2/ -29	-32/ -35	-29/-30	-27.8/ -31
85	-39.4/ -44	-39.6/ -45	-34.5/ -38	-43.4/ -46	-43.2/-	-41.4/ -46
100	-37.1/ -39	-37.6/ -41	-30.2/ -37	-38.5/ -42	-37.2/- 38	-37.4/ -40

The format of column entries is 'Average / Minimum' and indicate fitness function values. The *X* implies *Crossover* operation. Thus, *X(a)* indicates *Crossover(a)* as described above, and so on. CSA and UGA indicate Conformational Space Annealing Algorithm [18] and Unger's GA [57], respectively. **Bold** entries indicate the row-wise best values obtained.

We abstracted the general form of the CSA algorithm by removing the heuristic-based special moves, keeping the generic form intact, to provide a fair comparison in our experiment. Comparison with CSA algorithm is particularly important for our work, since the CSA approach has recently been applied in the PSP software ROSETTA [35, 71-74]. Both UGA and CSA ran 2000 GA generation equivalent runs per iteration.

3.2. Experiments using 3D HP FCC model

As mentioned earlier in section 3, we kept the same value for GA parameters, i.e., population size 200, crossover rate 85%, mutation rate 5% and elite rate 5%. We have also used the same set of benchmark HP sequences shown in Table 1 as used for the 2D square HP lattice model. For each sequences, the simulation ran for 10 iterations, however each run has been executed maximum of 1500 GA generations. The results obtained are shown in Table 3.

Table 3

Average and minimum-fitness-value from run-results of 10 iterations on each PSP sequence using 3D FCC model.

Len.	$X(a)$	$X(b)$	$X(c)$	CSA	$X(d)$	UGA
50	-70.3/ -77	-69.8/ -76	-59.6/ -69	-72.4/ -77	-78.4/ -82	-70.2/ -74
60	-145.8/ -157	-149.8/ -161	-115.1/ -121	-140.8/ -149	-152.7/ -161	-148.1/ -158
64	-138.2/ -147	-139.6/ -145	-119.2/ -130	-129/ -135	-140.8/ -147	-133.3/ -141
85	-191.6/ -204	-193.1/ -203	-169.8/ -187	-188.5/ -200	-197.9/ -212	-189.6/ -202
100	-180.3/ -190	-185.6/ -193	-142.7/ -154	-182.8/ -193	-191.1/ -204	-171.2/ -183

The format of column entries is ‘Average / Minimum’. The X implies *Crossover* operation. Thus, $X(i)$ indicates *Crossover*(i) where $i = 'a'$ to $'d'$. CSA and UGA indicate Conformational Space Annealing Algorithm and Unger’s GA, respectively. **Bold** entries indicate the row-wise best values obtained.

The performance using the 3D FCC model remains consistent when compared with the previous experiments performed using the 2D HP model in section 3.1. The $X(d)$ algorithm consistently performed the best and $X(c)$ performed the worst and slowest amongst all other search algorithms. The minimum conformations achieved in this experiment using 3D FCC model have been shown in Figure 8.

4. Discussion on the experimental results

We have introduced the concept of finding potential partial pathways using a depth-first search (DFS) strategy when a converging potential sub-conformation in a crossover failed to find a matching counterpart to produce a valid (i.e., having a self-avoiding-walk) conformation. Crossover variation $X(c)$ gave the worst results (see Table 2 and Table 3). $X(c)$ involves applying DFS constantly at the same rate as the crossover operation to generate the other half of the crossover portion, which is misleading the optimum results more than guiding them compared to the other strategies applied. $X(a)$ represents the crossover-only approach, that is, crossover with DFS, and $X(b)$ is the

variant where DFS is applied whenever a crossover fails. $X(b)$ demonstrates a slight improvement over $X(a)$. $X(d)$ performed the best, with results comparable to the UGA and CSA algorithms. This is most likely due to the fact that in $X(d)$, crossover was applied exhaustively by allowing a failed crossover to search for more counterparts to match and when there was no improvement at all in the whole population for consecutive few generations, the failed crossover is combined with DFS to generate possible potential pathways. It is interesting to note that, in our experiment we find DFS has *zero* failure in finding pathways. Thus, a constantly failing sub-conformation in a crossover operation, which is likely to have few possible pathways, can be salvaged using DFS to unravel the hidden paths effectively. As an alternative to DFS, breadth-first search (BFS) [19] could have been used; however, BFS is both memory and time intensive.

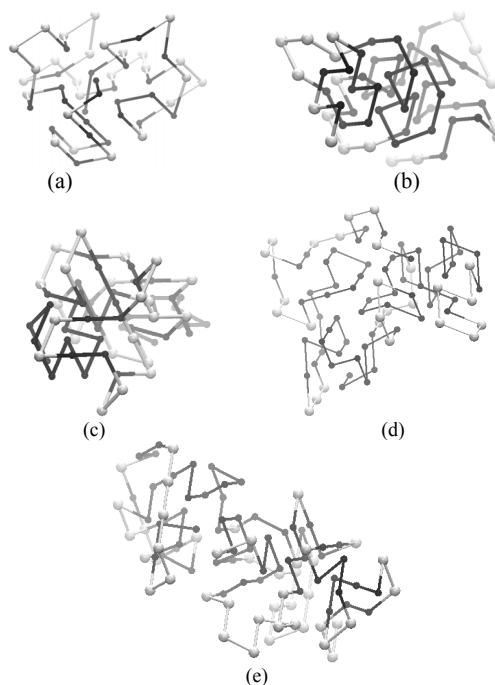


Fig 8. Minimum conformations found in the $X(d)$ simulation using the 3D FCC configuration. The figure from (a) to (e) correspond to the sequence lengthen from 50 to 100 (see Table 1) in the same order. The figures are drawn using VMD tool with ‘orthographic’ and ‘depth cueing’ options chosen.

The idea of the conformational space annealing (CSA) algorithm seemed appealing as it maintained the diversity based on maintaining distinguished conformation in a separate bank of population implying the division of the search space into a manageable finite number of banks (say, 50) which can represent possible distinguished (non-overlapped) region of the fitness landscape and for each region a best representative is maintained. However, with the immense complexity and vastness of the fitness

landscape associated with the PSP problems, dividing the landscape into finite region (say 50) does not make any difference: consider a very large number divided by 50, then numbers of possible conformations within each division almost remain very close to that very large numbers. Therefore, the real effort would be to find the best representative for each region with the large and convoluted search landscape which is again the main goal of the other entire search algorithms. Thus, for the convoluted nature of the PSP problem, it is not convincing that CSA can reasonably maintain the region representative sufficiently, and so practically it has not performed well.

5. Supplementary applications of DFS in PSP

It is important to remember that *ordinary random conformation generation*[†] takes exponential time, fitted to a curve given by the following Equation (1),

$$y = 2.8723 e^{0.0326x} \quad (1)$$

The square of coefficient of determination of Equation (1), $R^2 = 0.9832$, with increasing sequence length using the random-move-only (RMO) approach. In contrast, the runtime for *random conformation generation by DFS* remains a quadratic fitted curve (see Figure 9), as shown in Equation (2),

$$y = 0.02 x^2 - 0.5717x + 54.789 \quad (2)$$

with $R^2 = 0.9996$ (for Equation (2)). Although Equation (2) is mathematically quadratic, the coefficient of x^2 in Equation (2) being close to zero, the actual relationship can be considered to be almost linear.

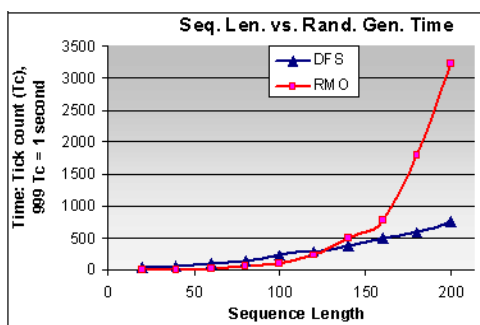


Fig 9. Random conformation generation: DFS approach versus random-move-only (RMO) approach. An average of 100 iterations is taken for a particular length of a single random conformation generation (not from benchmark sequences).

The application of *random conformation generation by DFS* may have a generally lower impact because totally

[†]Terms in **bold** and *italic* are explained in section 2.6.

random conformations are only generated for initialization of the population. However, to maintain diversity many GA approaches replenish the population for a considerable amount and at frequent intervals [75, 76]. For example, Hoque *et al.* have shown removal of chromosomes having 80-90% or greater similarity from a GA population helps it to perform better [75, 77]. After removal it is necessary to replenish the population by random conformations of 16 to 40% in each generation as indicated in Table 4. Thus, in such a case, for longer sequences, *random conformation generation by DFS* would make the GA search far more efficient.

Table 4

For various PSP sequences, entries indicate percentage of average chromosome removal per generation while percentage in column header (except #1) indicates equal or above percentage of similarity between two chromosomes while one of them are removed from the population to maintain diversity.

Len.	100%	90%	80%	70%	60%	50%
50	17.2	33.6	39.5	45.4	53.2	76.9
60	11.8	25.7	31.1	38.3	45.0	67.9
64	13.4	23.0	28.5	36.7	43.7	64.5
85	6.2	16.7	24.5	32.0	40.4	50.9
100	6.2	19.7	25.4	33.6	41.2	49.8

Len. indicates length.

6. Conclusions

A depth-first search (DFS) strategy at a low rate has been applied in combination with a powerful crossover operation. Together they revealed convoluted and microscopic pathways for solving protein structure prediction problem. Experiments using a variety of longer, standard benchmark sequences from the literature have demonstrated the efficacy and improved performance characteristics of this approach empirically using two different HP model configurations. The search strategy developed was inspired by the pathway hypothesis. Further work will be directed at exploring the biological significance and relevance of this novel approach.

Acknowledgments

Support from Australian Research Council (grant no DP0557303) is thankfully acknowledged.

References

- [1] J. Pietzsch, The importance of protein folding, Nature, <http://www.nature.com/horizon/proteinfolding/background/importance.html>, accessed: Jan (2009).

- [2] S. Petit-Zeman, Treating protein folding diseases, Nature, <http://www.nature.com/horizon/proteinfolding/background/treating.html>, accessed: Jan (2009).
- [3] J. Lee, S. Wu, Y. Zhang, *Ab Initio* Protein Structure Prediction. In: Rigden, D.J. (ed.): From Protein Structure to Function with Bioinformatics. Springer Netherlands (2009) 3–25.
- [4] D. Chivian, T. Robertson, R. Bonneau, D. Baker, AB INITIO METHODS. In: Bourne, P.E., Weissig, H. (eds.): Structural Bioinformatics. Wiley-Liss, Inc. (2003) 547–557.
- [5] F. Allen, *et al.*, Blue Gene: A vision for protein science using a petaflop supercomputer. IBM System Journal 40 (2001) 310-327.
- [6] M.T. Hoque, M. Chetty, A. Lewis, A. Sattar, DFS Based Partial Pathways in GA for Protein Structure Prediction Pattern Recognition in Bioinformatics (PRIB). LNCS, Springer, Melbourne, Australia (2008) 41-53.
- [7] K.A. Dill, Theory for the Folding and Stability of Globular Proteins. Biochemistry 24 (1985) 1501-1509.
- [8] R. Backofen, S. Will, A Constraint-Based Approach to Fast and Exact Structure Prediction in Three-Dimensional Protein Models. Constraints Journal 11 (2006) 5-30.
- [9] P. Crescenzi, D. Goldman, C. Papadimitriou, A. Piccolboni, M. Yannakakis, On the complexity of protein folding (extended abstract). the second annual international conference on Computational molecular biology. ACM (1998) 597-603.
- [10] B. Berger, T. Leighton, Protein Folding in the Hydrophobic-Hydrophilic (HP) Model is NP-Complete. Journal of Computational Biology 5 (1998) 27-40.
- [11] R. Schiemann, M. Bachmann, W. Janke, Exact Enumeration of Three – Dimensional Lattice Proteins. Computer Physics Communications, Elsevier Science. 166 (2005) 8-16.
- [12] A. J. Guttmann, Self-avoiding walks in constrained and random geometries. In: B.K. Chakrabarti (Ed.), Statistics of Linear Polymers in Disordered Media, Elsevier (2005) 59-101.
- [13] R. Unger, J. Moulton, On the Applicability of Genetic Algorithms to Protein Folding. The Twenty-Sixth Hawaii International Conference on System Sciences, Vol. 1 (1993) 715-725.
- [14] U. Bastolla, H. Frauenkron, E. Gerstner, P. Grassberger, W. Nadler, Testing a new Monte Carlo Algorithm for Protein Folding. National Center for Biotechnology Information 32 (1998) 52-66.
- [15] F. Liang, W.H. Wong, Evolutionary Monte Carlo for protein folding simulations. J. Chem. Phys 115 (2001) 3374-3380.
- [16] T. Jiang, Q. Cui, G. Shi, S. Ma, Protein folding simulations of the hydrophobic-hydrophilic model by computing tabu search with genetic algorithms. J. Chem. Phys. 119 (2003) 4592-4596.
- [17] A. Shmygelska, H.H. Hoos, An ant colony optimization algorithm for the 2D and 3D hydrophobic polar protein folding problem. BMC Bioinformatics 6:30 (2005) 1-22.
- [18] J. Lee, Conformational space annealing and a lattice model Protein. Journal of the Korean Physical Society 45 (2004) 1450-1454.
- [19] T.H. Cormen, C.E. Leiserson, R.L. Rivest, Introduction to Algorithms. MIT Press (1998).
- [20] L. Toma, S. Toma, Folding simulation of protein models on the structure based cubo-octahedral lattice with the Contact interactions algorithm. Protein Science 8 (1999) 196-202.
- [21] D. Baker, A surprising simplicity to protein folding. Nature 405 (2000) 39-42.
- [22] V.S. Pande, D. Rokhsar, Folding pathway of a lattice model for proteins. Proc. Natl. Acad. Sci. (PNAS) USA, Biochemistry 96 (1999) 273–1278.
- [23] C. Levinthal, Are there pathways for protein folding? Journal of Chemical Physics 64 (1968) 44-45.
- [24] M. T. Hoque, M. Chetty, L.S. Dooley, A Guided Genetic Algorithm for Protein Folding Prediction Using 3D Hydrophobic-Hydrophilic Model. Special session in WCCI / IEEE Congress on Evolutionary Computation (CEC) (2006) 2339-2346.
- [25] M. T. Hoque, M. Chetty, L.S. Dooley, Significance of Hybrid Evolutionary Computation for *Ab Initio* Protein Folding Prediction In: Grosan, C., Abraham, A., Ishibuchi, H. (eds.): Hybrid Evolutionary Algorithms, Vol. 75. Springer-Verlag, Berlin (2007) 241-268.
- [26] M. T. Hoque, M. Chetty, L.S. Dooley, A Hybrid Genetic Algorithm for 2D FCC Hydrophobic-Hydrophilic Lattice Model to Predict Protein Folding. 19th ACS Australian Joint Conference on Artificial Intelligence. LNAI, Springer (2006) 867-876.
- [27] M. T. Hoque, M. Chetty, L.S. Dooley, Fast computation of the fitness function for protein folding prediction in a 2D hydrophilic-hydrophobic model. Journal published in the special issue of the International Journal of Simulation Systems, Science and Technology 6 (2005) 27-37.
- [28] M. T. Hoque, M. Chetty, L.S. Dooley, A New Guided Genetic Algorithm for 2D Hydrophobic-Hydrophilic Model to Predict Protein Folding. IEEE Congress on Evolutionary Computation (CEC), Edinburgh, UK (2005) 259-266.
- [29] M.T. Hoque, M. Chetty, A.Sattar, Protein Folding Prediction in 3D FCC HP Lattice Model Using Genetic Algorithm Bioinformatics special session, IEEE Congress on Evolutionary Computation (CEC), Singapore (2007) 4138-4145.
- [30] K. Ghosh, S.B. Ozkan, K.A. Dill, The Ultimate Speed Limit to Protein Folding Is Conformational Searching. Journal of American Chemical Society 129 (2007) 11920-11927.
- [31] K. F. Lau, K.A. Dill, A lattice statistical mechanics model of the conformational and sequence spaces of proteins Macromolecules 22 (1989) 3986-3997.
- [32] K.A. Dill, S. Bromberg, K. Yue, K. M. Fiebig, D. P. Yee, P.D. Thomas, H.S. Chan, Principles of protein folding – A perspective from simple exact models. Protein Science 4 (1995) 561-602.
- [33] K.A Dill, S.B. Ozkan, T.R. Weikel, J.D Chodera, V. A Voelz, The protein folding problem: when will it be solved? Current Opinion in Structural Biology 17 (2007) 342-246.
- [34] D. W. Corne, G.B Fogel, An Introduction to Bioinformatics for Computer Scientists. In: Fogel, G.B., Corne, D.W. (eds.): Evolutionary Computation in Bioinformatics (2004) 3-18.
- [35] D. Baker, Prediction and design of macromolecular structures and interactions. Phil. Trans. R. Soc. B 361 (2006) 459-463.
- [36] O. Schueler-Furman, C. Wang, P. Bradley, K. Misura, D. Baker, Progress in Modeling of Protein Structures and Interactions. Science 310 (2005) 638-642.
- [37] Y. Xia, E. S. Huang, M. Levitt, R. Samudrala, Ab Initio Construction of Protein Tertiary Structures using a Hierarchical Approach. J. Mol. Biol. 300 (2000) 171-185.
- [38] R. Backofen, S. Will, P. Clote, Algorithmic approach to quantifying the hydrophobic force contribution in protein folding. Pacific Symp. On Biocomputing 5 (2000) 92-103.
- [39] K. Yue, K. A. Dill, Sequence-structure relationships in proteins and copolymers Phys. Rev. E 48 (1993) 2267 – 2278.
- [40] L. Toma, S. Toma, Contact interactions methods: A new Algorithm for Protein Folding Simulations. Protein Science 5 (1996) 147-153.
- [41] E. Bornberg-Bauer, Chain Growth Algorithms for HP-Type Lattice Proteins. RECOMB, Santa Fe, NM, USA (1997) 47-55.
- [42] M.T. Hoque, M. Chetty, L.S Dooley, Non-Isomorphic Coding in Lattice Model and its Impact for Protein Folding Prediction Using Genetic Algorithm. IEEE Computational Intelligence in Bioinformatics and Computational Biology (CIBCB). IEEE, Toronto, Canada (2006) 1-8.
- [43] R. Samudrala, Y. Xia, M. Levitt, A Combined Approach for ab initio Construction of Low Resolution Protein Tertiary Structures from Sequence Pacific Symposium on Biocomputing (PSB) 4 (1999) 505-516.
- [44] T. Hoque, M. Chetty, A. Sattar, Extended HP Model for Protein Structure Prediction. Journal of Computational Biology 16 (2009) 85-103.
- [45] T. C. Hales, A proof of the Kepler conjecture. Annals of Mathematics 162 (2005) 1065-1185.

- [46] R. Backofen, S. Will, E. Bornberg-Bauer, Application of constraint programming techniques for structure prediction of lattice proteins with extended alphabets *Bioinformatics* 15 (1999) 234-242.
- [47] O.S. University, The Structure of Crystalline Solids Online Lectures: Materials Science and Engineering, Ohio State University, http://www.matsceng.ohio-state.edu/mse205/lectures/chapter3/index_chap3.htm, accessed: Dec (2008).
- [48] M. Chen, K.Y. Lin, Universal amplitude ratios for three-dimensional self-avoiding walks. *Journal of Physics A: Mathematical and General* 35 (2002) 1501-1508.
- [49] I.K. Roterman, M.H. Lambert, K.D. Gibson, H. Scheraga, A comparison of the CHARMM, AMBER and ECEPP potentials for peptides. II. Phi-psi maps for N-acetyl alanine N'-methyl amide: comparisons, contrasts and simple experimental tests. *J. Biomol. Struct. Dynamics* 7 (1989) 421-453.
- [50] W.D. Cornell, P. Cieplak, C.I. Bayly, I.R. Gould, Jr, K.M.M., D.M. Ferguson, D.C. Spellmeyer, T. Fox, J.W. Caldwell, P.A. Kollman, A second generation force field for the simulation of proteins and nucleic acids. *J. Am. Chem. Soc.* 117 (1995) 5179-5197.
- [51] Y. Duan, P. A. Kollman, Computational protein folding: From lattice to all-atom. *IBM Systems Journal* 40 (2001) 297-309.
- [52] T. Head-Gordon, S. Brown, Minimalist models for protein folding and design. *Current Opinion in Structural Biology* 13 (2003) 160-167.
- [53] R. Wroe, E. Bornberg-Bauer, H.S. Chan, Comparing Folding Codes in Simple Heteropolymer Models of Protein Evolutionary Landscape: Robustness if the Superfunnel Paradigm. *Biophysical Journal* 88 (2005) 118-131.
- [54] R. Santana, P. Larrañaga, J.A. Lozano, Protein Folding in Simplified Models With Estimation of Distribution Algorithms. *IEEE Transactions on Evolutionary Computation* 12 (2008) 418-438.
- [55] V. Cutello, G. Nicosia, M. Pavone, J. Timmis, An Immune Algorithm for Protein Structure Prediction on Lattice Models. *IEEE Transactions on Evolutionary Computation* 11 (2007) 101-117.
- [56] O. Takahashi, H. Kita, S. Kobayashi, Protein Folding by A Hierarchical Genetic Algorithm. 4th Int. Symp. AROB (1999) 334-339.
- [57] R. Unger, J. Moult, Genetic Algorithms for Protein Folding Simulations. *Journal of Molecular Biology* 231 (1993) 75-81.
- [58] R. Unger, J. Moult, Genetic Algorithm for 3D Protein Folding Simulations. 5th International Conference on Genetic Algorithms (1993) 581-588.
- [59] R. König, T. Dandekar, Refined Genetic Algorithm Simulation to Model Proteins. *Journal of Molecular Modeling* 5 (1999) 317-324.
- [60] J. Lee, H. A. Scheraga, S. Rackovsky, New Optimization Method for Conformational energy Calculations on Polypeptides: Conformational Space Annealing. *Journal of Computational Chemistry* 18 (1997) 1222-1232.
- [61] R.L Haupt, S.E. Haupt, *Practical Genetic Algorithms* (2004).
- [62] J. G. Digalakis, K.G. Margaritis, An experimental Study of Benchmarking Functions for Genetic Algorithms Intern. *J. Computer Math.* 79 (2002) 403-416.
- [63] S.D. Flores, J. Smith, Study of Fitness Landscapes for the HP model of Protein Structure Prediction. *IEEE CEC* (2003) 2338-2345.
- [64] N. Mousseau, G. T. Barkema, Exploring High-Dimensional Energy Landscape. *Computing in Science & Engineering*, IEEE 1 (1999) 74-80, 82.
- [65] U.H.E. Hansmann, Protein Folding in Silico: An Overview. *IEEE Computing in Science & Engineering* 5 (2003) 64-69.
- [66] J. Skolnick, A. Kolinski, Computational Studies of Protein Folding *IEEE Computing in Science & Engineering* 3 (2001) 40-50.
- [67] Y. Cui, W.H. Wong, E. Bornberg-Bauer, H.S. Chan, Recombinatoric exploration of novel folded structures: A heteropolymer-based model of protein evolutionary landscapes. *PNAS* 99 (2002) 809-814.
- [68] K. Schreiner, Distributed Project Tackle Protein Mystery. *Computing in Science & Engineering*, IEEE 3 (2001) 13-16.
- [69] N. Lesh, M. Mitzenmacher, S. Whitesides, A Complete and Effective Move Set for Simplified Protein Folding. *RECOMB*, Berlin, Germany (2003) 188-195.
- [70] W.E. Hart, S. Istrail, HP Benchmarks, http://www.cs.sandia.gov/tech_reports/compbio/tortilla-hp-benchmarks.html, accessed: Aug (2005).
- [71] Yiliu, Rosetta 2.1.0., 2007-2008 The Rosetta Commons, <http://www.rosettacommons.org/tiki/tiki-index.php?page=Change+Log>, Accessed: March (2008).
- [72] R. Bonneau, J. Tsai, I. Ruczinski, D. Chivian, C. Rohl, C.E.M. Strauss, D. Baker, Rosetta in CASP4: Progress in Ab Initio Protein Structure Prediction. *PROTEINS: Structure, Function, and Genetics* 5 (2001) 119-126.
- [73] P. Bradley, D. Chivian, J. Meiler, K.M.S Misura, C.A. Rohl, W.R. Schief, W.J. Wedemeyer, O. Scueler-Furman, P. Murphy, J. Schonbrun, C.E.M Strauss, D. Baker, D. Rosetta Predictions in CASP5: Success, Failure, and Prospects for Complete Automation. *PROTEINS: Structure, Function, and Genetics* 53 (2003) 457-468.
- [74] K.T. Simons, R. Bonneau, I. Ruczinski, D. Baker, Ab Initio Protein Structure Prediction of CASP III Target Using ROSETTA. *PROTEINS: Structure, Function, and Genetics* 3 (1999) 171-176.
- [75] M.T. Hoque, M. Chetty, L.S. Dooley, Generalized Schemata Theorem Incorporating Twin Removal for Protein Structure Prediction. *Pattern Recognition in Bioinformatics*. Springer, Singapore (2007) 84-97.
- [76] V.K. Koumoussis, C.P. Katsaras, A Saw-Tooth Genetic Algorithm Combining the Effects of Variable Population Size and Reinitialization to Enhance Performance. *IEEE Transaction on Evolutionary Computation* 10 (2006) 19-28.
- [77] M.T. Hoque, M. Chetty, A. Lewis, A. Sattar, Twin Removal in Genetic Algorithms for Protein Structure Prediction Using Low Resolution Model. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* (2009) doi:10.1109/TCBB.2009.34.



Md Tamjidul Hoque received both his B. Sc. Engg. and M.Sc. Engg. degrees in Computer Science and Engineering (CSE) from Bangladesh University of Engineering and Technology in 1998 and 2002 respectively and received his PhD degree in IT from Monash University (Australia) in 2008. He was a lecturer at the CSE department, Ahsanullah

University of Science and Technology, 1998-99. He was in the technical management being IT incharge and DGM at Bashundhara Group, Dhaka, Bangladesh from December 1999-04. Currently he is a research fellow at Griffith University (Australia) in Discovery Biology, Eskitis and a member of IIIS. His research focus is on 'ab initio protein structure prediction' and 'high content image analysis and algorithm development'. His research interests include Evolutionary Computation, Bioinformatics, Networking, Communication, Database System, Compiler Design, Automata Theory, Distributed Systems and Parallel Computing, Computer Architecture, Petri Net Theory, Security and Operating System.

Dr Madhu Chetty has been with Monash University, Australia since 1995 and is currently the Deputy Head of Gippsland school of Information Technology. His research interests include bioinformatics, optimization, computational intelligence, and modeling complex systems. Dr Chetty has authored over 100 scientific articles which include book chapters and articles in journals and international conferences. He is Senior Member of IEEE and Fellow, Institution of Engineers (India). He is currently serving as Chair of technical committee (TC-20) of International Association for Pattern Recognition (IAPR) on



bioinformatics and was General Chair of PRIB'08 (Pattern Recognition in Bioinformatics) conference. He has also served as Vice Chair of the IEEE CIS Technical Committee on Bioinformatics and Bioengineering. He is serving as the Associate Editor of the Elsevier's Neurocomputing journal and is on the editorial board of three other journals in bioinformatics. Prior to his career at Monash, Dr. Chetty worked at VRCE (now VNIT), Nagpur India (1980-1993), and University of Melbourne (1993-1995).



Dr Andrew Lewis is a Senior Research Specialist in Research Computing Services and an Adjunct Senior Lecturer in ICT at Griffith University. Prior to this appointment he worked in industrial applied research with BHP Billiton. His research interests include: parallel optimisation algorithms for large numerical simulations, including gradient descent, direct search methods, evolutionary programming, particle swarm and ant colony systems, multi-objective optimisation techniques for engineering design, and parallel, distributed and grid computing methods. He has numerous publications in the area of optimisation algorithms and applications.



Professor Abdul Sattar is the founding Director of the Institute for Integrated and Intelligent Systems (IIS) and a Professor of Computer Science and Artificial Intelligence at Griffith University. He is also a Research Leader at National ICT Australia (NICTA) Queensland Research Lab (QRL), where he has held the positions of QRL Education Director (2006-08) and Leader of the Smart Applications for Emergencies (SAFE) project (2005-08), and is currently leading the QRL node of NICTA's largest project, Advanced Technologies for Optimisation and Modelling in Constraints (ATOMIC). He has been an academic staff member at Griffith University since February 1992 as a lecturer (1992-95), senior lecturer (1996-99), and professor (2000-present) within the School of Information and Communication Technology. Prior to his career

at Griffith University, he was a lecturer in Physics in Rajasthan, India (1980-82), and a research scholar at Jawaharlal Nehru University, India (1982-85), the University of Waterloo, Canada (1985-87), and the University of Alberta, Canada (1987-1991).

Associate Professor Vicky Avery obtained her PhD in 1994 from Flinders University of South Australia, and was awarded an Australian NHMRC Postdoctoral Fellowship which was undertaken at the University of Adelaide. Between 1998-2004, A/Prof Avery was based at Active



Biotech AB, Sweden, and she held several positions including Section Head for Protein Interaction and Drug Discovery; Scientific Project Leader to identify the molecular target of 'Laquinimod', a novel oral treatment for relapsing multiple sclerosis, that has successfully concluded Phase IIb trials; Director of Biochemistry and Molecular Biology and Director, Business Development.

Also of significance, she was responsible for the development of assays for FDA to assess efficacy of a cholera vaccine designed and developed assays to identify immuno-modulatory compounds against CD80, which led to RhuDex®, an oral treatment for RA in clinical trials. A/Prof Avery specializes in high-throughput and high content screening. As the Head of Discovery Biology for the AstraZeneca - Griffith University collaboration, she was responsible for more than 49 HTS campaigns conducted between 2004-2007. The Discovery Biology team has also successfully designed and implemented HTS assays for Malaria (MMV) and African Sleeping Sickness (DnDi), being awarded MMV Project of the Year (2007) for innovative use of technology to identify new anti-malarials. As a Programme Leader for the recently established CRC Cancer Therapeutics, she has played an active role in the acquisition of funds and establishment of the Bioactive Discovery (HTS) programme.