

An Enhanced Genetic Algorithm for *Ab initio* Protein Structure Prediction

Mahmood A. Rashid, Firas Khatib, Md Tamjidul Hoque, and Abdul Sattar

Abstract—*In-vitro* methods for protein structure determination are time-consuming, cost-intensive, and failure-prone. Because of these expenses, alternative computer-based predictive methods have emerged. Predicting a protein’s three-dimensional structure from only its amino acid sequence—also known as *ab initio* protein structure prediction—is computationally demanding because the search space is astronomically large and energy models are extremely complex. Some successes have been achieved in predictive methods but these are limited to small sized proteins (around 100 amino acids); thus, developing efficient algorithms, reducing the search space, and designing effective search guidance heuristics are necessary to study large sized proteins. An on-lattice model can be a better ground for rapidly developing and measuring the performance of a new algorithm, and hence we consider this model for larger proteins (>150 amino acids) to enhance the genetic algorithms framework. In this paper, we formulate protein structure prediction as a combinatorial optimization problem that uses three-dimensional face-centered-cubic lattice coordinates to reduce the search space and hydrophobic-polar energy model to guide the search. The whole optimization process is controlled by an enhanced genetic algorithm framework with four enhanced features: *i*) an exhaustive generation approach to diversify the search; *ii*) a novel hydrophobic core-directed macro-mutation operator to intensify the search; *iii*) a per-generation duplication elimination strategy to prevent early convergence; and *iv*) a random-walk technique to recover from stagnation. On a set of standard benchmark proteins, our algorithm significantly outperforms state-of-the-art algorithms. We also experimentally show that our algorithm is robust enough to produce very similar results regardless of different parameter settings.

Index Terms—Genetic Algorithms, Combinatorial Optimization, Macro-mutation, Random-walk, Protein Structure Prediction, HP Model, FCC Lattice.

I. INTRODUCTION

PROTEINS are amongst the most important macromolecules in all living organisms. More than half of the dry weight of a cell is protein [1], [2]. Proteins are the sequential-chains of amino acids connected together by single peptide bonds. These connected chains fold into their functional three-dimensional (3D) structures [3] and regulate the cellular activities in living organisms to keep them alive. The function of a protein greatly depends on its folded 3D

M. A. Rashid is with the School of CIMS, University of the South Pacific, Suva, Fiji; IIIS, Griffith University, Brisbane, QLD, Australia and Department of CIS, University of Massachusetts Dartmouth, MA, USA. Email: mahmood.rashid@usp.ac.fj.

F. Khatib is with the Department of CIS, University of Massachusetts Dartmouth, MA, USA. Email: fkhatib@umassd.edu.

M. T. Hoque is with the Department of CS, University of New Orleans, LA, USA. Email: thoque@uno.edu.

A. Sattar is with the Institute for Integrated and Intelligent Systems (IIIS), Griffith University, Brisbane, QLD, Australia. Email: a.sattar@griffith.edu.au

Copyright © 2012 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

structure (also known as *native* structure), which has the lowest possible free energy—the approximation of interaction energies amongst the amino acids in a protein [4]. There are, however, some exceptions such as proteins of *prion domain* that have multiple functional structures [5]. Many fatal diseases such as prion disease, Alzheimer’s disease, Huntington’s disease, Parkinson’s disease, diabetes, and cancer are associated with the aggregation of non-functional proteins due to misfolding [6], [7], [8], [9]. The 3D structures of proteins are decidedly important in rational drug design [10], [11], protein engineering [12], [13], and biotechnology [14], [15]; thus, the protein structure prediction (PSP) has emerged as an important multi-disciplinary research problem.

Among the *in-vitro* methods, x-ray crystallography, nuclear magnetic resonance (NMR) spectroscopy, and electron microscopy are widely used for determining protein structures. X-ray crystallography, which is considered to be one of the available relatively accurate methods, requires the target proteins to undergo a complex crystallization process. As a result, x-ray crystallography often remains less feasible for proteins which are hard to crystallize such as membrane proteins [16]. In addition, these experimental methods are time consuming, cost-intensive and failure-prone and are not sufficient to fill the gap between the number of known protein sequences ($\approx 50M^a$) and the number of solved structures ($\approx 0.1M^b$). Therefore, developing efficient algorithms for the predictive methods is thought to be a possible solution for bridging this gap. Again, predicting the 3D structure of a protein from only its amino acid sequence is computationally demanding [17] because the search space is astronomically large and the conformational-energy space is hugely complex. Some successes are achieved in high resolution predictive methods [18], [19], [20], [21], [22], but those are limited to mostly of protein-length <100 amino acids and some cases ≈ 150 amino acids [23], [24]. However, the average length of proteins could be more for some important domains to study, such as 270 ± 9 for archaeal, 330 ± 5 for bacteria, and 449 ± 25 for eukaryotes [25], [26]. This limitation of high resolution models eventually motivates us to consider low resolution-based models for PSP. Because of the simplicity, a lattice-based model helps develop and validate new algorithms fast particularly for larger proteins (>150 amino acids). The high performing algorithms on lattice models could then be tested by implementing on high-resolution models. In a work Higgs *et al.* [27] applied a similar strategy by replacing the Monte Carlo (MC) or Conformational Space Annealing (CSA) used in Rosetta with a genetic algorithm and found improved

^aThe EBI: <http://www.ebi.ac.uk/uniprot/TrEMBLstats/> (as of 22-Jul-2015)

^bThe PDB: <http://www.pdb.org/> (as of 22-Jul-2015)

results.

According to *Levinthal's paradox* [28], it is impossible for a protein to go through all of its possible conformations to reach the correct native state, since it would take an astronomically large time, while in reality proteins take only a few seconds or less to reach their native state. This implies that there should be a folding pathway. In addition, *Anfinsen's thermodynamic hypothesis* [29] states that, at least for small globular proteins, the native structure is determined only by the protein's amino acid sequence. This also states that, at normal conditions (temperature, solvent concentration and composition, etc.), when folding occurs, the native structure is a unique, stable, and kinetically accessible minimum of the free energy.

Levinthal's paradox and Anfinsen's hypothesis are the bases of formulating *ab initio* PSP as an optimization problem. Given a protein's amino acid sequence, the problem is to find the 3D structure of a protein such that the total interaction energy amongst the amino acids in the sequence is minimized. For this combinatorial optimization problem, we challenged ourselves to propose an enhanced genetic algorithm to find a combination of points on the 3D face-centered-cubic (FCC) lattice space to map the amino acids by minimizing the interaction energy.

The state-of-the-art results on the face-centered cubic (FCC) lattice based hydrophobic-polar (HP) energy model have been achieved by local search (LS) methods in [30], [31], [32], [33]. On the other hand, genetic algorithms (GA) [34] and tabu search [35] found promising results on 2D and 3D hexagonal lattice based HP models. In general, the success of either GA- or LS-based methods crucially depends on the balance of the diversification and the intensification of the sampling over the search space. Moreover, these algorithms often get stuck in local minima and perform poorly on large sized proteins. To progress further with these algorithms will require us to address the aforementioned issues appropriately, and an on-lattice low resolution model can provide such a feasible paradigm to design efficient sampling algorithms.

In this paper, we propose a population-based algorithm within the GA framework, named GAPlus, for backbone-only PSP with the hope of studying larger proteins. We use an HP energy model and a 3D FCC lattice to simplify the problem. In GAPlus, we introduce *i*) an exhaustive generation approach to diversify the search; *ii*) a novel hydrophobic core-directed macro-mutation operator to intensify the search; *iii*) a per-generation duplication elimination strategy to prevent early convergence; and *iv*) a random-walk technique to recover from stagnation. On a set of benchmark proteins, the GAPlus significantly outperforms state-of-the-art algorithms for PSP on the same models. In conjunction with the comparative results, we present a detailed empirical study on GA operators and GA parameters.

The rest of the paper is organized as follows: Section II reviews background knowledge; Section III discusses related work on PSP; Section IV describes our GAPlus framework; Section V presents the experimental results and analyses; Section VI presents a discussion on the enhanced features of GAPlus and their suitability of adaptation in the other classes of optimization problems; and finally, Section VII draws the

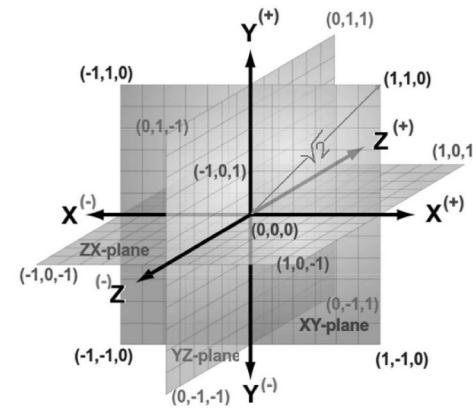


Figure 1. 3D FCC lattice space with 12 basis vectors on the Cartesian coordinates.

conclusions and outlines our future research.

II. BACKGROUND

The computer-based predictive methods are mostly available for template based modeling (TBM) and *ab initio* modeling. In template based modeling, if a solved structure is found in the protein data bank (PDB) that has a significant sequential identity ($\geq 35\%$) with the target [36], [37], it is taken as a template. Then, different algorithms [38], [39], [40], [41] are applied to refine the template to obtain the structure of the target protein.

However, *ab initio* modeling depends only on the amino acid sequence of the protein. The high resolution *ab initio* models [42], [18] consider all atoms of the backbone with their side chains, a complex energy model (such as molecular dynamics), and an internal or angular coordinates in 3D space. On the other hand, the low resolution models [30], [31], [32], [33] consider an amino acid with its side-chain as one residue of the backbone, a relatively simpler energy function (such as hydrophobic-polar), and on-lattice or off-lattice representation. Our work is based on low resolution *ab initio* modeling that uses the HP energy model for evaluating the structures and the FCC lattice points on 3D space for mapping the structures maintaining a self-avoiding-walk constraint within an enhanced GA framework.

A. Self-Avoiding Walk

In lattice-based protein representations, the amino acids of a given sequence are mapped on lattice points to build the protein backbone satisfying a self-avoiding-walk constraint. This constraint ensures no re-visitation of any lattice point during the sequence mapping.

B. 3D FCC Lattice

The FCC lattice has the highest packing density compared to the other existing lattices [43]. The hexagonal closed pack (HCP) lattice—known as cuboctahedron [34]—also has 12 neighbors corresponding to 12 *basis vectors* with real-numbered coordinates that causes the loss of structural precision during structure mapping. Therefore, we do not consider the HCP lattice in this work.

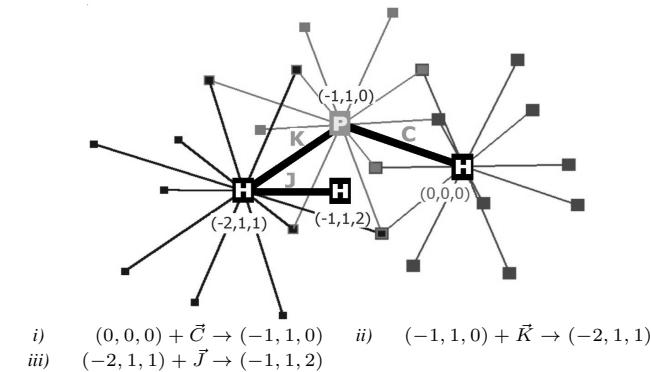


Figure 2. The degree of freedom in 3D FCC lattice space. The backbone mapping of a 4-residue (HPHH) sequence segment is shown by the thicker lines. \vec{C} , \vec{K} , and \vec{J} are the basis vectors as stated in Section II-B. The encoded string for the segment is CKJ. The encoding and decoding procedures are discussed in Section IV-B2

1) Topological Neighbors of a 3D FCC lattice point:

In FCC lattice space, each lattice point has 12 neighbors with corresponding 12 basis vectors as shown in Figure 1. The following 12 basis vectors, which are the topological neighbors of the origin $(0, 0, 0)$, are the building blocks (genes) of our encoded solutions (chromosomes or individuals):

$$\begin{array}{llll} A = (1, 1, 0) & B = (-1, -1, 0) & C = (-1, 1, 0) & D = (1, -1, 0) \\ E = (0, 1, 1) & F = (0, -1, -1) & G = (0, 1, -1) & H = (0, -1, 1) \\ I = (-1, 0, -1) & J = (1, 0, 1) & K = (-1, 0, 1) & L = (1, 0, -1) \end{array}$$

We map the conformations on the 3D FCC lattice points using a sequence of basis vectors (see Figure 2).

2) Degree of Freedom and Computational Complexity: The local degree of freedom (DOF) on 3D FCC lattice points varies based on the amino acid positions in the sequence. Starting from a random point (any of the twelve FCC lattice points), the DOF of the first and the second amino acid in the series is 12. However, the DOF for all other amino acids starting from 3 onward is 11 (see Figure 2). For sampling all possible conformations of a sequence with N amino acids, an algorithm needs to explore a search space that consists of $12^2 \times 11^{(N-2)} \approx 11^N$ valid combinations.

C. HP Energy Model

The 20 primary constituent amino acids of proteins are broadly divided into two categories, based on the hydrophobicity: (a) the hydrophobic amino acids are denoted as H (Gly, Ala, Pro, Val, Leu, Ile, Met, Phe, Tyr, Trp); and (b) the hydrophilic or polar amino acids are denoted as P (Ser, Thr, Cys, Asn, Gln, Lys, His, Arg, Asp, Glu). In the HP

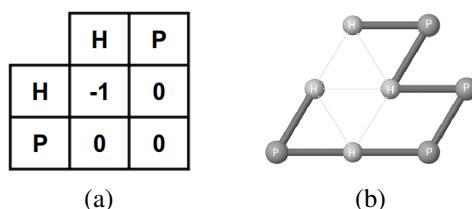


Figure 3. (a) HP energy model [44] and (b) H-H contacts in HP model for a random sequence HPHPPHPH on a 2D triangular lattice for convenience of depiction.

Algorithm 1: A Genetic Algorithm

```

1 initialize population and evaluate individuals;
2 while (!stopCondition) do
3   select the best-fit individuals for reproduction;
4   breed new individuals through crossover and
      mutation operations;
5   evaluate new generated individuals;
6   replace least-fit individuals with the new better
      individuals;

```

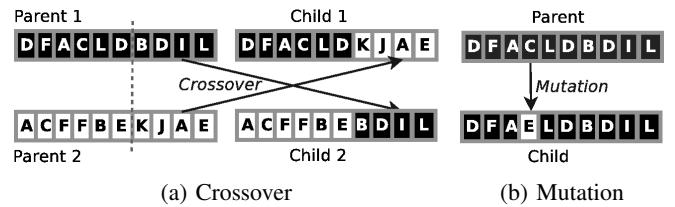


Figure 4. Typical (a) crossover and (b) mutation operators

model [44], when two non-consecutive hydrophobic amino acids become topological neighbors, they contribute a certain amount of negative energy, which, for simplicity, is shown as -1 in Figure 3a. The total energy (E) of a conformation based on the HP model is the sum of the contributions from all pairs of non-consecutive hydrophobic amino acids as shown in Equation (1). In Figure 3b, the number of such H-H contacts is 5. Therefore, the fitness of the structure of the sequence HPHPPHPH is -5 in the HP energy model.

$$E = \sum_{i < j-1} c_{ij} \times e_{ij} \quad (1)$$

Here, $c_{ij} = 1$ if amino acids i and j are non-consecutive neighbors on the lattice, otherwise 0; and $e_{ij} = -1$ if the i_{th} and j_{th} amino acids are hydrophobic, otherwise 0. Note that Equation (1) is used as the minimization function for our genetic algorithm based combinatorial optimization.

D. The Genetic Algorithms

A genetic algorithm maintains a set of solutions known as population (Algorithm 1 Line 1). In each *generation*, the GA generates a new population from the current population using a given set of genetic operators known as *crossover* and *mutation* (Algorithm 1 Line 4). The algorithm then replaces the inferior solutions with the newly generated superior solutions to get a better current population (Algorithm 1 Line 6). A typical crossover operator splits two solutions at a randomly selected crossover point and exchanges parts between them (Figure 4a). Conversely, a typical mutation operator alters a solution at a random point (Figure 4b). Here the conformations are regarded as the individuals (chromosomes) of the GA population.

III. RELATED WORK

In the brief history of *ab initio* predictive methods, some successes have been achieved for high resolution-based real protein modeling [18], [19], [20], [21]. However, these successes are limited to the small sized proteins. Our low-resolution model is designed to target large sized proteins and

therefore, we did not compare our outcomes with the high resolution models at this stage.

Different types of meta-heuristics have been used in solving the protein structure prediction problem. These include Monte Carlo Simulation [45], Simulated Annealing [46], Genetic Algorithms (GA) [47], [48], Tabu Search with GA [35], Tabu Search with Hill Climbing [49], Ant Colony Optimization [50], Particle Swarm Optimization [51], [52], Immune Algorithms [53], Tabu-based Stochastic Local Search [30], [32], and Constraint Programming [54], [31].

Cebrian *et al.* [30] used tabu-based local search, and Shatabda *et al.* [32] used memory-based local search with tabu-based heuristics and achieved state-of-the-art results. However, Dotu *et al.* [31] used constraint programming and found promising results. To develop the Constraint-based Protein Structure Prediction (CPSP) tools, Mann *et al.* [54] applied an exact and complete algorithm for conformational search. The CPSP tools can find the optimal solution if the target protein has a matching hydrophobic-core stored in the CPSP database.

Amongst the population based approaches, Unger and Moult [47] applied genetic algorithms to PSP and found their method to be more promising than Monte Carlo based methods [45]. GAs have been used by Hoque *et al.* [34] for cubic and 3D HCP lattices. A twin-removal operator was introduced in [55] to remove duplicates from the GA population. The few other simplified methods for PSP are also available [56], [57], [58], [59] but they use different energy models [60], [61] and thus, we keep these methods out of our comparisons.

In 3D FCC lattice and HP energy based PSP, the recent state-of-the-art results have been achieved by tabu-guided local search algorithms [30], [32], [33] and by memetic algorithms [62]. In order to confirm the efficacy of our proposed algorithm, we have compared our results with the results of these state-of-the-art approaches.

IV. OUR APPROACH

In this section, we present an enhanced genetic algorithm (GAPlus) with detailed implementation of *i*) an exhaustive generation approach to diversify the search; *ii*) a novel hydrophobic core-directed macro-mutation operator to intensify the search; *iii*) a per-generation duplication elimination strategy to prevent early convergence; and *iv*) a random-walk technique to recover from stagnation.

A. GA Variants

The relations amongst the GA variants are shown in Figure 5. We have implemented five variants of GA for PSP to demonstrate the effectiveness of the enhanced features as listed below:

- 1) **RGA** (**R** denotes **R**andomness) is a typical implementation of a GA. It uses crossover and mutation operators. The crossover operator includes only single-point crossovers. The mutation operators include rotation, diagonal moves, pull moves, and tilt moves. The operator selection process is random. The crossover and mutation points are also selected randomly. RGA has a duplicate

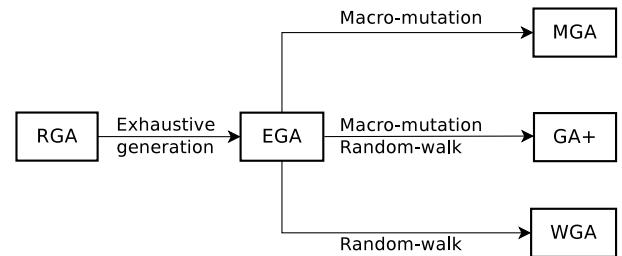


Figure 5. The relations amongst the GA variants. The RGA is the typical implementation of GA with random generation approach; the EGA inherits all of the features of the RGA with an additional exhaustive generation approach; the MGA inherits all of the features of the EGA with an additional macro-mutation operator; the WGA inherits all of the features of the EGA with an additional random-walk based stagnation recovery technique; and the GAPlus inherits all of the features of the EGA with the additional macro-mutation operator and random-walk based stagnation recovery technique.

elimination policy (see Section IV-B5). Further details are in Section IV-B.

- 2) **EGA** (**E** denotes **E**xhaustiveness) inherits all of the features of the RGA. Additionally, a novel exhaustive generation approach is used. In EGA, the random selection of crossover and mutation points is replaced with an exhaustive selection approach as described in Algorithm 9 and 10. The EGA is the base of WGA, MGA, and GAPlus. This variant demonstrates the effectiveness of our exhaustive generation approach. Further details are in Section IV-C2.
- 3) **MGA** (**M** denotes **M**acro-mutation) inherits all of the features of the EGA. Additionally, a novel macro-mutation operator is applied. The macro-mutation operator is used as a new mutation operator like other mutation operators (see Figure 6b-e). The macro-mutation operator helps form hydrophobic-cores that hide the hydrophobic amino acids from water and expose the polar amino acids to the surface to be in contact with the surrounding water molecules [63]. This variant demonstrates the effectiveness of our hydrophobic-core directed macro-mutation operator. Further details are in Section IV-C3.
- 4) **WGA** (**W** denotes **random-Walk**) inherits all of the features of the EGA. Additionally, at the point of stagnation, a novel random-walk based stagnation recovery technique is applied. This variant demonstrates the effectiveness of our random-walk based stagnation recovery approach. Further details are in Section IV-C4.
- 5) **GAPlus** inherits all of the features of EGA. Additionally, it includes the hydrophobic-core directed macro-mutation operator and the random-walk based stagnation recovery technique. The inclusion of these enhanced features helps our GAPlus produce state-of-the-art results (see Section V-A) for lattice-based PSP. The further details are presented in Section IV-C6.

The executable binaries of the GA variants are available online at <http://cs.uno.edu/~tamjid/Software.html>, under the heading: GAPlus.

Algorithm 2: RGA()

```

1 //op: Operators; c, c': Conformations
2 //opR: Operator selection rate
3 //curP,newP: Current & new population
4 //N: No. of amino acids in the sequence
5 //pos: Mutation or crossover point
6 curP ← initialize(repeat);
7 foreach (generation until timeout) do
8   selectOperator(op, opR);
9   if (mutation(op)) then
10     foreach (c ∈ curP) do
11       pos ← randMutationPoint(0, N - 1);
12       c' ← mutate(c, pos);
13       add(newP, c')
14   else
15     while (¬ full(newP)) do
16       c1, c2 ← randConfs(curP);
17       pos ← randCrossoverPoint(1, N - 2);
18       c'1, c'2 ← crossOver(c1, c2, pos);
19       add(newP, c'1, c'2)
20   curP ← newP
21 return bestConformation(curP)

```

B. The RGA framework

A typical randomized GA framework is implemented and denoted as the RGA. The *pseudocode* of RGA is presented in Algorithm 2. The RGA starts by generating an initial population (Line 6). The operator selection process is random in the RGA (Line 8). For mutation operators, each individual undergoes the offspring generation process (Line 10). For the crossover operator, the two parents are selected randomly (Line 16). The mutation (Line 11) and crossover (Line 17) points are also selected randomly. The randomized mutation or crossover operators are applied at Line 12 or Line 18 respectively. The duplicate elimination process is performed in Line 13 or Line 19. No explicit stagnation recovery technique is implemented in the RGA. The implementation details of the RGA are presented below:

1) *Implementation of the primitive operators:* In our GA, we implemented one primitive crossover operator and four primitive mutation operators (see Figure 6). The mutation operators are rotation, diagonal moves, pull moves, and tilt moves. The primitive operators are described below.

a) *Crossover operators:* The crossover operators are applied on two selected parent conformations to exchange their parts to generate child conformations. At a given crossover point (dotted circle in Figure 6a), two parent conformations exchange their parts and generate on the two child conformations. The success rate of the crossover operator decreases with increasing the compactness of the structure.

b) *Mutation operators:* The mutation operators are applied on a single conformation. The alteration can occur either on a single amino acid or on a series of amino acids of the conformation. The primitive mutation operators (as shown in Fig 6 (b-e)) are described below:

Algorithm 3: encode(*c*)

```

1 //c: Conformation; AA: Amino acid array
2 //S: Encoded string; s: Vector symbol
3 //N: No. of amino acids in the sequence
4 S ← {}//the length of S is (N - 1);
5 AA ← getAminoAcid(c);
6 for (i ← 1 to N - 1) do
7   vec ← getAbsoluteVector(AA[i], AA[i - 1]) ;
8   s ← getSymbol(vec) ;
9   S ← concat(S, s);
10 return S

```

Algorithm 4: decode(*S*)

```

1 //c: Conformation; AA: Amino acid array
2 //S: Encoded string; s: Vector symbol
3 //N: No. of amino acids in the sequence
4 N ← Length(S) + 1;
5 AA[0] ← (0, 0, 0);
6 for (i ← 1 to N - 1) do
7   s ← getSymbol(S, i);
8   vec ← getAbsoluteVector(s);
9   point ← getNextPoint(point, vec);
10  AA[i] ← point;
11  c ← buildConformation(AA);
12 evaluate(c);
13 return c ;

```

- 1) **Rotation:** One part of a given conformation is rotated around a selected point (Figure 6b). This move is mostly effective at the beginning of the search.
- 2) **Diagonal Move:** Given three consecutive amino acids at lattice points *A*, *B*, and *C*, a diagonal move at position *B* takes the corresponding amino acid diagonally to a free position (6c). Diagonal moves are very effective on FCC lattice [30], [31] points.
- 3) **Pull Moves:** The amino acids at points *A* and *B* are pulled to the free points (Figure 6d) and the connected amino acids are pulled as well to get a valid conformation. Pull moves [64] are local, complete, and reversible. Pull moves are very effective when the conformation becomes compact.
- 4) **Tilt Moves:** Two or more consecutive amino acids connected in a straight line are moved by a tilt move to immediately parallel lattice positions [48]. Tilt-moves pull the conformation from both sides until a valid conformation is found. In Figure 6e), the amino acids at points *C* and *D* are moved and, subsequently, other amino acids from both sides are moved as well.

- 2) *Conformation encoding and decoding:* We represent conformations with 3D coordinates and relative encodings. Each time, when a valid conformation is generated by applying genetic operators, we encode the conformation in a string of 12 different characters [A–L] representing the 12 absolute directions on a 3D FCC lattice space (see Section II-B). For each valid conformation, we store both the coordinates of the

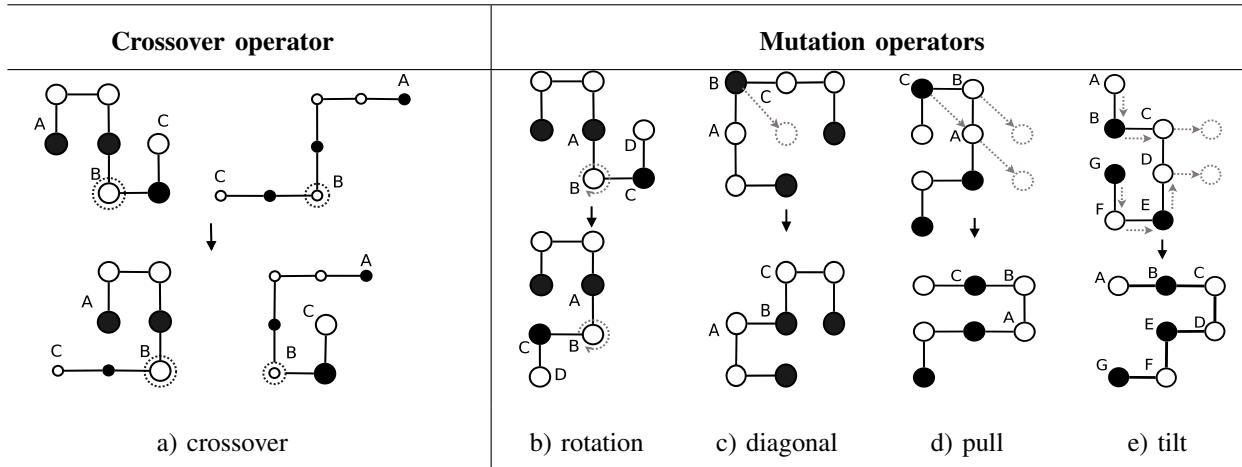


Figure 6. The operators that are used in our GAPlus on 3D FCC lattice space. For simplification, the figures are presented in 2D space. The solid circles represent the hydrophobic amino acids and the hollow circles represent the hydrophilic amino acids.

Algorithm 5: initialize(*repeat*)

```

1 //c: Conformation; AA: Amino acid array
2 //initP: Initial population
3 //N: No. of amino acids in the sequence
4 while ( $\neg$  full(newP)) do
5   AA[0]  $\leftarrow$  AminoAcid(0,0,0);
6   while (try  $\leq$  repeat) do
7     for (i  $\leftarrow$  1 to N - 1) do
8       k  $\leftarrow$  getRandom(1,12);
9       point  $\leftarrow$  AA[i-1] + basisVec[k - 1];
10      if ( $\neg$  free(point)) then break;
11      AA[i]  $\leftarrow$  AminoAcid(point);
12    if (mapping succeeded) then
13      c  $\leftarrow$  buildConformation(AA);
14    else
15      c  $\leftarrow$  deterministicConformation();
16    encode(c);
17    evaluate(c);
18    add(initP, c);
19 return initP;

```

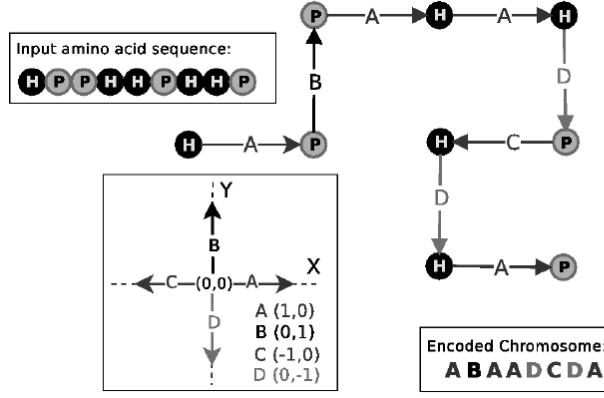
amino acids and the encoded string. While the coordinates help us determine whether a point on the lattice is free or not, the relative encodings help validate conformations and identify and eliminate duplicate ones.

a) *Conformation encoding*: In our implementation, we denote encoding as the procedure of generating a string to present a conformation using character notations, where each notation is the position of an amino acid in a 3D space. The first character of the encoded string will be the absolute direction of the second amino acid from the first amino acid. The second character of the encoded string will be the absolute direction of the third amino acid from the second amino acid, and so on. The length of the encoded string will be one less than the number of amino acids in the sequence (i.e., $N - 1$). Figure IV-B2b presents a simplified encoding procedure in 2D space. Algorithm 3 presents the encoding procedure.

b) *Conformation decoding*: In our implementation, we denote decoding as the procedure of extracting the coordinates of each amino acid from the encoded string. Considering (0, 0, 0) as the coordinate of the first amino acid of the sequence, the first character of the encoded string will provide the absolute direction of the second amino acid and hence the coordinate of the second amino acid is determined. The second character of the encoded string will provide the absolute direction of the third amino acid and hence the coordinate of the third amino acid is determined, and so on. Algorithm 4 presents the decoding procedures.

3) *Initializing the RGA population*: The RGA starts with an initial population, which is a set of feasible conformations. We generate initial conformations following a self-avoiding walk on the FCC lattice points. The *pseudocode* of the initialization algorithm is presented in Algorithm 5. It places the first amino acid at (0, 0, 0) (Line 5) and randomly selects a basis vector to place the next amino acid at a neighboring free lattice point (Lines 8 to 11). The amino acid position mapping continues until a self-avoiding walk is found for the whole protein sequence within a given number of iterations. The conformation is encoded (Line 16), evaluated (Line 17), and added into the initial population list (Line 18) if it does not exist in the list. If no valid conformation is found within the given number of iterations, a deterministic structure—an amino acid chain with unidirectional growth (e.g., AAAA... as the encoded chromosome)—is built (Line 15).

4) *Evaluating the solutions*: For each iteration, the conformation is evaluated by calculating the number of contacts (topological neighbors), where the two amino acids are non-consecutive and are both hydrophobic. The *pseudocode* in Algorithm 6 presents the procedure for calculating the contact energy of a given conformation. The contact potentials are found in the HP energy model as shown in Figure 3. The algorithm calculates the contacts based on Equation (1). The Euclidean distance between two topological neighbors on the FCC lattice space is $\sqrt{2}$. Therefore, all H amino acids which are non-consecutive and are $\sqrt{2}$ lattice distance apart are counted (Line 12) to get the fitness of the conformation.



Pos	Curr	Vec	Next
1		Start	H1(0, 0)
2	H1(0, 0)	+	A(1, 0) → P2(1, 0)
3	P2(1, 0)	+	B(0, 1) → P3(1, 1)
4	P3(1, 1)	+	A(1, 0) → H4(2, 1)
5	H4(2, 1)	+	A(1, 0) → H5(3, 1)
6	H5(3, 1)	+	D(0, -1) → P6(3, 0)
7	P6(3, 0)	+	C(-1, 0) → H7(2, 0)
8	H7(2, 0)	+	D(0, -1) → H8(2, -1)
9	H8(2, -1)	+	A(1, 0) → P9(3, -1)

Figure 7. The chromosome encoding procedures. For simplification, the figure is presented in 2D space with four possible vectors A-right, B-up, C-left, and D-down. In the chart, the *Pos*, *Curr*, *Vec*, and *Next* represent the amino acid positions, current amino acids, basis vectors, and next amino acids respectively. The next position (e.g. P2(1, 0)) is obtained adding the next direction vector (e.g. A(1, 0)) with the coordinates of the current amino acid (e.g. H1(0, 0)).

Algorithm 6: evaluate(*c*)

```

1 //c: Conformation; AA: Amino acid array
2 //N: No. of amino acids in the sequence
3 AA ← getAminoAcid(c);
4 fitness ← 0;
5 for (i ← 0 to N − 1) do
6   for (j ← i + 2 to N − 1) do
7     if (AcidType[i] = AcidType[i] = 'H') then
8       pointI ← AA[i];
9       pointJ ← AA[j];
10      sqrD ← getSqrDist(pointI, pointJ);
11      if (sqrD = 2) then
12        | fitness ← fitness − 1;
13 update(c, fitness);
14 return c;

```

Algorithm 7: add(*pop*, *c*)

```

1 //pop: population, c, c': Conformations
2 exist ← false;
3 foreach (c' ∈ pop) do
4   if (c' = c) then
5     | exist ← true;
6     | break;
7 if (¬exist) then pop ← add(conf);

```

5) *Removing the duplicate solutions:* In our implementation, no duplicate solution is allowed in the population in any instant of the search. Every time the RGA generates a solution, it is compared with the new population to check if it already exists. The new solution is included in the list if and only if the solution does not exist in the new population. In Algorithm 2, Line 13 or Line 19 prevents the identical—100% similar—solutions from being added to the new population by using Algorithm 7. This is how the duplicate-free new population becomes the current population in the next generation of the RGA (Line 20).

6) *Selecting the RGA operators:* The probability distribution to select operators is chosen intuitively. The crossover

and mutation operators are selected with the probabilities of 20% and 80% respectively. The four implemented primitive mutation operators are rotation, diagonal-move, pull-move, and tilt-move. For a particular generation, only one mutation operator is selected randomly from the four with a probability of 80%. At this stage of experimentation, we do not consider the optimality of operator selection probabilities. We present operator profiling in detail in Section V.

7) *Choosing the RGA population size:* The population size is chosen intuitively. For the time being, we use 100, 80, 60, 50 as the population sizes for the protein sequences having ≤ 50 , ≤ 100 , ≤ 200 , and ≤ 300 amino acids respectively. At this stage of experimentation we do not consider the optimality of GA population size. We present a rigorous study on GA population size in Section V.

8) *Housekeeping between generations:* The globally best (i.e. top elite) solution is always preserved. At the end of building a new generation, before shifting the new population to the current population, the best solution of the current population is compared with the globally best solution in terms of the free energy level. If the free energy of the current best is better than that of the globally best, then the globally best solution is replaced by the current best solution.

C. The GAPplus framework

We make three major enhancements on typical genetic algorithms: an exhaustive generation to diversify population, a macro-mutation operator to intensify the population, and a random-walk for stagnation recovery. We call the enhanced GA, GAPplus. The *pseudocode* of GAPplus is shown in Algorithm 8.

Note that our GAPplus is different from a typical GA (e.g., RGA) in a number of ways. The distinguishable difference between RGA and GAPplus are listed below:

- At Lines 11 and 18, the GAPplus applies crossover or mutation operators exhaustively at each amino acid position (Algorithm 9 Line 3, Algorithm 10 Line 3). This exhaustiveness is the core feature of EGA.

Algorithm 8: GAPlus()

```

1 //op: Operators, c, c': Conformations
2 //opR: Operator selection probabilities
3 //curP, newP: Current & new populations
4 //rwT: Non-improving generation counter
5 curP ← initialize(repeat);
6 foreach (generation until timeout) do
7   selectOperator(op, opR);
8   if (mutation(op)) then
9     foreach (c ∈ curP) do
10    if (¬ isMacro) then
11      | c' ← mutConf(c);
12    else
13      | c' ← macroMutation(c, repeat);
14    add(newP, c')
15   else
16     while (¬ full(newP)) do
17       | c1, c2 ← randomConfs(curP);
18       | c'1, c'2 ← crossConf(c1, c2);
19       | add(newP, c'1, c'2)
20   if (¬ improved(newP, rwT)) then
21     | curP ← randomWalk(newP)
22   else
23     | curP ← newP
24 return bestConformation(curP)

```

Algorithm 9: mutConf(c)

```

1 // c: Conformation; pos: Mutation point
2 addToList(list, c);
3 for (pos ← 0 to N – 1) do
4   | c' ← applyOperator(c, pos);
5   | addToList(list, c');
6 return bestConformation(list)

```

Algorithm 10: crossConfs(c₁, c₂)

```

1 // c1, c2, c'1, c'2: Conformations; pos: Crossover point
2 addToList(list, c1, c2);
3 for (pos ← 1 to N – 2) do
4   | c'1, c'2 ← applyOperator(c1, c2, pos);
5   | addToList(list, c1, c2, c'1, c'2);
6 return bestTwoConformations(list)

```

- 2) Besides primitive GA operators, at Line 13, the GAPlus introduces a macro-mutation operator (Algorithm 11) in a similar fashion of the other mutation operators. This macro-mutation is the core feature of MGA.
- 3) Besides duplicate elimination, at Line 21, the GAPlus uses a random-walk algorithm (Algorithm 12) to recover from stagnation. This random-walk is the core feature of WGA.

1) *Different selection strategies in GAPlus:* All of the following processes in the GAPlus framework involve their purpose-specific selection strategies:

a) *Population size selection:* For the time being, we choose the population size as used in RGA (see Sec-

tion IV-B7). We present a detailed analysis on the GA population size in Section V.

b) *Operator selection:* In our implementation, we intuitively assign a 20% selection probability for one crossover operator and 80% for five mutation operators (including the macro-mutation) with equal selection probabilities (20% of the remaining 80% for each mutation operator). Only the selected operator is applied to build the population for the next generation.

c) *Parent selection:* If the selected operator is mutation then every individual in the population has been passed through the mutation operator to build the population for the next generation. However, for the crossover operator, two parents are selected randomly until all individuals in the population have been passed through the crossover operator to build the population for the next generation.

d) *Point of operation selection:* The point of operation is chosen exhaustively starting from the first to the last amino acid of the sequence for the mutation operator and from the second to the second-last amino acid of the sequence for the crossover operator until finding the acceptable children.

e) *Next generation selection:* The fitness is the measure of selecting a solution for the next generation. When mutation is chosen to build the next generation, the parent and all the children resulting from the exhaustive mutation are added to a list. The fittest solution from the list then replaces the parent. Similarly, when crossover is chosen, the two parents and all the children resulting from the exhaustive crossover are added to a list. The two best fit solutions from the list then replace the two parents.

f) *Random-walk timing selection:* We evaluate and compare the current best solution with the so-far global-best solution at the end of each generation. When the global-best solution remains the same for a certain number of generations (a random value between 5 and 10 inclusive is considered), a random-walk is triggered.

2) *Exhaustive generation based approach:* To reduce the randomness, we introduce an exhaustive generation approach in the parent selection as well as in the mutation and crossover point selection.

a) *Exhaustive mutation:* For mutation operators, our algorithm adds one resultant conformation to the new population for each of the conformations in the current population. In Algorithm 9, the child conformations are generated by applying the genetic operator at each of the positions of the parent conformation (Line 4). The resultant conformation of a mutation operation is either the parent conformation itself or a child depending on the quality (fitness and structure) of the conformations.

b) *Exhaustive crossover:* For crossover operators, two resultant conformations are added to the new population for the two randomly selected parents from the current population. Crossover operators (Algorithm 10) generate child conformations by exhaustively selecting the crossover points on the parent conformations. In the GAPlus, a number of child conformations are generated and listed (Line 5). The two best conformations in the list then become the resultant conformations (Algorithm 10 Line 6).

c) *Effectiveness of exhaustive generation:* To show the effectiveness of exhaustive generation experimentally, we tested our approach on 10 small and medium sized benchmark proteins taken from the literature. In Table I, we present the results obtained from two variants of GA: RGA that uses the random selection approach and EGA that uses the exhaustive generation approach. From the experimental results, it is clear that the EGA significantly outperforms the RGA and, in some cases, the EGA is able to reach the lower bound of the free energy as denoted by "*" in Column *Best*, under the header *EGA*. In Table I, we present the global minimum energies in Column *Best*, the average of minimum energies in Column *Mean*, and the standard deviation of the minimum energies in Column *STD*. These values are obtained from 50 different trials in identical settings for all proteins. For both cases, the EGA performs better than RGA. This consistent performance demonstrates the effectiveness of the exhaustive generation.

3) *Hydrophobic-core directed macro-mutation:* In this section, we present a macro-mutation operator guided by the hydrophobic property of amino acids. The macro-mutation operator compresses the conformation and quickly forms the hydrophobic-core. In GAPlus, the macro-mutation operator is applied in a similar manner of the primitive mutation operators (Algorithm 8 - Line 13).

a) *Hydrophobic-core:* Protein structures have hydrophobic-cores that hide the hydrophobic amino acids from water and expose the polar amino acids to the surface to be in contact with the surrounding water molecules [63]. The hydrophobic-core formation is an important objective of HP energy model based protein structure prediction.

b) *Hydrophobic-core center:* The hydrophobic-core center is a virtual point in a 3D space which is the centroid of the core formed by the hydrophobic amino acids in a given protein sequence. In the macro-mutation (Algorithm 11), the hydrophobic core center (HCC) is calculated by finding the arithmetic means of x , y , and z coordinates of all hydrophobic amino acids as shown in Equation (2).

$$x_{\text{hcc}} = \frac{1}{n_h} \sum_{i=1}^{n_h} x_i, \quad y_{\text{hcc}} = \frac{1}{n_h} \sum_{i=1}^{n_h} y_i, \quad z_{\text{hcc}} = \frac{1}{n_h} \sum_{i=1}^{n_h} z_i \quad (2)$$

where n_h is the number of H amino acids in the protein.

Table I

THE ENERGY VALUES OBTAINED BY RUNNING RGA AND EGA TO DEMONSTRATE THE EFFECTIVENESS OF THE EXHAUSTIVE GENERATION. COLUMN *LB-FE* PRESENTS THE LOWER BOUND OF FREE ENERGIES.

Protein details			RGA			EGA			Time
Seq	Size	Hs	Best	Mean	STD	Best	Mean	STD	mins
H1	48	24	-69	-65	-58	2.63	*	-67	1.05
H2	48	24	-69	-63	-57	3.16	*	-67	0.93
H3	48	24	-72	-64	-57	3.49	-71	-69	1.29
H4	48	24	-71	-66	-58	3.55	*	-69	1.31
H5	48	24	-70	-65	-58	3.85	*	-68	1.10
F90_1	90	50	-168	-133	-120	5.73	-159	-144	6.54
F90_2	90	50	-168	-132	-124	6.09	-155	-142	5.84
F90_3	90	50	-167	-138	-124	8.00	-158	-146	5.78
F90_4	90	50	-168	-141	-124	7.01	-158	-146	5.03
F90_5	90	50	-167	-131	-120	7.02	-158	-144	6.43

* denotes the lower bound of free energy is found.

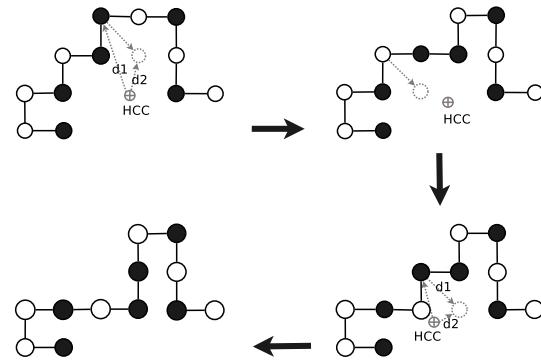


Figure 8. A macro-mutation operator comprising a series of diagonal moves. For simplification, the figures are presented in 2D space.

The HCC is used to calculate the Euclidean distance of all H amino acids using Equation (3).

$$d_i = \sqrt{(x_i - x_{\text{hcc}})^2 + (y_i - y_{\text{hcc}})^2 + (z_i - z_{\text{hcc}})^2} \quad (3)$$

where d_i is the Euclidean distance of the i th amino acid at (x_i, y_i, z_i) from the HCC at $(x_{\text{hcc}}, y_{\text{hcc}}, z_{\text{hcc}})$.

c) *Macro-mutation operator:* The macro-mutation operator is a composite operator that performs a series of diagonal moves on a given conformation to build the hydrophobic-core around the HCC as shown in Figure 8. The *pseudocode* of the macro-mutation operator is shown in Algorithm 11.

The diagonal moves are applied repeatedly either at each P- or H-type amino acid position. Whether to apply the diagonal move on P- or H-type amino acids is determined by using a *Bernoulli* distribution (Line 5) with probability p (intuitively we use $p = 20\%$ for P-type amino acids). For a P-type amino acid, the first successful diagonal move is considered (Line 13). However, for an H-type amino acid, the first successful diagonal move (Line 24) that does not increase the Euclidean distance (Line 23) of the amino acid from the HCC, is taken. All the amino acids are traversed and the successful moves are considered as a single composite move. In this way, the macro-mutation operator compresses the conformation and quickly forms the hydrophobic-core by repeating the procedure (Line 4).

d) *Effectiveness of macro-mutation:* In Table II, we present the results obtained from two variants of GA: EGA that uses an exhaustive generation approach only and MGA that uses all features of EGA along with the macro-mutation operator. From the experimental results, it is clear that the MGA significantly outperforms the EGA. In Table II, we present the global minimum energy (Column *Best*), the average of minimum energies (Column *Mean*), and the standard deviation of the minimum energies (Column *STD*) obtained from 50 different runs of identical settings for all proteins using these two GA variants. For both cases, the MGA performs better than the EGA. This consistent performance demonstrates the effectiveness of the macro-mutation operator in GA.

4) *Random-walk based stagnation recovery:* Search-based optimization algorithms often suffer from stagnancy [65], [66],

Algorithm 11: macroMutation(c , $repeat$)

```

1 //  $c$ : Conformation
2  $hAA \leftarrow$ listHydrophobicAminoAcid( $c$ );
3  $pAA \leftarrow$ listPolarAminoAcid( $c$ );
4 for ( $i \leftarrow 1$  to  $repeat$ ) do
5   if ( $bernoulli(p)$ ) then
6     |  $T \leftarrow P$ 
7   else
8     |  $T \leftarrow H$ 
9   if ( $T = P$ ) then
10    | foreach ( $aa \in pAA$ ) do
11      | |  $point \leftarrow$ findFreePoint( $aa$ );
12      | | if ( $\neg empty(point)$ ) then
13        | | | applyDiagonalMove( $aa, point$ );
14        | | | updateConf( $c, aa, point$ );
15        | | | break;
16   else
17     | |  $hcc \leftarrow findHCC()$ ;
18     | foreach ( $aa \in hAA$ ) do
19       | |  $d_{old} \leftarrow$ getDistance( $aa, hcc$ );
20       | |  $point \leftarrow$ findFreePoint( $aa$ );
21       | | if ( $\neg empty(point)$ ) then
22         | | |  $d_{new} \leftarrow$ getDistance( $point, hcc$ );
23         | | | if ( $d_{new} \leq d_{old}$ ) then
24           | | | | applyDiagonalMove( $aa, point$ );
25           | | | | updateConf( $c, aa, point$ );
26           | | | | break;
27   encode( $c$ );
28 evaluate( $c$ );
29 return  $c$ ;

```

Table II

THE ENERGY VALUES OBTAINED BY RUNNING EGA AND MGA TO DEMONSTRATE THE EFFECTIVENESS OF THE MACRO-MUTATION OPERATOR. COLUMN LB-FE PRESENTS THE LOWER BOUND OF FREE ENERGIES.

Protein details			EGA			MGA			Time	
Seq	Size	Hs	Best	Mean	STD	Best	Mean	STD	mins	
H1	48	24	-69	*	-67	1.05	*	-67	0.85	30
H2	48	24	-69	*	-67	0.93	*	-68	0.83	30
H3	48	24	-72	-71	-69	1.29	*	-71	1.06	30
H4	48	24	-71	*	-69	1.31	*	-70	1.06	30
H5	48	24	-70	*	-68	1.10	*	-69	0.75	30
F90_1	90	50	-168	-159	-144	6.54	-165	-159	2.69	120
F90_2	90	50	-168	-155	-142	5.84	-164	-158	2.82	120
F90_3	90	50	-167	-158	-146	5.78	-163	-158	2.69	120
F90_4	90	50	-168	-158	-146	5.03	-165	-160	2.30	120
F90_5	90	50	-167	-158	-144	6.43	-166	-160	2.47	120

* denotes the lower bound of free energy is found.

[67]. For population based algorithms, it occurs when the algorithms fail to build new individuals [68]. While exploring a search landscape, stagnation occurs when a search algorithm gets stuck or stalls in a local minimum, gets trapped in a valley, or gets lost in a plateau. This happens often in solving hard optimization problems such as protein structure prediction, particularly for large sized proteins. Therefore, developing an effective stagnation recovery technique is crucial for the conformational search.

a) *Random-walk*: We define a random-walk as a process of exploring feasible solutions in the vicinity of a base solution with minimal changes made to it. Both good and bad solutions are taken into consideration during those changes. The process is repetitive and increases the dissimilarity of the new solution with the base solution.

b) *Stagnation in GA*: As the search progresses, the GA population becomes richer and richer from one generation to the next. From experimental results, it has been found that after a certain number of generations, a significant portion of the individuals in the population become identical or almost identical. At this situation, the GA operators fail to generate new improved solutions from similar or almost similar parents. Therefore, the search gets stuck and stagnation arises. We evaluate and compare the current best solution with the so-far global-best solution at the end of each generation. When the global-best solution found remains the same for a number of generations (Line 20 in Algorithm 8), we consider this stagnation.

c) *Random-walk at stagnation*: In a stagnation situation, a random-walk technique is applied to perform unconditional pull moves (see Figure 6d) on each conformation of the new population to diversify the population. The *pseudocode* of the random-walk is shown in Algorithm 12. A local minimum or relative minimum in PSP is encountered when a premature hydrophobic-core is formed. To break the premature hydrophobic-core, we apply a random-walk algorithm. In random-walk, we use pull-moves (Line 6). During pulling, energy levels and structural diversification are observed to maintain a balance among the two. In a controlled random-walk (Line 7), we allow the energy level to change between 5% and 10% and structural change between 10% and 75% of the original. We accept the conformation that has the minimum differential energy level, but has the maximum structural diversity with respect to the original one.

d) *Effectiveness of random-walk*: In Table III, we present the results obtained from two variants of GA: the EGA and the WGA. From the experimental results, it is clear that the WGA significantly outperforms the EGA. In Table III, we present the global minimum energy (Column *Best*), the average of minimum energies (Column *Mean*), and the standard deviation of the minimum energies (Column *STD*) obtained from 50 different runs of identical settings for all proteins using these two GA variants. For both the cases, the WGA performs better than EGA. This consistent performance demonstrates the

Algorithm 12: randomWalk(pop)

```

1 //  $pop$ : Population;  $c$ : Conformation
2 foreach ( $c \in pop$ ) do
3    $isFound \leftarrow$  false;
4   while ( $\neg isFound$ ) do
5     | for ( $pos \leftarrow 0$  to  $N - 1$ ) do
6       | | applyPullMove( $c, pos$ );
7       | |  $isFound \leftarrow$  checkDiversity( $c$ );
8       | | updatePopulation( $pop, c$ );
9   return  $pop$ ;

```

Table III

THE ENERGY VALUES OBTAINED BY RUNNING EGA AND WGA TO DEMONSTRATE THE EFFECTIVENESS OF RANDOM-WALK BASED STAGNATION RECOVERY TECHNIQUE. COLUMN LB-FE PRESENTS THE LOWER BOUND OF FREE ENERGIES.

Protein details			EGA			WGA			Time	
Seq	Size	Hs	LB-FE	Best	Mean	STD	Best	Mean	STD	mins
H1	48	24	-69	*	-67	1.05	*	-68	0.67	30
H2	48	24	-69	*	-67	0.93	*	-68	0.96	30
H3	48	24	-72	-71	-69	1.29	*	-70	0.87	30
H4	48	24	-71	*	-69	1.31	*	-70	0.86	30
H5	48	24	-70	*	-68	1.10	*	-69	0.73	30
F90_1	90	50	-168	-159	-144	6.54	-165	-161	2.53	120
F90_2	90	50	-168	-155	-142	5.84	-166	-161	2.26	120
F90_3	90	50	-167	-158	-146	5.78	-165	-161	2.36	120
F90_4	90	50	-168	-158	-146	5.03	-166	-161	2.42	120
F90_5	90	50	-167	-158	-144	6.43	-166	-162	2.48	120

* denotes the lower bound of free energy is found.

effectiveness of the random-walk based stagnation recovery technique within the GA framework.

5) *Housekeeping between generations*: As mentioned in Section IV-B8, if an improvement has occurred, the non-improving time tracker and non-improving generation counter are reset. Conversely, if the current best solution is worse or comparable to the globally best one, the non-improving time tracker and non-improving generation counter are incremented. Finally, if the non-improving time tracker or non-improving generation counter reach their threshold values, the random-walk process is activated to create diversity in the population.

6) *Combined efforts of all of the enhancements*: In our multi-step enhancement, the final outcome is GAPlus. The GAPlus inherits all features of the EGA. Additionally, it includes the hydrophobic-core directed macro-mutation operator and the random-walk based stagnation recovery technique. The inclusion of the exhaustive generation approach, the random-walk technique, and the macro-mutation operator helps our GAPlus produce the state-of-the-art results in HP based on-lattice protein structure prediction.

7) *Effectiveness of the enhanced features*: In Tables IV and V, we present the results obtained from four variants of GA: EGA, WGA, MGA, and GAPlus. From the experimental results, it is clear that the GAPlus significantly outperforms the other variants. The GAPlus performs better than others in terms of the global minimum free energy and the average free energy for all of the 35 benchmark proteins. This consistent performance demonstrates the effectiveness of the combined efforts of the exhaustive generation, the random-walk, and the macro-mutation operator in GA.

8) *Observing GAPlus features*: In Figure 9, when random-walk occurs the population becomes more diversified. Because of this, the average free energy for the specific generation increases. Similarly, when the macro-mutation applied, the average quality of the individuals improves quickly. Because of this, the average free energy for the specific generation decreases.

D. Performance hierarchies amongst the GA variants

The performance hierarchies amongst the GA variants are presented in Figure 10. Because of the exhaustive generation

technique, the EGA performs better than the RGA. Because of the exhaustive generation in combination with the macro-mutation operator, the MGA performs better than the RGA and the EGA. Because of the exhaustive generation in combination with the random-walk based stagnation recovery technique, the WGA performs better than the RGA, the EGA, and the MGA. Finally, because of the combination of all three enhancements, the GAPlus performs better than the others.

For further investigation, we perform rigorous analysis on the results obtained from the GA variants and also compare our experimental results with the state-of-the-art approaches for the same models in Section V.

V. DETAILED EXPERIMENTATION AND ANALYSES

In this section, we evaluate the performance of the GAPlus by comparing the results with that of the state-of-the-art approaches. We also perform an extensive experimental study on the GAPlus parameters to find an optimum set of parameters.

A. Experimental settings

Our implementation is based on the *Java Programming Language*. We ran the experiments on the NICTA [69] cluster. This cluster consists of a number of identical Dell PowerEdge R415 computers, each equipped with $2 \times$ AMD 6-Core Opteron 4184 processors, 2.8GHz clock speed, 3M L2/6M L3 Cache, and 64 GB memory and running Rocks OS (a Linux variant for cluster).

B. The benchmark proteins used in GAPlus evaluation

In our experiment, the protein instances that we use as benchmarks are taken from the literature. The H, F90, S, F180, and R instances are taken from the Peter Clote laboratory website [70]. These benchmarks are used in [30], [31], [32], [33] to test their algorithms. Additionally, we also use another six large sized sequences that are taken from the Critical Assessment of Protein Structure Prediction (CASP) competition [71]. The CASP9 targets are 3mse, 3mr7, 3mqz, 3no6, 3no3, and 3on7 which are also used in [32], [33]. To fit in the HP model, the CASP targets are converted into the HP sequences based on the hydrophobic properties (see Section II-C) of the constituent amino acids.

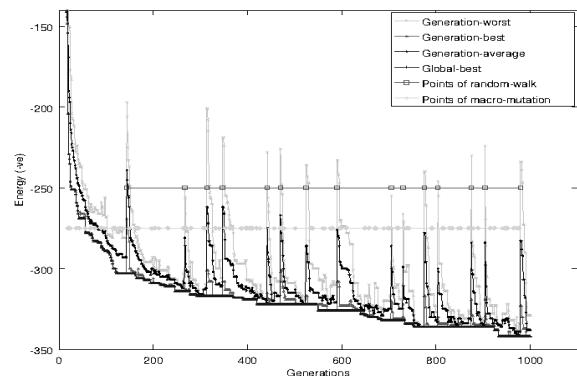


Figure 9. The activities while GAPlus is in operation. The figure presents the progressive search of a sample run of 1000 generation for Protein R1.

Table IV

FOR SMALL SIZED PROTEINS, THE EXPERIMENTAL RESULTS OF GA VARIANTS, AND THE LOCAL SEARCH ALGORITHM. THE COLUMNS *Size* AND *LBFE* PRESENT THE LENGTH OF THE AMINO ACID SEQUENCES AND LOWER BOUND OF FREE ENERGIES RESPECTIVELY.

Protein details			GA variants energy values												State-of-the-art			Time (min)
			RGA		EGA		WGA		MGA		GAPlus			TLS				
Seq	Size	Hs	LBFE	Best	Avg	Best	Avg	Best	Avg	Best	Avg	STD	Best	Avg	STD	Best	Avg	
H1	48	24	-69	-65	-58	*	-67	*	-68	*	-67	*	-69	0.27	-68	-66	1.13	30
H2	48	24	-69	-63	-57	*	-67	*	-68	*	-68	*	-69	0.24	=	-65	1.60	30
H3	48	24	-72	-64	-57	-71	-69	*	-70	*	-71	*	-72	0.42	-69	-66	1.63	30
H4	48	24	-71	-66	-58	*	-69	*	-70	*	-70	*	-71	0.14	-70	-65	2.05	30
H5	48	24	-70	-65	-58	*	-68	*	-69	*	-69	*	-70	0.00	-68	-65	1.70	30
H6	48	24	-70	-63	-58	-69	-67	*	-69	*	-68	*	-69	0.45	-69	-66	1.63	30
H7	48	24	-70	-66	-58	*	-68	*	-69	*	-69	*	-70	0.00	-69	-66	1.88	30
H8	48	24	-69	-62	-55	*	-66	*	-68	*	-68	*	-69	0.14	-67	-64	1.76	30
H9	48	24	-71	-65	-59	*	-69	*	-70	*	-70	*	-71	0.27	-68	-66	1.26	30
H10	48	24	-68	-63	-57	*	-67	*	-67	*	-67	*	-68	0.00	*	-65	1.65	30
F90_1	90	50	-168	-133	-120	-159	-144	-165	-161	-165	-159	*	-166	1.62	-164	-160	2.70	120
F90_2	90	50	-168	-132	-124	-155	-142	-166	-161	-164	-158	*	-165	1.38	-165	-158	3.78	120
F90_3	90	50	-167	-138	-124	-158	-146	-165	-161	-163	-158	*	-164	1.37	-165	-159	4.13	120
F90_4	90	50	-168	-141	-124	-158	-146	-166	-161	-165	-160	*	-165	1.40	-165	-159	3.58	120
F90_5	90	50	-167	-131	-120	-158	-144	-166	-162	-166	-160	*	-166	0.96	-165	-159	4.12	120

* denotes the lower bound of free energy is found.

Table V

FOR LARGE SIZED PROTEINS, THE EXPERIMENTAL RESULTS OF GA VARIANTS, AND THE LOCAL SEARCH ALGORITHM. THE COLUMNS *Size* AND *LBFE* PRESENT THE LENGTH OF THE AMINO ACID SEQUENCES AND LOWER BOUND OF FREE ENERGIES RESPECTIVELY.

Protein details			GA variants energy values												State-of-the-art			Time (min)				
			RGA		EGA		WGA		MGA		GAPlus			TLS [30]			MLS1 [32]		MLS2 [33]			
Seq	Size	Hs	LBFE	Best	Avg	Best	Avg	Best	Avg	Best	Avg	STD	Best	Avg	STD	Best	Avg	Best	Avg			
S1	135	100	-357	-300	-278	-332	-312	-344	-336	-352	-342	-355	-348	3.51	-351	-341	7.72	-	-	-	120	
S2	151	100	-360	-298	-258	-332	-300	-346	-334	-351	-339	-356	-349	4.64	-355	-343	9.97	-	-	-	120	
S3	162	100	-367	-290	-250	-322	-297	-347	-334	-347	-335	-361	-348	6.93	-355	-340	10.62	-	-	-	120	
S4	164	100	-370	-273	-246	-317	-291	-350	-333	-348	-337	-364	-352	6.72	-354	-343	9.05	-	-	-	120	
F180_1	180	100	-378	-256	-229	-309	-280	-330	-315	-347	-333	-351	-341	4.60	-338	-326	5.74	-360	-335	-361	-341	300
F180_2	180	100	-381	-262	-236	-314	-283	-341	-328	-350	-335	-362	-346	5.24	-345	-333	6.21	-362	-340	-368	-350	300
F180_3	180	100	-378	-261	-236	-315	-288	-352	-334	-356	-340	-361	-350	5.354	-352	-338	6.86	-357	-343	-365	-355	300
R1	200	100	-384	-249	-221	-295	-273	-345	-326	-353	-331	-355	-345	6.074	-332	-318	6.80	-353	-325	-359	-339	300
R2	200	100	-383	-262	-219	-302	-276	-345	-327	-351	-333	-360	-346	7.11	-337	-324	6.75	-351	-329	-361	-343	300
R3	200	100	-385	-250	-220	-299	-274	-340	-323	-344	-330	-363	-344	8.21	-339	-323	8.25	-352	-329	-354	-340	300
3MSE	179	84	-323	-198	-178	-221	-206	-271	-249	-278	-268	-292	-278	6.19	-268	-251	7.32	-278	-255	-288	-271	300
3MR7	189	93	-355	-226	-208	-252	-233	-299	-288	-325	-309	-330	-316	6.05	-304	-287	5.90	-311	-292	-320	-304	300
3MQZ	215	120	-474	-303	-276	-338	-310	-389	-369	-412	-398	-427	-412	7.023	-404	-384	7.96	-415	-387	-430	-404	300
3NO6	229	116	-455	-274	-251	-340	-310	-395	-372	-399	-384	-423	-402	9.59	-390	-372	7.88	-400	-374	-419	-397	480
3NO3	258	122	-494	-288	-250	-359	-316	-394	-375	-406	-387	-421	-404	8.51	-388	-359	9.79	-379	-361	-406	-376	480
3ON7	279	146	<i>u/k</i>	-359	-315	-428	-389	-482	-458	-494	-473	-519	-490	12.94	-491	-461	9.32	-499	-463	-514	-476	480

u/k denotes unknown.

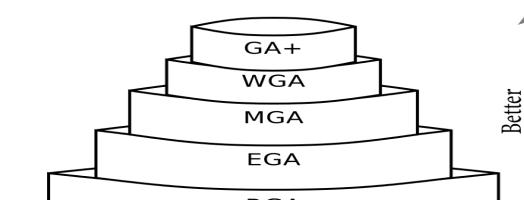


Figure 10. The performance hierarchy within the GA variants. The higher the position the better the performance.

C. State-of-the-art methods used to evaluate GAPlus

In order to measure the efficiency of our enhanced GA framework (GAPlus), we compare the results with a tabu-guided local search algorithm [30] which is denoted as TLS, a memory-based tabu-guided local search [32] which is denoted as MLS1, an enhanced memory-based tabu-guided local search [33] which is denoted as MLS2, and a population-based

memetic algorithm [62] which is denoted as MA in the rest of the sections. We tried to run the algorithm proposed in [31], but, unfortunately for most of the proteins, the program aborted showing the memory related errors. Any effective comparison in this case is therefore not possible.

D. Comparative analyses for evaluation

We ran the five variants of GA: RGA, EGA, MGA, WGA, and GAPlus (see Section IV-A for details). For each protein, we run each algorithm 50 times with identical experimental settings and with the time limits specified in Tables IV and V under Column *Time*.

1) *Comparing our results with the state of the art results:* Based on the size of the protein sequences, we present the comparative results in two different tables: Table IV presents the small sized proteins (*size* < 100) and Table V presents the large sized proteins (*size* > 100). The lower bounds of the free energy (LBFE) values (Column *LBFE* in Tables IV and V) are obtained from [30], [32]; however, there are some

Table VI

COMPARATIVE RESULTS OF GAPLUS WITH THE RECENTLY PUBLISHED REPORTED RESULTS USING MEMETIC ALGORITHM OVER 8 SMALL SIZED PROTEINS

Protein details	Our Approach				State-of-the-art						
	GAPlus				MA						
Seq	Size	Hs	LBFE	Best	Mean	STD	Worst	Best	Mean	STD	Worst
B1	20	12	-28	-28	-28.00	± 0.00	-28	-23	-23.00	± 0.00	-23
B2	24	10	-23	-23	-23.00	± 0.00	-23	-23	-23.00	± 0.00	-23
B3	25	9	-17	-17	-17.00	± 0.00	-17	-17	-17.00	± 0.00	-17
B4	36	16	-38	-38	-38.00	± 0.00	-38	-38	-38.00	± 0.00	-38
B5	48	25	-68	-68	-67.96	± 0.20	-67	-68	-66.94	± 0.28	-65
B6	50	32	-97	-99	-97.32	± 1.12	-95	-71	-69.15	± 0.37	-67
B7	60	32	-98	-130	-129.72	± 0.61	-128	-128	-125.65	± 0.34	-124
B8	64	42	-130	-132	-131.00	± 1.12	-129	-128	-124.67	± 0.45	-123

unknown values of lower bounds of free energy for larger sequences (presented as u/k in Table V).

a) *Comparing on small sized benchmark proteins:* The experimental results in Table IV show that for H instances ($Size = 48$): amongst the GA variants, the EGA is able to achieve the *LBFE* for 8 proteins out of 10 and the WGA, MGA, and GAPlus are able to hit the *LBFE* for all 10 proteins. Moreover, the average energy obtained by GAPlus demonstrates that our final version of GA is able to reach native energy levels in almost every single run for H instances. On the other hand, for F90 instances ($Size = 90$): amongst the GA variants, WGA and GAPlus performs better than the state-of-the-art approaches, however, only the GAPlus is able to hit the *LBFE* for all 5 proteins. The performances of the MGA and the TLS are very similar.

b) *Comparing on large sized benchmark proteins:* Notice that the GAPlus, which benefits from all our three novel techniques, clearly outperforms RGA, EGA, MGA, and WGA. Nevertheless, we observe that the results obtained by GAPlus are better than those of TLS and MLS1 with wide margins for all (16 out of 16) benchmark proteins. TLS and MLS1 are also outperformed by WGA, but they outperform EGA and RGA; the results of TLS and MGA are very close. On the other hand, MLS2 produces competitive results for few cases (3 out of 16), but is outperformed by the GAPlus for the rest of the cases.

Furthermore, as shown in Table VI, we ran GAPlus for another set of benchmark proteins taken from [62]. We compared our results with the reported values in [62] produced by a memetic algorithm (MA). For a set of very small proteins ($Size < 50$), the MA produces competitive results but fails to compete with the GAPlus for a slightly larger set of proteins ($Size > 50$).

2) *Relative improvement:* As the predicted energy approaches the *LBFE*, further improvement becomes more difficult. The relative improvement (RI) explains how close the predicted results (target) are to the *LBFEs* with respect to the energy values scored by the state-of-the-art (reference) approaches. In Table VII Columns *RI*, we present a comparison of (%) improvements on the average conformation quality. We compare target GAPlus with references RGA, TLS, MLS1, and MLS2. For each protein, the relative improvement of the target (*t*) with respect to the reference (*r*) is calculated using Equation (4). We present the relative improvements only for

Table VII

THE RELATIVE IMPROVEMENTS (*RI* COLUMNS) OF GAPLUS OVER RGA, TLS, AND MLS2. THE VALUES ARE CALCULATED USING THE FORMULA EXPLAINED IN EQUATION (4). COLUMN *LBFE* PRESENTS THE LOWER BOUND OF FREE ENERGIES.

Protein details	GA+			RGA		TLS		MLS1		MLS2	
	(<i>t</i>)	(<i>r</i>)	[30] (<i>r</i>)	Avg	RI	Avg	RI	Avg	RI	Avg	RI
H1	48	24	-69	-68	-57	92%	-66	67%			
H2	48	24	-69	-68	-57	92%	-65	75%			
H3	48	24	-72	-71	-56	94%	-65	86%			
H4	48	24	-71	-70	-57	93%	-65	83%			
H5	48	24	-70	-70	-57	100%	-65	100%			
F90_1	90	50	-168	-165	-119	94%	-159	67%			
F90_2	90	50	-168	-165	-123	94%	-157	72%			
F90_3	90	50	-167	-164	-124	93%	-158	67%			
F90_4	90	50	-168	-164	-123	91%	-159	56%			
F90_5	90	50	-167	-165	-120	96%	-158	78%			
S1	135	100	-357	-348	-278	89%	-341	44%			
S2	151	100	-360	-349	-258	89%	-343	35%			
S3	162	100	-367	-348	-250	84%	-340	30%			
S4	164	100	-370	-352	-246	85%	-343	33%			
F180_1	180	100	-378	-341	-229	75%	-326	29%	-335	14%	-341
F180_2	180	100	-381	-346	-236	76%	-333	27%	-340	15%	-350
F180_3	180	100	-378	-350	-236	80%	-338	30%	-343	20%	-355
R1	200	100	-384	-345	-221	76%	-318	41%	-325	34%	-339
R2	200	100	-383	-346	-219	77%	-324	37%	-329	31%	-343
R3	200	100	-385	-344	-220	75%	-323	33%	-329	27%	-340
3MSE	179	84	-323	-278	-178	69%	-251	36%	-255	34%	-271
3MR7	189	93	-355	-316	-208	73%	-287	43%	-292	38%	-304
3MQZ	215	120	-474	-412	-276	69%	-384	31%	-387	29%	-404
3NO6	229	116	-455	-402	-251	74%	-372	36%	-374	35%	-397
3NO3	258	122	-494	-404	-250	63%	-359	33%	-361	32%	-476

× denotes "no improvement".

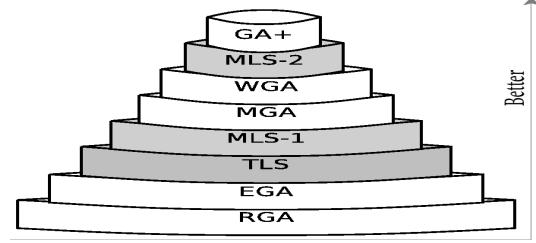


Figure 11. Performance hierarchy within the GA variants and the state-of-the-art: TLS, MLS1 and MLS2 methods. The higher the position, the better the performance.

the proteins having known lower bounds of free energy.

$$RI = \frac{E_t - E_r}{E_l - E_r} * 100\% \quad (4)$$

where E_t and E_r denote the average energies of the target and the reference algorithms and E_l is the *LBFE*.

3) *Performance hierarchy:* The performance hierarchy amongst the GA variants and the state-of-the-art approaches is presented in Figure 11. Experimentally, we show that the three variants of GA—the MGA, WGA, and GAPlus—outperform the local search TLS and MLS1. However, another memory-based local search, MLS2, outperforms the MGA and WGA but is outperformed by GAPlus.

4) *Search progress:* To demonstrate the search progress, we periodically find the best energy values obtained so far in each run for all approaches. For a given period, we then calculate the average energy values obtained for that period over 50 runs. We used 20 minute time intervals to take the

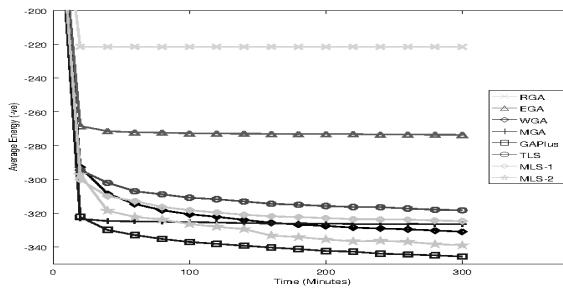


Figure 12. Search progress of different approaches for Protein R1.

average energy values. Figure 12 presents the chart of search progress with progressive time in terms of average energy values obtained by each algorithm for the protein R1.

We observe that after very small initial progress, the EGA and RGA become flat with almost no improvements. The MGA achieves very good progress initially compared to the others, but becomes almost flat later on. The WGA, the TLS, and the MLS1 perform equally initially compared to each other, but later the WGA makes more progress than the TLS and MLS1. The progress of MLS2 is better than that of TLS, MLS1 and other GA variants except GAPlus. Initially, the GAPlus achieves the same progress as the MGA, and later it mostly benefits from the random-walk. The combined positive efforts of the macro-mutation and the random-walk make GAPlus the most efficient conformational search framework. The difference between the performances of WGA and GAPlus remains roughly the same as the initial boosted progress made by the macro-mutation operator.

E. Parameter tuning for GAPlus

In this section, we work with GAPlus parameters. We experimentally try to work out an optimal population size for GAPlus. We also try to find out the optimal rate of using crossover and mutation operators.

1) Tuning GAPlus population size: To observe the effect of population size, we ran GAPlus by setting population sizes 10, 50, 100, 150, 200, 250, 300, 400, and 500. We trace the number of generations, the total number of solutions explored, the total number of stagnations, and the minimum energy level achieved for each population size. We present the average outcomes for the traced values over 50 different runs with further analysis.

a) Population size versus the number of generations: We run GAPlus with different population sizes for proteins F180 and R. We calculate the average number of generations for different population sizes over 50 different runs of identical settings (such as computer, period of running) for each protein. We plot the average number of generations with corresponding population sizes as shown in Figure 13. The experimental results show that the number of generations is proportionately decreasing when increasing the population sizes. Because of the exhaustive generation approach, the total number of conformations explored in each generation are $O(P \times N)$, where P is the population size and N is the number of amino acids in the protein sequence.

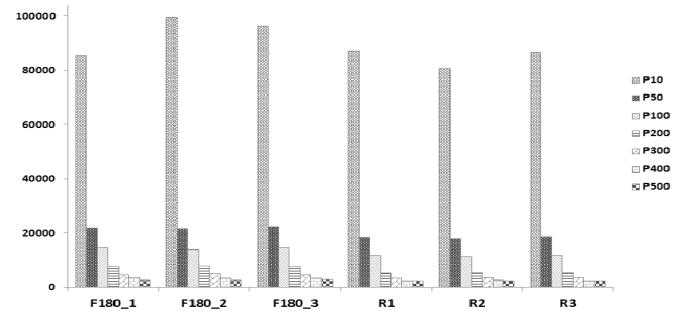


Figure 13. The effect of population size on number of generations.

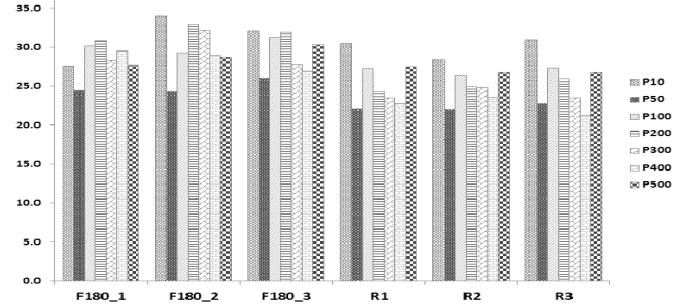


Figure 14. The effect of population size on explored conformations.

b) Population size versus explored conformations: We plot the average number of explored conformations (in millions) with corresponding population sizes as shown in Figure 14. The experimental results show that the number of explored conformations does not vary much with the population sizes. Because the exhaustive generation approach in GAPlus does not remember the previously explored solutions, it may re-explore early explored conformations throughout the period of running. Nevertheless, it is notable that the number of explored conformations is consistently at a minimum for population size 50.

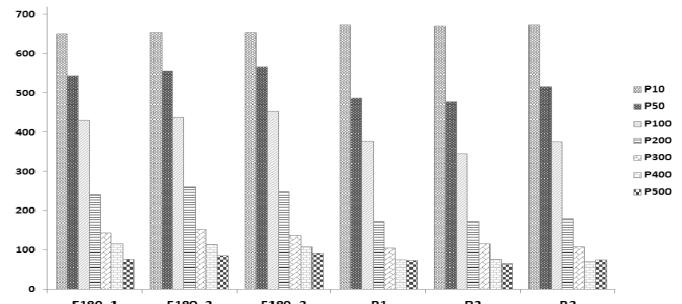


Figure 15. The effect of population size on average number of stagnations.

c) Population size versus the number of stagnations: We plot the average number of stagnations with corresponding population sizes as shown in Figure 15. The experimental results show that the number of occurred stagnations is proportionately decreasing with increasing the population sizes. The result is expected because at stagnation random-walk is applied and the population is diversified. Therefore for larger population size, it takes more time to reach the stagnation

situation by intensifying the population again. The frequency of occurring stagnation decreases while the population size increases.

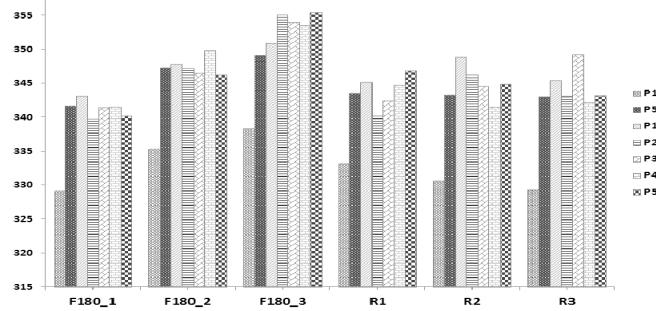


Figure 16. The effect of population size on average energy values.

d) Population sizes versus free energy levels: We plot the average energy values with corresponding population sizes as shown in Figure 16. The experimental results show that the levels of average free energies do not vary much with the population sizes. However, it is notable that for population size 10 (very small), the average energy value for each protein is consistently the minimum amongst the population sizes. Because, for a very small population size, the diversification amongst the individuals is narrow in comparison to other reasonably large population sizes.

In conclusion, observing the effect of population sizes, it is notable that the number of generations and the number of stagnation situations decrease proportionately with increasing population sizes; the impacts on the number of explored conformations with the change of population size is not significant; and the impact of population size on the lower energy level is interesting – for population size 10, the GAPlus performs poorly (due to insufficient diversity) for each protein, however, for other population sizes, the results are very close for each protein.

2) Profiling GAPlus operators: In this section, we experimentally tried to work out an optimal crossover (or mutation) rate for GAPlus. We run GAPlus for crossover rates 0%, 10%, 20%, 30%, 40%, 50%, 60%, 70%, 80%, 90%, and 100% (corresponding mutation rates are 100%, 90%, 80%, 70%, 60%, 50%, 40%, 30%, 20%, 10%, and 0%). We have observed the variations on average energy levels achieved for different crossover or mutation rates in the experimental results.

We further trace some characteristics during operation: in the case of mutation operators, we use all primitive mutation operators such as rotation, diagonal moves, pull moves, and tilt moves (see Figure 6) and present the combined results under the mutation operator. All primitive mutation operators are chosen with equal selection probability.

a) Effects of crossover rates on energy levels: From our experimental results on 11 different crossover rates, we plot a chart of average energy levels with the crossover rates as shown in Figure 17. The chart shows a clear impact of crossover rates on average energy levels. The smaller sizes of crossover rates are more effective than that of the larger sizes. In our experiment, the crossover rate in the range of

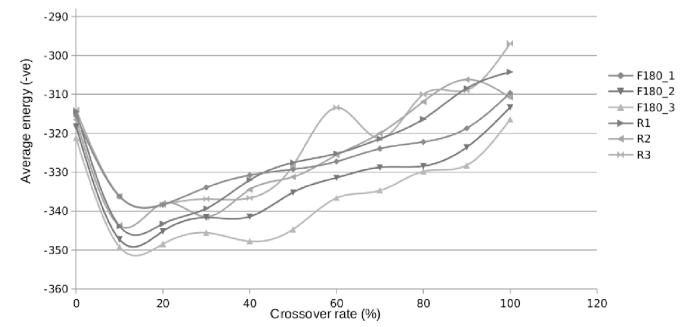


Figure 17. The average energy values corresponding to crossover rates.

10% – 20% (i.e., mutation rates: 80% – 90%) produces better results in GAPlus.

Table VIII
THE EFFECT OF THE APPLICATIONS OF OPERATORS IN THEIR VARIOUS RATES. THE COLUMN XO DENOTES CROSSOVER.

Seq [size]	XO Rate [%]	Crossover Operators			Mutation Operators			Avg energy level
		Attempt [mill]	Succeed [%]	Gain [%]	Usage [%]	Attempt [mill]	Succeed [%]	
0	0	0	0	0	89.67	31.89	1.08	100 -333.4
10	1.5	32.53	30.83	4	84.65	31.15	0.72	96 -348.1
20	4.17	48.34	19.14	16	71.54	30.93	0.72	84 -348.0
30	9.14	54.6	8.66	38	51.1	30.99	0.74	58 -346.3
40	14.18	58.15	5.18	60	36.07	31.24	0.72	40 -347.0
<i>S1</i>	14.86	60.75	4.17	74	21.45	30.94	0.83	26 -342.8
<i>{135}</i>	60	17.33	63.02	3.57	83	14.66	31	0.82
	70	19.3	63.2	3.37	89	9.8	31.02	0.87
	80	18.09	64.8	3.66	94	5.16	30.63	1.02
	90	18.9	67.43	5.13	97	2.37	29.97	1.09
	100	19.54	70.61	10.54	100	0.11	0	0 -330.4
	0	0	0	0	76.77	29.46	1.04	100 -315.5
	10	1.12	26.54	33.78	4	64.05	28.75	0.7
	20	2.7	43.24	30.22	13	57.43	28.68	0.68
	30	6.03	50.61	14.95	34	43.5	28.86	0.71
	40	8.56	54.39	9.61	54	28.21	29.14	0.79
<i>F1</i>	50	11.33	56.4	6.84	70	19.92	28.46	0.75
<i>{180}</i>	60	12.9	58.43	6.05	77	12.87	.30	0.84
	70	12.44	60.99	5.95	87	7.53	29.29	0.95
	80	14.04	62.62	5.93	93	4.51	29.89	0.89
	90	14.77	65.54	7.32	97	2.14	29.95	1.09
	100	12.87	69.6	10.01	100	0.08	0	0 -309.65
	0	0	0	0	54.29	29.13	1.01	100 -314.4
	10	0.85	30.89	42.91	4	52.85	28.79	0.57
	20	2.34	42.86	30.29	13	51.29	28.02	0.56
	30	5.36	49.74	13.98	33	39.06	28.37	0.57
	40	8.27	53.81	8.73	54	26.21	28.6	0.63
<i>R1</i>	50	10.18	57.54	6.99	70	17.52	29.26	0.72
<i>{200}</i>	60	11.13	58.92	5.82	80	11.26	29.26	0.75
	70	12.18	61.86	5.16	88	7.1	29.22	0.8
	80	11.73	66.03	5.47	93	3.88	30.65	0.86
	90	12.68	66.66	6.02	97	1.82	30.54	0.94
	100	12.01	71.57	7.11	100	0.07	0	0 -304.3

b) Effects of crossover rates on other factors: Other GA factors such as operator selection, success of operator application, impact of the applied operator, and total usages of any operator are traced and the summarized results are presented in Table VIII. For each operator, we observed the number of times an operator is selected to apply (Column *Attempt*), how many times the attempts are successful (Column *Succeed*), how many times it improves locally (Column *Gain*), and how much time the operator consumed throughout the life of the run (Column *Usage*). The presented averages are the outcomes over 30 different runs. The experiments are conducted with the population size 100.

Table VIII shows that successful attempts and local energy gain on successful attempts are significantly higher for

crossover in comparison to mutation. However, from the experimental results, it is observed that these local energy gains have little or no impact on the global minimum energy. Another important observation from the experimental results is that crossover consumes almost double operational-time than mutation operators. The time consumption of the crossover operator is reasonable because the crossover operator generates two resultant solutions but mutation generates one.

Observing the experimental results obtained by using different crossover rates (and therefore different mutation rates) in GAPlus, we can conclude that the small crossover rates (i.e., high mutation rates) are better for finding lower energy structures on face-centered-cubic lattice models.

F. Summary of the study

After implementing the proposed new features within GA framework and studying the GA parameters, we have the following finding:

- 1) **Feature enhancement:** The exhaustiveness of individual generation helps improve over random individual generation. The macro-mutation operator and the random-walk technique separately obtain more improvements, with their combination obtaining even further improvements.
- 2) **Parameter tuning:** Using any population size > 10 , the performance of GAPlus does not vary much. The crossover consumes more operational time in comparison to mutation and leads to poorer performance when high crossover rates ($> 20\%$) are used in GAPlus.

VI. DISCUSSION

The proposed GAPlus framework is initially tested on lattice-based protein structure prediction. The problem is chosen because it is a hard combinatorial optimization problem with an astronomically large search space. However, the enhanced features could easily be adopted by other branches of optimization problems such as scheduling, traveling salesman, and vehicle routing. The following discussion would be the future directions of this work:

- 1) The *exhaustive generation* strategy implemented in GAPlus is not a complete or full exhaustive sampling. In contrast, this approach selects successively one amino acid at a time and explores exhaustively in the vicinity of the selected amino acid. It does not sample all possible combinations considering all the amino acids at a time (See Section IV-C2 for details). This strategy might be adopted in any optimization problem that uses some sort of randomness in exploring neighbor solutions.
- 2) The *macro-mutation* operator quickly improves an individual by applying multiple mutations in a single pass. This operator could be applied in other classes of optimization problems to improve a solution either by considering a segment of the solution or by considering the solution as a whole. This operator is guided by a separate heuristic—distance from the hydrophobic core center—other than the objective function (energy function for PSP). Thus, designing a problem specific heuristic is a key to implement this operator.

- 3) The random restarting procedure could replace the *random-walk* technique regardless of the class of optimization problem. Handling stagnation or dealing with local minima is an integral part of an optimization algorithm. Random-restart is an effective strategy in such situations; however, we experimentally showed that the random-walk is more effective for the GAPlus framework.
- 4) By encoding the conformation with angular coordinates (ϕ and ψ), the GAPlus might be applied in high-resolution protein structure prediction. While the minimizing energy function is highly complex (such as molecular dynamics), a simple guidance heuristic (such as hydrophobic property or exposed surface area) could be used to guide the macro-mutation operator. Within GAPlus framework, the macro-mutation operator could be applied optimizing the segments of secondary structures (α -helix and β -sheet).
- 5) More realistic models with complicated energy functions increase the complexity of the problem, however, our approach has a macro-mutation operator that performs as a local optimizer. Based on this, our approach can easily divide the whole optimization process into two stages guided by two energy models with different complexities. The macro-mutation operator can be guided by simpler energy models such as distance from hydrophobic core, exposed surface area, hydrophobicity of amino acids, hydropathy index of the amino acids, and so on. Conversely, the main objective function can be more realistic such as molecular dynamics based energy models. This two-stage optimization will reduce the overall computational complexities. As a result, our framework has a good chance to succeed in more realistic models even for large sized proteins.
- 6) The on-lattice PSP problem can be applied directly to the other classes of optimization problems such as computer organization [72], VLSI design [73], and network topology designing [74] where nodes or processors are connected in a daisy chain but a subset of those need to be topologically compact or closer to have efficient dedicated connections. The PSP application can also be extended to solve quadratic assignment problems using lattice model encoding in genetic algorithms [75], solving job shop scheduling problems in a multi-agent environment [76], and so on.

VII. CONCLUSION

In this paper, we presented five variants of genetic algorithms that individually and in a combined way use three different enhancement techniques: *i*) an exhaustive generation approach; *ii*) a novel hydrophobic-core directed macro-mutation operator; and *iii*) a random-walk based stagnation recovery technique. We found that our final algorithm, GAPlus, which uses a combination of all the enhancement techniques, significantly outperforms all current approaches to simplified protein structure prediction. We also showed that our algorithm is robust enough to produce very similar results under different parameter settings. In the future, we intend to apply GAPlus in

high resolution protein structure prediction. In particular, we want to configure GAPlus within the Rosetta *Codebase* so that it could be used as a separate protocol in place of the Monte Carlo simulation method. We also aim to test GAPlus for other branches of optimization such as scheduling, traveling salesman, and vehicle routing.

ACKNOWLEDGMENT

Mahmood Rashid and Firas Khatib would like to acknowledge the department of Computer and Information Science at the University of Massachusetts Dartmouth. Md Tamjidul Hoque acknowledges the Louisiana Board of Regents through the Board of Regents Support Fund, LEQSF(2013-16)-RD-A-19. Mahmood Rashid and Abdul Sattar would like to express their great appreciation to National ICT Australia and Institute for Integrated and Intelligent Systems at Griffith University, Australia.

SPECIAL THANKS

We would like to express our sincere thanks and appreciation to *Nicholas Trieber* for his excellent proofreading and editing service.

REFERENCES

- [1] H. J. Morowitz, *Energy flow in biology*. Academic Press, 1968.
- [2] A. Stouthamer, "A theoretical study on the amount of ATP required for synthesis of microbial cell material," *Antonie van Leeuwenhoek*, vol. 39, no. 1, pp. 545–565, 1973.
- [3] B. Alberts, A. Johnson, J. Lewis, M. Raff, K. Roberts, and P. Walters, "The shape and structure of proteins," in *Molecular Biology of the Cell, Fourth Edition*. New York and London: Garland Science, 2002.
- [4] N. Horton and M. Lewis, "Calculation of the free energy of association for protein complexes," *Protein science : a publication of the Protein Society*, vol. 1, no. 1, pp. 169–181, 1992.
- [5] O. D. King, A. D. Gitler, and J. Shorter, "The tip of the iceberg: RNA-binding proteins with prion-like domains in neurodegenerative disease," *Brain Research*, vol. 1462, pp. 61–80, 2012.
- [6] Adam Smith, "Protein misfolding," *Nature Reviews Drug Discovery*, vol. 426, no. 6968, pp. 78–102, December 2003.
- [7] C. M. Dobson, "Protein folding and misfolding," *Nature*, vol. 426, no. 6968, pp. 884–890, 2003.
- [8] F. Chiti and C. M. Dobson, "Protein Misfolding, Functional Amyloid, and Human Disease," *Annu Rev Biochem*, vol. 75(1), pp. 333–366, 2006.
- [9] Jucker Mathias and Walker Lary C., "Self-propagation of pathogenic protein aggregates in neurodegenerative diseases," *Nature*, vol. 501, no. 7465, pp. 45–51, 2013.
- [10] P. Veerapandian, Ed., *Structure-based drug design*. New York: Marcel Dekker, 1997.
- [11] A. Breda, N. F. Valadares, O. N. de Souza, and R. C. Garratt, "Protein Structure, Modelling and Applications," 2008.
- [12] P. Kast and D. Hilvert, "3D structural information as a guide to protein engineering using genetic selection," *Current Opinion in Structural Biology*, vol. 7, no. 4, p. 470479, 1997.
- [13] S. M. Berry and Y. Lu, *Protein Structure Design and Engineering*. John Wiley & Sons, Ltd, 2001.
- [14] H. Lilie, "Designer proteins in biotechnology," *EMBO Reports*, vol. 4, no. 4, pp. 346–351, 2006.
- [15] R. Jaenicke and R. Sterner, "Protein Design at the Crossroads of Biotechnology, Chemistry, Theory, and Evolution," *Angewandte Chemie International Edition*, vol. 42, no. 2, pp. 140–142, 2003.
- [16] A. Yonath, "X-ray crystallography at the heart of life science," *Current Opinion in Structural Biology*, vol. 21, no. 5, pp. 622–626, 2011.
- [17] The Science Editorial, "So much more to know," *The Science*, vol. 309, no. 5731, pp. 78–102, July 2005.
- [18] A. Leaver-Fay and et al., "Rosetta3: An Object-Oriented Software Suite for the Simulation and Design of Macromolecules," ser. Method enzymol. Academic Press, 2011, vol. 487, pp. 545–574.
- [19] Kemenge Kyle E, Hickey John M, Lovell Scott, Battaile Kevin P, Zhang Yang, and Hefty P Scott, "Ab Initio Structural Modeling of and Experimental Validation for Chlamydia trachomatis Protein CT296 Reveal Structural Similarity to Fe(II) 2-Oxoglutarate-Dependent Enzymes," *J. of Bacteriology*, vol. 193, no. 23, pp. 6517–6528, 2011.
- [20] Xu Dong and Zhang Yang, "Ab Initio Protein Structure Assembly Using Continuous Structure Fragments and Optimized Knowledge-based Force Field," *Proteins*, vol. 80, no. 7, pp. 1715–1735, apr 2012.
- [21] K. Molloy, S. Saleh, and A. Shehu, "Probabilistic search and energy guidance for biased decoy sampling in *Ab Initio* protein structure prediction," *IEEE/ACM Trans. Comput. Biology Bioinform.*, vol. 10, no. 5, pp. 1162–1175, 2013.
- [22] B. Olson and A. Shehu, "An evolutionary-inspired algorithm to guide stochastic search for near-native protein conformations with multiobjective analysis," in *AAAIW*, 2013.
- [23] P. Bradley, K. M. S. Misura, and D. Baker, "Toward high-resolution *de novo* structure prediction for small proteins," *Science*, vol. 309, no. 5742, p. 1868, 2005.
- [24] D. E. Kim, B. Blum, P. Bradley, and D. Baker, "Sampling bottlenecks in *De novo* protein structure prediction," *Journal of Molecular Biology*, vol. 393, no. 1, pp. 249 – 260, 2009.
- [25] J. Zhang, "Protein-length distributions for the three domains of life," *Trends in Genetics*, vol. 16, no. 3, pp. 107 – 109, 2000.
- [26] L. Brochieri and S. Karlin, "Protein length in eukaryotic and prokaryotic proteomes," *Nucleic Acids Res.*, vol. 33(10), pp. 3390–3400, 2005.
- [27] T. Higgs, B. Stantic, M. T. Hoque, and A. Sattar, "Applying feature-based resampling to protein structure prediction," in *BiCoB*. Curran Associates, Inc, 2012, pp. 239–244.
- [28] C. Levinthal, "Are there pathways for protein folding?" *Journal of Medical Physics*, vol. 65, no. 1, pp. 44–45, 1968.
- [29] C. B. Anfinsen, "The principles that govern the folding of protein chains," *Science*, vol. 181, no. 4096, pp. 223–230, 1973.
- [30] M. Cebríán, I. Dotú, P. Van Hentenryck, and P. Clote, "Protein structure prediction on the face centered cubic lattice by local search," in *Proceedings of the 23rd national conference on Artificial intelligence - Volume 1*. AAAI Press, 2008, pp. 241–246.
- [31] I. Dotú, M. Cebríán, P. Van Hentenryck, and P. Clote, "On lattice protein structure prediction revisited," *IEEE Trans. on Comp. Bio and Bioinfo.*, 2011.
- [32] S. Shatabda, M. A. H. Newton, D. N. Pham, and A. Sattar, "Memory-based local search for simplified protein structure prediction," in *The ACM Conference on Bioinformatics, Computational Biology and Biomedicine (ACM-BCB)*. Orlando, FL, USA: ACM, 2012.
- [33] S. Shatabda, M. A. H. Newton, M. A. Rashid, D. Pham, and A. Sattar, "The road not taken: retreat and diverge in local search for simplified protein structure prediction," *BMC Bioinformatics*, vol. 14, 2013.
- [34] M. T. Hoque, M. Chetty, and A. Sattar, "Protein folding prediction in 3D FCC HP lattice model using genetic algorithm," vol. 2007. IEEE Congress on Evolutionary Computation, 2007, pp. 4138–4145.
- [35] H.-J. Böckenhauer, D. Ullah, L. Kapsokalivas, and K. Steinmögel, "A local move set for protein folding in triangular lattice models," in *WABI*, vol. 5251. Springer, 2008, pp. 369–381.
- [36] T. Schwede, J. Kopp, N. Guex, and M. C. Peitsch, "Swiss-model: an automated protein homology-modeling server," *Nucleic Acids Res.*, vol. 31, no. 13, pp. 3381–3385, 2003.
- [37] W. T. Des Higgins, *Bioinformatics: Sequence, Structure and Databanks*. Oxford University Press, 2001, 1st Ed.
- [38] J. Skolnick, D. Kihara, and Y. Zhang, "Development and large scale benchmark testing of the PROSPECTOR_3 threading algorithm," *Proteins*, vol. 56, no. 3, pp. 502–518, 2004.
- [39] H. Zhou and Y. Zhou, "Fold recognition by combining sequence profiles derived from evolution and from depth-dependent structural alignment of fragments," *Proteins*, vol. 58, no. 2, pp. 321–328, 2005.
- [40] S. Wu and Y. Zhang, "MUSTER: Improving protein sequence profile-profile alignments by using multiple sources of structure information," *Proteins*, vol. 72, no. 2, pp. 547–556, 2008.
- [41] L. A. Kelley and M. J. E. Sternberg, "Protein structure prediction on the web: a case study using the PHYRE Server," *NATURE PROTOCOLS*, vol. 4, no. 3, pp. 363–371, 2009.
- [42] S. Wu, J. Skolnick, and Y. Zhang, "Ab initio modeling of small proteins by iterative TASSER simulations," *BMC Biology*, vol. 5, no. 1, 2007.
- [43] T. C. Hales, "A proof of the Kepler conjecture," *The Annals of Mathematics*, vol. 162, no. 3, pp. 1065–1185, 2005.
- [44] K. F. Lau and K. A. Dill, "A lattice statistical mechanics model of the conformational and sequence spaces of proteins," *Macromolecules*, vol. 22, no. 10, pp. 3986–3997, 1989.
- [45] C. Thachuk, A. Shmygelska, and H. H. Hoos, "A replica exchange Monte Carlo algorithm for protein folding in the HP model," *BMC bioinformatics*, vol. 8, no. 1, p. 342, 2007.
- [46] A.-A. Tantar, N. Melab, and E.-G. Talbi, "A grid-based genetic algorithm combined with an adaptive simulated annealing for protein structure

- prediction," *Soft Computing-A Fusion of Foundations, Methodologies and Applications*, vol. 12, no. 12, pp. 1185–1198, 2008.
- [47] R. Unger and J. Moult, "A genetic algorithm for 3D protein folding simulations," Morgan Kaufmann Publishers. The 5th International Conference on Genetic Algorithms, 1993, p. 581.
- [48] M. T. Hoque, "Genetic Algorithm for *ab initio* protein structure prediction based on low resolution models," Ph.D. dissertation, Monash University, Victoria, Australia, Sep. 2007.
- [49] G. W. Klau, N. Lesh, J. Marks, and M. Mitzenmacher, "Human-guided tabu search," in *The Eighteenth National Conference on Artificial Intelligence (AAAI-02)*. Edmonton, AB, Canada: The Eighteenth National Conference on Artificial Intelligence (AAAI-02), 2002.
- [50] C. Blum, "Ant colony optimization: Introduction and recent trends," *Physics of Life reviews*, vol. 2, no. 4, pp. 353–373, 2005.
- [51] I. Kondov and R. Berlich, "Protein structure prediction using particle swarm optimization and a distributed parallel approach," in *ACM workshop on Biologically inspired algorithms for distributed systems*, ser. BADS '11, 2011, pp. 35–42.
- [52] N. Mansour, F. Kanj, and H. Khachfe, "Particle swarm optimization approach for protein structure prediction in the 3D HP model." *Interdiscip Sci*, vol. 4, no. 3, pp. 190–200, Sep 2012.
- [53] V. Cutello, G. Nicosia, M. Pavone, and J. Timmis, "An immune algorithm for protein structure prediction on lattice models," *IEEE Trans. on Evolutionary Computation*, vol. 11, no. 1, pp. 101–117, 2007.
- [54] M. Mann, S. Will, and R. Backofen, "CPSp-tools – Exact and complete algorithms for high-throughput 3D lattice protein studies," *Bmc Bioinformatics*, vol. 9, no. 1, p. 230, 2008.
- [55] M. T. Hoque, M. Chetty, A. Lewis, and A. Sattar, "Twin removal in genetic algorithms for protein structure prediction using low-resolution model," *IEEE Trans. on Comp. Bio and Bioinfo.*, vol. 8, no. 1, pp. 234–245, 2011.
- [56] L. Kapsokalivas, X. Gan, A. A. Albrecht, and K. Steinhöfel, "Population-based local search for protein folding simulation in the MJ energy model and cubic lattices," *Comput. Biol. Chem.*, vol. 33, no. 4, pp. 283–294, Aug. 2009.
- [57] D. Ullah and K. Steinhöfel, "A hybrid approach to protein folding problem integrating constraint programming with local search," *BMC bioinformatics*, vol. 11, no. Suppl 1, p. S39, 2010.
- [58] A. Dal Palù, A. Dovier, F. Fogolari, and E. Pontelli, "Exploring protein fragment assembly using CLP." AAAI Press, 2011, pp. 2590–2595.
- [59] S. Shatabda, M. A. H. Newton, and A. Sattar, "Mixed heuristic local search for protein structure prediction," in *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence*. AAAI Press, 2013.
- [60] S. Miyazawa and R. L. Jernigan, "Estimation of effective interresidue contact energies from protein crystal structures: quasi-chemical approximation," *Macromolecules*, vol. 18, no. 3, pp. 534–552, 1985.
- [61] M. Berrera, H. Molinari, and F. Fogolari, "Amino acid empirical contact energy definitions for fold recognition in the space of contact maps," *BMC Bioinformatics*, vol. 4, no. 1, p. 8, 2003.
- [62] M. K. Islam and M. Chetty, "Clustered memetic algorithm with local heuristics for ab initio protein structure prediction," *IEEE Trans. on Evolutionary Computation*, vol. 17, no. 4, pp. 558–576, 2013.
- [63] K. Yue and K. A. Dill, "Sequence-structure relationships in proteins and copolymers," *Physical Review E*, vol. 48, no. 3, p. 2267, 1993.
- [64] N. Lesh, M. Mitzenmacher, and S. Whitesides, "A complete and effective move set for simplified protein folding," in *Research in Computational Molecular Biology*. ACM, 2003, pp. 188–195.
- [65] M. Rocha and J. Neves, "Preventing premature convergence to local optima in genetic algorithms via random offspring generation," in *Int. Conf. on Ind. and Eng. Applications of Artificial Intelligence and Expert Systems*. Springer-Verlag New York, Inc., 1999, pp. 127–136.
- [66] T. F. Gonzalez, *Handbook of Approximation Algorithms and Metaheuristics*. Chapman & Hall/CRC, 2007, ch 14 : page 12.
- [67] K. E. Parsopoulos and M. N. Vrahatis, *Particle Swarm Optimization and Intelligence: Advances and Applications*. Information Science Reference - Imprint of: IGI Publishing, 2010, page 83.
- [68] M. Dorigo and T. Stützle, *Ant Colony Optimization*. Scituate, MA, USA: Bradford Company, 2004, ch 3 : page 86.
- [69] NICTA, "National ICT Australia research centre of excellence," [www.nicta.com.au].
- [70] P. Clote, "Peter Clote Laboratory website," last accessed on November 23, 2013, [bioinformatics.bc.edu/clotelab/FCCproteinStructure/].
- [71] CASP, "The critical assessment of protein structure prediction," [http://predictioncenter.org].
- [72] H. J. Patterson DA, *Computer Organization and Design, The Hardware/Software Interface*. Elsevier, 2013, 5th ed 800 p.
- [73] S. D. Kugelmass, R. Squier, and K. Steiglitz, "Performance of vlsi engines for lattice computations," *Complex Systems*, 1987.
- [74] J. Degila and B. Sanso, "Topological design optimization of a yottabit-per-second lattice network," *Selected Areas in Communications, IEEE Journal on*, vol. 22, no. 9, pp. 1613–1625, 2004.
- [75] R. C. Mitsuo Gen, *Genetic Algorithms and Engineering Design*. Wiley Online Library, 2007.
- [76] X. Duan, J. Liu, L. Zhang, and L. Jiao, "Multi-objective job shop scheduling based on multiagent evolutionary algorithm." Springer Berlin Heidelberg, 2010, vol. 6457, pp. 543–552.



Mahmood A. Rashid is a Fellow in the School of Computing, Information & Mathematical Sciences at the University of the South Pacific, Suva, Fiji. He was a Postdoctoral Fellow in the Department of Computer and Information Science at the University of Massachusetts Dartmouth, USA. He received his Ph.D. in Computer Science from Griffith University, Australia in 2014 and since then he has been appointed as an Adjunct Research Fellow in the Institute for Integrated and Intelligent Systems at Griffith University. He was a Graduate Researcher

at NICTA during his Ph.D. candidature. He also earned his MSc degree in Computer Science and BSc degree in Mechanical Engineering. His research focuses on the optimization algorithms; the parallel and distributed computing; the search guidance heuristics; and the computational biology. He spent several years (2000–2009) in information technology industries as an IT specialist.



Firas Khatib is an Assistant Professor in the Computer and Information Science Department at the University of Massachusetts Dartmouth. He received his bachelor's degree in Applied Mathematics from UC Berkeley in 2001, and his Ph.D. in Bioinformatics from UC Santa Cruz in 2008. Prior to joining UMass Dartmouth, Firas was a Senior Fellow in the Molecular Engineering & Sciences Institute at the University of Washington. Firas was awarded a Postdoctoral Research Fellowship in Biology from the National Science Foundation to work on the

protein-folding game Foldit. The long term goal of this project is to utilize the combined power of humans and computers in order to build accurate models of disease-related proteins.



Md Tamjidul Hoque is an Assistant Professor of Computer Science Department at the University of New Orleans. He was a postdoctoral fellow (Feb/2011-July/2012) at IUPUI, Indiana and was a research fellow (Jan/2007 to Jan/2011) at Griffith University, Australia. He obtained his Ph.D. in IT from Monash University, Australia in 2008. His research focus predominantly in the area of machine learning, evolutionary computation and artificial intelligence, applying towards hard optimization problems especially for biology and bioinformatics such

as protein-structure-prediction, disorder-predictor, and energy function. He also developed efficient algorithms for high-throughput screening and high content image analysis. Hoque was a deputy general manager and head of IT in a large industrial-group (Dec/1999 to Mar/2004), where he developed a number of commercial software for the industries.



Abdul Sattar is the founding director of the IIIS and a professor of CS and AI at Griffith University, Australia. He is also one of the research leaders at the NICTA Queensland Research Lab (QRL). He has been an academic staff member at Griffith University since February 1992 as a Lecturer (1992–1995), Senior Lecturer (1996–1999), and Professor (2000– present) within the School of Information and Communication Technology. Prior to join at Griffith University, he was a Lecturer of Physics in Rajasthan, India (1980–1982), and a Research

Scholar at Jawaharlal Nehru University, India (1982–1985), the University of Waterloo, Canada (1985–1987), and the University of Alberta, Canada (1987–1991). Abdul earned his PhD in Computer Science from the University of Alberta, Canada, in 1990 majoring Artificial Intelligence.