

Heuristics, Optimizations, and Parallelism for Protein Structure Prediction in $\text{CLP}(\mathcal{FD})$ *

Alessandro Dal Palù¹, Agostino Dovier¹, and Enrico Pontelli²

¹ Dip. di Matematica e Informatica, Univ. di Udine
(dalpalu|dovier)@dimi.uniud.it

² Department of Computer Science, New Mexico State University.
epontell@cs.nmsu.edu

Abstract. The paper describes a constraint-based solution to the protein folding problem on *face-centered cubic lattices*—a biologically meaningful approximation of the general protein folding problem. The paper improves the results presented in [14] and introduces new ideas for improving efficiency: (i) proper reorganization of the constraint structure; (ii) development of novel, both general and problem-specific, heuristics; (iii) exploitation of parallelism. Globally, we obtain a speed up in the order of 60 w.r.t. [14]. We show how these results can be employed to solve the folding problem for large proteins containing subsequences whose conformation is already known.

1 Introduction

Proteins are responsible for nearly every function required for life. The sequence of elements (amino acids) identifying a protein is known as the primary (1D) structure. A functional protein can be thought of as a properly folded chain of amino acids in 3-dimensional (3D) space. The 3D structure of a protein characterizes its function. A folded protein interacts three-dimensionally with other proteins (e.g., lock and key arrangements) and this interaction determines the functions of the organism. In fact, an organism is essentially determined by the three-dimensional interactions between proteins and substrates. Thus, without knowing the 3D structure of the proteins coded in a genome, we cannot completely understand the phenotype and functioning of living organisms. Understanding how protein folds has profound implications—e.g., towards the theoretical design of exact drugs, the improvement of proteins functionality, and the precise modeling of cells.

In recent decades, most scientists have agreed that the answer to the folding problem lies in the concept of the *energy state* of a protein. The predominant strategy in solving the protein folding problem has been to determine a state of the amino acid sequence in the 3D space with minimum energy. According to this theory, the 3D conformation that yields the lowest energy state represents the protein's natural shape (a.k.a. the *native conformation*). The energy of a conformation can be modeled using *energy functions*, that determine the energy level based on the interactions between any pairs of amino acids [7]. Thus, we can reduce the protein folding problem to an optimization problem, where the energy function has to be minimized under a collection of constraints (e.g., derived from known chemical and physical properties) [12].

We employ *Constraint Logic Programming (CLP)*, in particular, constraint logic programming over *finite domains* ($\text{CLP}(\mathcal{FD})$), to model and solve a tractable representation of the protein folding problem—i.e., protein folding in the context of face-centered cubic

* This paper has been submitted and accepted in PPDP 2005.

lattices [28]. The choice of CLP is natural for a variety of reasons (see, e.g., [2]). CLP is a paradigm that is highly suitable to address optimization problems; it provides declarative and high-level modeling, combined with effective built-in search and resolution strategies. Furthermore, the high-level modeling offered by CLP allows us to easily *add* new constraints and to plug-and-play different search strategies and heuristics.

A preliminary approach to the use of CLP(\mathcal{FD}) for the protein folding problem has been presented in [14]. In this paper, we elaborate such proposal to consider the issue of *efficiency* and *scalability*. The ultimate objective of this paper is to demonstrate that

- the modeling and optimization capabilities of CLP(\mathcal{FD}) are highly suitable to tackle the protein folding problem;
- the high-level modeling capabilities allow us to easily add or modify constraints as they become available, and to explore the use of different heuristics and search strategies;
- efficiency and scalability can be achieved for realistic problems.

The first step in this process lies in a remodeling of the constraint problem, aim at making the modelization more suitable to the capabilities of current CLP(\mathcal{FD}) solvers. We also introduce new high-level heuristics for guiding the exploration of the search space, leading to a more effective pruning and enhanced scalability. In particular, we introduce a heuristic called *Bounded Block Fails*. Where the built-in strategies are insufficient to achieve acceptable levels of performance, we introduce the use of *parallelism*, easily exploitable from the high-level search structure generated by the CLP execution. Using a novel structure, heuristics, and a 14 processors parallel machine, we obtain a speed-up in the order of 60 w.r.t. the performance of the code presented in [14]—running on a single processor of the same parallel machine. The investigation proposed here pushes the built-in capabilities of CLP solvers to, what we believe, are the limits for the problem at hand. The proposed results indicate also that *better performance could be accomplished*, but only at the price of building some of the proposed heuristics at a lower level—i.e., as an ad-hoc constraint solver, and bypassing the built-in strategies of CLP(\mathcal{FD}). Finally, we demonstrate what, we believe, is one of the greatest application areas for our technology: folding large proteins containing subsequences whose native conformation is already known (e.g., by homology)—e.g., macro blocks linked by a neutral coil. This type of situation is very common; in our framework, the known structures can be directly added as constraints (with *no modifications* to the rest of the constraint model), allowing us to tackle large problems with excellent performance.

1.1 Related Works

The bibliography on the protein folding problem is extensive [8, 25]; the problem has been recognized as a fundamental challenge [22], and it has been addressed with a variety of approaches (e.g., comparative modeling through homology, fold recognition through threading, ab initio fold prediction).

An abstraction of the problem, that has been recently investigated, is the protein folding problem in the *HP* model, where amino acids are separated into two classes (*H*, hydrophobic, and *P*, hydrophilic). The goal is to search for a conformation produced by an *HP* sequence, such that most HH pairs are neighboring in a predefined lattice. The problem has been studied on 2D square lattices [13, 21], 2D triangular lattices [1], 3D square models [19], and face-centered cubic lattices (fcc) [23]. Backofen and Will have extensively studied this last problem [3–5]. The approach is suited for globular proteins, since the main force driving the folding process is the electrical potential generated by *H*s and *P*s, and the fcc lattices are one of the best and simplest approximation of the 3D space (Sect. 2.2). Compared to the work of Backofen and Will, our approach refines the energy contribution model, extending the interactions between classes *H* and *P* to interactions between each pair of amino

acids [7]. Moreover, we introduce the possibility to model secondary structure elements, that cannot be reproduced correctly using only a simple energy model as the one adopted by other researchers.

The use of constraint programming technology in the context of the protein folding problem has been fairly limited. Backofen and Will have made use of constraints over finite domains in the context of the *HP* problem [5]. Clark et al. employed Prolog to implement heuristics in pruning a exhaustive search for predicting α -helix and β -sheet topology from secondary structure and topological folding rules [11]. Distributed search and continuous optimization have been used in ab initio structure prediction, based on selection of discrete torsion angles for combinatorial search of the space of possible protein foldings [18].

2 Problem modeling

2.1 The Protein Folding Problem

The *Primary* structure of a protein is a sequence of linked units (*amino acids* or *residues*) of a given length. The amino acids can be identified by an alphabet \mathcal{A} of 20 different symbols, associated to specific chemical-physical properties. A protein has a high degree of freedom, and its 3D conformation is named *Tertiary* structure.

From the *energy* point of view, the molecule tends to reach a conformation with a minimal value of free energy (*Native* conformation). Native conformations are largely built from *Secondary Structure elements* (i.e., helices and sheets) often arranged in well-defined motifs. α -helices are constituted by residues arranged in a regular right-handed helix with 3.6 residues per turn. β -sheets are constituted by extended strands. Each strand is made of contiguous residues, but strands participating in the same sheet are not necessarily contiguous in sequence. Algorithms, e.g., based on neural networks, that have been developed, are capable to predict the secondary structure with high accuracy (75% [8]).

Another important structural feature of proteins is the capability of cysteine residues of covalently bind through their sulphur atoms, thus forming disulfide bridges, which impose important constraints on the structure (also known as *ssbonds*). This kind of information is often available, either through experiments or predictions.

Several models have been proposed for reasoning about the 3D properties of proteins—i.e., dealing with the *Tertiary Structure*. Given a primary sequence $S = s_1 \cdots s_n$, with $s_i \in \mathcal{A}$, let us represent with $\omega(i)$ the *position* of a point representing the amino acid s_i in space; $\omega(i)$ is a vector $\langle x_i, y_i, z_i \rangle \in \mathcal{D}$, for a given space domain \mathcal{D} . The values x_i, y_i , and z_i can be real numbers—in models in which proteins are free of taking any positions in space—or integer numbers—in models where amino acids can assume only a finite number of positions within a suitable lattice structure. We call \mathcal{D} the set of admissible points.

Given two points $\omega_1, \omega_2 \in \mathcal{D}$, we indicate with $\text{next}(\omega_1, \omega_2)$ the fact that the two points are admissible positions for two amino acids that are contiguous in the primary sequence. We assumed that consecutive amino acids are separated by a fixed distance.

We also employ the binary predicate `contact`, which is used to describe the fact that two amino acids are sufficiently close to be able to interact, and thus they contribute to the energy function: two non-consecutive amino acids s_i and s_j in the positions $\omega(i)$ and $\omega(j)$ are in contact (denoted by `contact`($\omega(i), \omega(j)$)) when their distance is less than a given threshold.

Given a primary sequence $S = s_1 \cdots s_n$, with $s_i \in \mathcal{A}$, a *folding* of S is a function $\omega : \{1, \dots, n\} \rightarrow \mathcal{D}$ such that:

1. `next`($\omega(i), \omega(i+1)$) for $i = 1, \dots, n-1$, and
2. $\omega(i) \neq \omega(j)$ for $i \neq j$ (namely, ω introduces no loops).

A simplified evaluation of the energy of a folding can be obtained by observing the *contacts* present in the folding. In particular, every time a contact between a pair of amino

acids is detected, a specific energy contribution is applied towards the global energy. These contributions can be obtained from tables developed using statistical methods applied to structures obtained from X-Rays and Nuclear Magnetic Resonance experiments [20, 7]; these tables associates an energy measure to each pair of non-consecutive amino acids when they are in contact. We denote with $\text{Pot}(s_i, s_j)$ the energy contribution associated to the amino acids s_i and s_j (the order does not matter).

The *protein structure prediction problem* can be modeled as the problem of finding the folding ω of S such that the following energy cost function is minimized:

$$E(\omega, S) = \sum_{1 \leq i < n} \sum_{i+2 \leq j \leq n} \text{contact}(\omega(i), \omega(j)) \cdot \text{Pot}(s_i, s_j).$$

With a slight abuse of notation predicate contact is here used as a Boolean function. This definition is sufficiently general to cover the case of several spatial models \mathcal{D} , such as the fcc lattice and the cubic lattice [12].

2.2 Lattice Models

Lattice models have long been used for protein structure prediction (see [26] for a survey). In [23] it is shown that the *Face-Centered Cubic Lattice* (fcc) model is a well-suited, realistic model for 3D conformations of proteins. The model is based on cubes of size 2, where the central point of each face is also admitted. The domain \mathcal{D} consists of the set of triples $\langle x, y, z \rangle$, where $x, y, z \in \mathbb{N}$ such that $x + y + z$ is even (see Fig. 1). Points at Euclidean distance $\sqrt{2}$ are linked; their distance is called *lattice unit*. Observe that, for linked points i and j , it holds that $|x_i - x_j| + |y_i - y_j| + |z_i - z_j| = 2$.

Each point is adjacent to 12 neighboring points. Thus, we define the predicate next as follows: $\text{next}(\omega(i), \omega(i+1))$ holds iff

- $|x_i - x_{i+1}| \in \{0, 1\}, |y_i - y_{i+1}| \in \{0, 1\}, |z_i - z_{i+1}| \in \{0, 1\}$,
- $|x_i - x_{i+1}| + |y_i - y_{i+1}| + |z_i - z_{i+1}| = 2$.

In fcc lattices, the angle between three adjacent residues may assume values 60° , 90° , 120° , and 180° . Volumetric constraints and energetic restraints in proteins make values 60° and 180° infeasible. Therefore, in our model, we retain only the 90° and 120° angles [28, 17]. No similar restriction exists on torsional angles among four adjacent residues. In detail, let $\mathbf{v}_{i-1,i} = \omega(i) - \omega(i-1)$ and $\mathbf{v}_{i,i+1} = \omega(i+1) - \omega(i)$. To impose that the angle between them can only be of 90° and 120° , we use the scalar product between these two vectors: $\mathbf{v}_{i-1,i} \cdot \mathbf{v}_{i,i+1} = |\mathbf{v}_{i-1,i}| |\mathbf{v}_{i,i+1}| \cos(\theta)$. Thus, since $|\mathbf{v}_{i-1,i}| = |\mathbf{v}_{i,i+1}| = \sqrt{2}$ we only need to impose that: $\mathbf{v}_{i-1,i} \cdot \mathbf{v}_{i,i+1} \in \{1, 0\}$.

A *contact* between two non-adjacent residues in fcc occurs when their separation is two lattice units—i.e., viewing the lattice as a graph whose edges connect adjacent points in the lattice, the positions of the residues are connected by a path of length 2.

Physically, two amino acids in contact cannot be at the distance of a single lattice unit, because their volumes would overlap. Consequently, we impose the constraint that two non-consecutive residues s_i and s_j must be separated by more than one lattice unit, by adding the constraint (called *non_next*): $(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2 > 2$. With these additional constraints, we can define: $\text{contact}(\omega(i), \omega(j))$ iff $|x_i - x_j| + |y_i - y_j| + |z_i - z_j| = 2$.

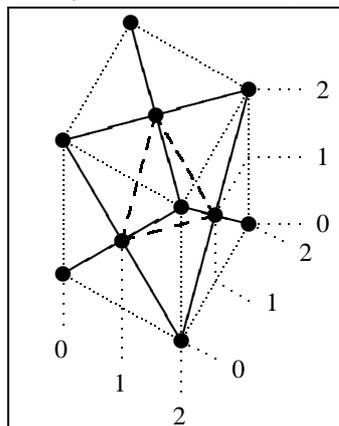


Fig. 1: A cube of the fcc lattice. Thick lines link connected points.

3 A CLP(\mathcal{FD}) Implementation

The formalization of the protein structure prediction problem in `fcc` has been instantiated in a declarative program in CLP over finite domains—and it follows the basic structure outlined in [15, 14]. In this section, we describe the main predicates employed in such implementation. We wrote the system in CLP(\mathcal{FD}) of SICStus PROLOG 3.12.0 [9] and the library `ic` of ECLiPSe 5.8 [10]. Complete code and other related material can be found in: www.dimi.uniud.it/dovier/PF.

The program has the classical *Constrain & Generate* [2] structure. The `constrain` predicate deterministically adds the constraints for the variables involved, while `labeling` searches for the solution in the search space (through chronological backtracking). The inputs are `Primary` (a list of amino acids), `Secondary` (a list of known secondary structure components), and `Matrix` (the matrix of energy contributions, see later in the Section). `Tertiary` is the output list of positions (a flat list of triples of integers) of the conformation and `Energy` is the output value of energy associated to such conformation. We do not discuss here the more technical parameters. The predicate `constrain` is defined here:

```
constrain(Primary, Secondary, Energy, Tertiary, Matrix,
         Compact, OrdTertiary):-
    length(Primary,N),
    generate_tertiary(N, Tertiary), domain_bounds(Tertiary,N),
    avoid_self_loops(Tertiary,N), next_constraints(Tertiary),
    distance_constraints(Tertiary), angles(Tertiary),
    compact_constraints(Tertiary,N,Compact),
    secondary_info(Secondary, Tertiary),
    avoid_symmetries(Secondary, Tertiary, SSP),
    define_variables_order(Secondary,Tertiary, SSP, OrdTertiary),
    energy_constraints(Primary, Tertiary, Secondary, Energy, Matrix).
```

Given the input list of amino acids `Primary` = $[s_1, \dots, s_N]$, `generate_tertiary` creates the list `Tertiary` = $[X_1, Y_1, Z_1, \dots, X_N, Y_N, Z_N]$ of $3N$ variables. The predicate `domain_bounds` specifies the domains for the X_i, Y_i, Z_i variables (the range $0 \dots 2 * N$), and adds the constraints forcing $X_i + Y_i + Z_i$ to be even. `avoid_self_loops` forces all triples to be distinct. `next_constraints` imposes that the points $[X_i, Y_i, Z_i]$ and $[X_{i+1}, Y_{i+1}, Z_{i+1}]$ are adjacent in the lattice. `distance_constraints` forces two non-consecutive points to be at a lattice distance greater than 1. `angles` defines the admissible angles formed by three consecutive amino acids. The predicate `compact_constraints` introduces a user-defined maximal distance between amino acids (called *compact factor*). `secondary_info` encodes the Secondary Structure information as constraints in the program. The secondary structure (`Secondary`) information is retrieved from the Protein Data Bank [6]; it is easy to modify the code to obtain such information from secondary structure prediction programs (e.g., [24]). `avoid_symmetries` removes equivalent admissible conformations modulo symmetries and/or rotations. `define_variables_order` makes use of the distribution of secondary structure elements to sort the variables in `Tertiary` for labeling purposes (see Section 4.1).

As described above, two amino acids s_i, s_j are in contact if a path of length 2 lattice units links them. The contact information is maintained in a symmetric matrix M , such that $M[i, j]$ is the energy contribution provided by s_i and s_j . The following code constrains the contact contribution to the matrix element:

```
table(si,sj,Pot), M[i,j] in {0,Pot},
2 #= abs(Xi-Xj) + abs(Yi-Yj) + abs(Zi-Zj) #<=> M[i,j] #= Pot.
```

where `table` reports the value $\text{Pot}(s_i, s_j)$ as computed in [7]. Values of `Pot` have been scaled w.r.t. [7] to have only integer values. The global energy is the sum of the elements in M . The optimal folding is reached when the global energy is minimal. During the *labeling* phase, the information stored in M is used to control the minimization process and to prune the search tree.

Various heuristics have been developed, aimed at simplifying the computation of the contact matrix M and at the search pruning based on the dynamic analysis of M .

To reduce the search space, it is possible to ignore the contact contributions for all pairs s_i, s_j such that $j < i + \text{min_aa_dist}$, where `min_aa_dist` is a parameter. In those cases, the constraint is not applied and $M[i, j] = 0$. When we increase the value `min_aa_dist`, we generate a simpler global constraint and, at the same time, we consider only contributions from distant amino acids (considered, from the biological perspective, more relevant). For more details about other heuristics we refer to [14, 15].

4 The new implementation

4.1 Constraints Redesign

In this section we present the improvements w.r.t. the code in [14], briefly described in the previous section, to seek better performance and scalability to larger proteins.

In our previous experiments, we observed that one of the causes behind the lack of scalability of the constraint system of Sect. 3 to large proteins is the limited propagation achieved in presence of non-linear constraints (e.g., constraints related to torsion angles, as they require vector scalar products). The first stage of redesign implied removing as many non-linear constraints as possible and enhancing the propagation structure.

Angle-free Encoding The constraint structure proposed in [14] to solve the problem, imposed an alternative local coordinate system, used for the description of the tertiary structure; this coordinate system is composed of a sequence of the torsion angles forming the protein. The secondary structure elements, due to their regular shape, were simply described by a regular sequence of angles. We noticed that, the resulting constraints caused a bottleneck during the search because of poor propagation.

In the current version, we opted to remove the torsion angles description of secondary structure elements and restrict our description of the tertiary structure only to Cartesian coordinates. We simplify and remodel the description of secondary structure elements by means of global coordinate constraints.

Variable Selection Strategy In the original constraint scheme in [14], the variable selection is dynamic. The strategy employed selects, for labeling, a variable which is adjacent to a ground subsequence of the primary sequence. Basically, the method tries to grow the ground part of protein until an acceptable solution is found.

The solution we propose here relies on precomputing the order of labeling of the variables, to avoid the run-time costs of the previous strategy. In the constrain phase, we select the longest secondary structure, and we assign to it ground values. The protein is then partitioned in five consecutive parts: *Left Tail*—containing no secondary structure elements; *Left Body*; *Longest Secondary Structure Element*—ground; *Right Body*; *Right Tail*—containing no secondary structure elements.

The exploration order used during labeling is the following: (i) *Left Body* first (in reverse order), (ii) *Right Body*, (iii) *Left Tail* (in reverse order), and *Right Tail*. The tails are free to assume every conformation—they do not contain any superimposed pattern. Since the main energy contribution is given by the protein body, the tails are instantiated at the end.

Contact revisited We introduce a modification in the computation of the energy function w.r.t. Sect. 3: for our computations we increased the parameter `min_aa_dist` to 3 from 2

(as in [14]). We noted that three consecutive amino acid, s_i, s_{i+1}, s_{i+2} can form different angles (90° or 120°), that cause different energy contributions due to the lattice structure and our energy model. To overpass this asymmetry in local conformations we decided to skip those energy contributions. The intuition of our choice is that the “important” energy contributions that characterize a folding are those provided by pairs of amino acids that are “far apart” in the primary sequence. Indeed, we observed that the contributions deriving from local subsequences tend to mask the global energy evaluation and thus bias our search heuristic.

Experimental Results We compared the new ideas implementation to the previous [14]. The experiments highlight a marked improvement in performance; only in one case we observed a slow-down, while in many cases we achieved from 4 to over 50 fold speedup. On the other hand, although the results are good for relatively small proteins, the new method are not sufficient to produce significant speedups on bigger proteins. In the next section, we introduce a new search heuristic that, combined with the ideas introduced here, is capable to efficiently explore the search tree generated by larger proteins.

4.2 Bounded Block Fails heuristic

In this section, we present a novel heuristic to guide the exploration of the search tree, called *Bounded Block Fails (BBF)*. This technique is general and can be applied to every search with a fixed ordering of variables—though it is particularly effective when applied to the protein folding problem at hand.

The heuristic involves the concept of *block*. Let us assume that V is a list $[V_1, \dots, V_n]$ of variables and constants. A block B_i is a sublist of V of size k composed of unbound variables. The concatenation of all the blocks $B_1 B_2 \dots B_\ell$ gives the ordered list of unbound variables present in V , where $\ell \leq \lceil \frac{n}{k} \rceil$. The blocks are selected dynamically, and they could exclude some of the original variables, that have already been instantiated due to constraint propagation. The number of blocks, thus, could be less than $\lceil \frac{n}{k} \rceil$ and it could be not constant during the whole search. In Figure 2 we depict a simple example for $k = 3$: we consider a list of 9 variables. The dark boxes represent ground assignments.

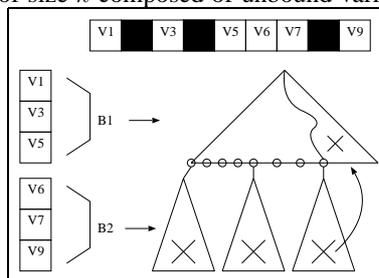


Fig. 2: The BBF heuristics

The heuristics consists of splitting the search among the ℓ blocks. Internally, each block B_i is individually labeled according to the desired labeling strategy—in our case, the same heuristic employed in [14]. When a block B_i has been completely labeled, the search moves to the successive block B_{i+1} , if any. If the labeling of the block B_{i+1} fails, the search backtracks to the block B_i . Here there are two options: if the number of times that B_{i+1} completely failed is below a certain threshold t_i , then the process continues, by generating one more solution to B_i and re-entering B_{i+1} . Otherwise, if too many failures have occurred, then the Bounded Block Fail heuristic generates a failure for B_i as well and backtracks to a previous block. Observe that the count of the number of failures includes both the regular search failures as well as those caused by the Bounded Block Failure strategy. The list t_1, \dots, t_ℓ of thresholds determines the behavior of the heuristic. In the Figure, we assume $t_1 = 3$; the figure shows that, after the third failure of B_2 , the search on B_1 fails as well.

The BBF heuristic is effective whenever the suboptimal solutions are spread sparsely in the search tree and/or for each admissible solution, there are many others with small differences in variables assignments and energy. In these cases, we can afford to skip solutions when generating block failure, because some others are going to be discovered following

other choices in some earlier blocks. For our problem, it is reasonable to keep high threshold values for the first blocks, while exploring only a small fraction of the search space present in the last blocks. Usually, many equivalent solutions can be found, just by changing the assignments in the last blocks, while the core of the problem lies in the first blocks.

In general, this technique can be effectively applied every time the variables are associated to some spatial properties and the corresponding physical object are related to each others, like in the case of a chain of amino acids. When assigning positions following the order of the amino acids on the chain, a failure in the current branch, means that the partial conformation does not allow to proceed without a collision. The BBF heuristic suggests to try to revise some earlier choices instead of exploring the whole space of possibilities depending on the block that collects failures. The high density and the great number of admissible solutions allow us to exclude some solutions, depending on the threshold values, and to still be able to find almost optimal solutions in shorter time.

4.3 Experimental Results

We run some specific tests to measure the benefits of BBF heuristic. On average the times are reduced by 2 times, while for *larger* proteins the benefits are more evident. Moreover the minima found are comparable to the ones obtained without activating the heuristic. Note that the current implementation of the BBF heuristic is performed in SICStus Prolog at a very high-level; e.g., when failing, many `fail` predicates are invoked, with a relatively high cost in the handling of the search tree. This makes the implementation rather inefficient. In spite of this, the BBF heuristics is designed to handle inputs with large sizes. For smaller proteins, the number of blocks tends to be small, and the heuristics accomplishes few block backtracks. In our experiments, we made use of blocks of size 3 (corresponding to the three coordinates of a single amino acid); attempts to reduce this value (i.e., increase the number of blocks) actually produced degradation of performance, due to the excessively small size of the blocks (that defeats the original idea of BBF). To handle larger proteins, it is possible to use larger blocks (e.g. $k = 9$), that include more nodes within the corresponding subtree.

In Table 1, we report a comparison between the results obtained in the original constraint structure [14] and the corresponding ones obtained using the ideas reported in Section 4.1 and the BBF heuristic. In the column *CF* we indicate whether a specific compact factor was used (see Section 2)—when blank, we used a (high) default value (see [14] for more details). In the *AA* column, we indicate the number of amino acids in the protein, while in the *Core Zone* column we identify the subsequence of the protein without the tails—since the tails are less stable and less relevant for a quality test. The *RMSD* column reports the RMSD values (between brackets we indicate the RMSD for the *Core Zone*) of the computed solution w.r.t. the native structure deposited in the PDB (c.f. Section 4.3). At the bottom of the table (marked with (**)), we report two larger proteins, in order to demonstrate the power of the BBF heuristic, using larger size of the BBF blocks ($k = 9$). The thresholds are set to $t_1 = 4$ and $t_{\lceil n/k \rceil} = 2$. The foldings of these proteins were beyond the capabilities of the original constraint structure, while the computation time is extremely low using BBF; furthermore, the RMSD errors are also sufficiently low, thus making this folding a reasonable input to a molecular dynamics refinement step.

RMSD: Discussion The *root-mean-square deviation (RMSD)* is the common measure of the *structural diversity* of two proteins and it measures the average distance between the atoms of two optimally aligned sets of amino acids. In our case, the full-atom reference model with real coordinates is taken from the sequences in the Protein Data Bank [6]. In order to compare this model to our prediction, we extract the backbone of the original protein and obtain a chain of α -carbons (the atoms that are considered by our model).

Protein	CF	AA	Results from [14]			New results with BBF			Core Zone	Speedup
			Time	Energy	RMSD	Time	Energy	RMSD		
1LE0		12	1.3s	-9040	2.8 (2.6)	0.80s	-6251	3.9 (3.2)	2-11	1.6
1KVG		12	7.3s	-14409	2.7 (2.4)	4.12s	-12661	4.1 (4.1)	2-15	1.8
1LE3		16	2.3s	-13653	3.0 (2.7)	1.53s	-11289	4.6 (3.4)	3-11	1.5
1EDP		17	20.4s	-19389	4.3 (1.1)	0.53s	-24492	4.0 (1.1)	9-15	38.4
1PG1		18	14.6s	-10126	6.0 (5.2)	0.83s	-27016	4.2 (3.2)	4-17	17.6
1E0N		27	7m 54s	-12029	5.2 (5.1)	3m	-22308	6.4 (4.5)	3-24	2.6
1ZDD		34	17m 25s	-22350	4.0 (4.0)	1m 51s	-18455	5.3 (5.2)	5-34	9.4
1VII		36	7h 42m	-26408	10.2 (7.8)	3h 40m	-25914	6.0 (5.6)	4-32	2.1
1VII	0.3	36	3h58m	-28710	8.0 (7.4)	22m 41s	-19181	8.2 (8.1)	4-32	10.5
1E0M		37	(*) 24h	-30163	8.9 (4.4)	4h 35m	-26745	9.2 (6.3)	8-29	≥ 5.2
2GP8		40	10h 39m	-29196	4.9 (3.5)	1m 33s	-15187	6.6 (3.6)	6-38	412.3
1ED0		46	9h 38m	-38218	8.0 (7.2)	1h 43m	-31565	6.9 (6.2)	3-40	5.6
1ENH		54	(*) 24h	-23373	9.9 (8.6)	55m 6s	-28559	11.1 (9.5)	8-52	≥ 26.1
6PTI	0.25	58	(*) 48h	-42096	9.7 (9.4)	(*) 48h	-52258	8.0 (7.9)	3-55	1
2IGD	0.17	60	4h 59m	-40588	12.6 (11.5)	2h 35m	-45462	10.6 (10.5)	6-59	1.9
2ERA	0.19	61	(*) 1000s	-38138	11.6 (10.6)	(*) 1000s	-45006	11.9 (11.7)	3-55	1
1SN1	0.25	63	(*) 10h	-47121	8.6 (8.1)	(*) 10h	-47650	9.1 (9.2)	2-51	1
1YPA	0.17	63	(*) 10h	-60244	12.9 (9.8)	(*) 10h	-45617	11.5 (10.5)	12-52	1
1FVS	0.15	72	(**)			11m 49s	-58587	13.1 (13.5)	13-70	-
1B4R	0.16	80	(**)			25m 47s	-78140	13.1 (13.1)	2-79	-

Table 1. Computational results and comparisons.

The RMSD measure contains some intrinsic and unavoidable parts, that derive from the discretization of the backbone on the fcc lattice. To quantify how this problem reflects on the RMSD measure, we run some test of mapping on the lattice the α -carbons of the backbone of some proteins in [6]. The placement fulfills our standard formalization of chain neighbors and allowed angles in the lattice. The RMSD errors due to the discretization range from 3.6Å and 6.6Å. These values are relatively high, even if the fcc lattice is considered one of the best known discrete approximations. These considerations stress the fact that our results are affected by this kind of systematic errors, which are inherent in the use of the fcc lattice and independent from the quality of our solution strategies. As shown in [14], it is possible to eliminate these errors by running some steps of molecular dynamics.

5 Comparison with other heuristics

In this section, we explore the impact of the different solvers on the performance of our protein folding problem solution. In particular, we tested the library `CLP(\mathcal{FD})` of SICStus Prolog 3.12.0 [9] and the libraries `ic` and `branch_and_bound` of ECLiPse 5.8 [10]. The tests have been performed using the optimized constraint system described in the previous sections. The aim is to compare the features of the two constraint solvers and, at the same time, understand which search strategy fits better to the problem.

As first test, we compared the efficiency of the built-in branch-and-bound solver. We coded the same program in SICStus and ECLiPse—in the first case using the `minimize` option offered by the built-in `labeling` predicate, that starts a branch and bound search; for ECLiPse we invoked a `complete` search of `ic` as parameter of the `bb_min` predicate in the `branch_and_bound` library. On average and uniformly, the SICStus solver performs 12.3 times faster than the ECLiPse solver.

On the other hand, ECLiPse offers a number of additional built-in search heuristics, that can be selected when solving a minimization problem. We tested the following built-in

strategies: the *Limited Discrepancy Search (LDS)*, the *Bounded Backtracking Search (BBS)*, and the *Depth Bounded Search (DBS)*—together with LDS.

In the Tables in Figure 3 we report the results obtained from benchmarks for different ECLiPSe heuristics. In these experiments we did *not* use the BBF strategy. Left Table compares SICStus to the LDS and the BBS strategies. Right Table focuses on the DBS strategy. Even if the SICStus solver is faster than ECLiPSe, it does not offer a comparable selection of built-in search strategies. Our pruning heuristic is implemented at a very high level, and thus not as efficient as possible. Nevertheless, it performs better in terms of time and best solution found. None of the tested ECLiPSe strategies is able to produce faster *or* better results than our heuristic in SICStus.

Heuristic	1PG1 (54 Vars)		1E0N (81 Vars)		1E0N	Y=0		Y=1		Y=2	
	Time	Energy	Time	Energy		Time	Energy	Time	Energy	Time	Energy
SICStus	0.73	-27016	78.0	-25493							
Lds(0)	0.22	-20767	No	No	X=1	No	No	21.5	-13715	166	-15382
Lds(1)	1.79	-21412	11.17	-13715	X=2	No	No	30.3	-12688	219	-15382
Lds(2)	8.8	-27082	93.0	-14958	X=3	No	No	45.9	-14052	326	-18046
Bbs(10)	0.3	-20676	1.3	-12017							
Bbs(100)	2.17	-22253	5.6	-12017							
Bbs(1000)	40.8	-27853	75.0	-17221							

Fig. 3. ECLiPSe times (in seconds)

One of the problems we observed is that the ECLiPSe heuristics do not properly fit the problem. Conceptually, we would like to explore the space considering that:

- small changes in the protein folding are not relevant to the global energy—which is the opposite of what the LDS strategy does;
- once a solution is found, it is likely that another interesting solution lies in a complete different part of the search tree—which is in contrast to the BBS strategy, that performs a fixed number of backtracks, and keeps the search in the proximity of the solution found;
- the complete exploration of the first levels of the tree (as done in DBS) is time consuming, and a more selective heuristic should be employed.

Nevertheless, the experiments performed with ECLiPSe highlight that some of the search strategies (e.g., the BBS strategy) can produce solutions that are suboptimal in terms of global energy, but in a significantly shorter period of time. These considerations suggest also that a more efficient search could be performed on a specific solver in which we handle the search tree at low level and implement some new ad-hoc heuristics.

6 A Parallel Solution

In this section we provide an overview of a parallel scheme we designed to solve the protein folding problem. The overall parallelization scheme has been designed as a customization of general or-parallelism techniques [16] to the structure of the problem at hand. The parallel scheme builds on the constraint structure described in the previous sections. The search for solutions is divided among a number of processors, and the search tree is fragmented into subtrees (called *tasks*) and distributed for parallel exploration.

There are two main issues related to the task assignment. The first is that the tasks should be *uniformly* distributed among processors during the execution; this is easier if the number of tasks is large. The second issue relates to constraint propagation and pruning of the search tree through problem-dependent heuristics, that are more effective when applied to large

search trees—i.e., fewer tasks. Hence, the task scheduling strategy should strike a proper balance between these two conflicting requirements.

6.1 Overview of the System

The system is composed of three components: a *loader*, a *scheduler*, and a set of *clients*. Figure 4 shows the components and the main interactions. The loader is a C program, in charge of creating the communication channels (realized using shared memory segments) between the scheduler and the clients. In addition, the loader is in charge of launching both the scheduler and the clients, as parallel processes. The system we developed makes use of a centralized scheduling mechanism. The scheduler is also a C program, that handles the dynamic distribution of tasks to the clients, and implements strategies for load balancing. Each client is a CLP program, that explores the subtree (i.e., task) assigned to the client by the scheduler. When the task is exhausted, the client will notify the scheduler that it is ready to receive an additional task.

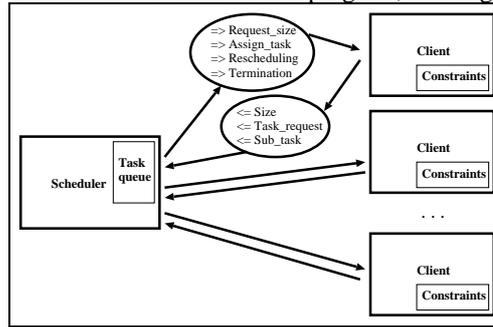


Fig. 4: The parallel system

6.2 Scheduling and Communication

The centralized scheduler implements a *direct scheduling* strategy. It relies on a static partitioning of the search tree, performed according to user defined parameters. During direct scheduling, tasks are assigned to clients upon request.

The scheduler determines the initial pool of tasks to be assigned (*task queue* in Fig. 4)—according to a static expansion of a user-specified number of levels (`Levels`) of the search tree. Since the scheduler is a C program, it does not have access to the collection of constraints; thus, the initial pool of tasks is generated by the clients, and retrieved by the scheduler during the initialization phase. Here, the first client is in charge of precomputing the expansion of `Levels` levels of the search tree and of returning the result of the expansion to the scheduler. The task queue is initialized with a set of subtrees of the search tree, all with roots at the same depth in the tree (`Level`). Each task is described by the list of nodes in the branch that connects the root of the search tree to the root of the task subtree. The scheduler assigns a task to a client whenever the client sends a `task_request` message. The task is assigned to the client by communicating the list of nodes describing it. The message that brings the task to the client is called `Assign_task`.

Due to the irregular structure of the search subtrees—because of the pruning performed by the constraint propagation process and by the heuristics—it is necessary to provide load balancing mechanisms. These mechanisms are employed when the scheduler has an empty task queue, and there is a mix of active and idle clients in the system. The purpose of load balancing is to dynamically generate new, smaller tasks, by further partitioning some of the tasks that are still active. Such smaller tasks can then be reassigned to the idle clients. The load balancing is implemented by a *rescheduling* procedure, activated by the scheduler every time there is at least one idle client and the task queue is empty. In this case, the scheduler selects, with a `Rescheduling` message, the client that has the estimated highest load of work. Due to lack of space, we do not provide technical details on this procedure.

The scheduler takes also care of detecting global termination. Termination occurs when the task queue is empty and task requests have been submitted from all clients. In such a situation, the scheduler returns a `Termination` message to each client.

6.3 Structure of the Client

Each client is a CLP program that implements the process of solving the constraints on a given subset of the search space—i.e., a subset of the domains of the variables in the problem. When launched, a client imports protein data, defines variable domains and applies constraints. After the initial loading, the first client communicates back to the scheduler the list of partial assignments (`Tasks`) obtained from the expansion of the search tree for `Levels` variables. After that, each client starts a loop that (i) delivers a `Task_request` to the scheduler, (ii) waits for the assignment of a task (`Assign_task`) and (iii) executes the task. The processing of a task is based on the CLP scheme described earlier. During each task execution, the client checks for eventual requests for `Rescheduling`—realized by breaking down the labeling process into labeling of smaller lists. If the request is received, the client stops the execution and communicates all the subtasks left to the scheduler. The client stops its execution when it receives the special termination signal.

6.4 Implementation Details

The parallel system has been implemented using a combination of C programming and Constraint Logic Programming (specifically SICStus Prolog 3.12.0 [27]). The activation of the processes implementing the scheduler and the clients is accomplished by standard `fork` calls of C issued by the loader. Separate processes are created for the scheduler and for the different clients.

The communication is realized using shared memory. The clients are written in Prolog, and encapsulate some low level C routines to access the shared memory, for a fast interaction between processes. The message exchange between clients and the scheduler is realized using shared memory queues.

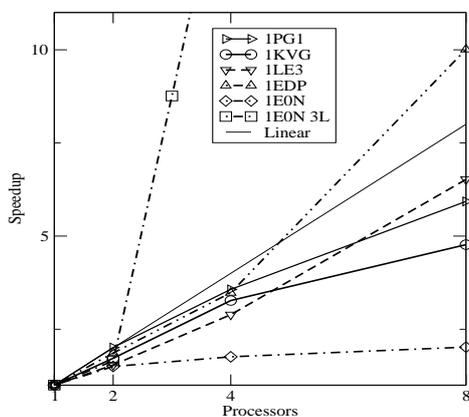
During the execution of Prolog predicates, we also relied on low-level C routines to maintain an efficient representation of the current state of execution of the task. This allows us to maintain a simplified representation of the state of the execution across the branches of the subtree (which are explored via backtracking by the Prolog system) and it simplifies communication with the scheduler during load balancing.

6.5 Experimental results

In this section, we report the experimental results obtained from the execution of our parallel system. The experiments have been performed using a HP RP8400 NUMA architecture, with 14 PA-RISC processors, 8GB RAM, and running HP-UX 11.1.0.

Figure 6.5 plots and shows the times, in seconds, of the parallel execution of the program, using the BBF strategy, no rescheduling and up to 8 parallel clients. We have parallelized the code on numbers of processors that are powers of 2. For each protein, the first 4 `Levels` of the tree are fully expanded. For the protein 1E0N we also included the results obtained using a value of `Levels` equal to 3—as this is sufficient to generate a sufficient number of tasks. The last column reports the execution time, on one CPU of the parallel machine, of the code from [14].

The performance results show that the scalability factors are strongly dependent on the specific problem. We discuss here some of the reasons that generate this behavior. First of all, the subdivision of the search tree into tasks does not imply an equal partition of the workload. In fact, since constraints and heuristics are applied to the trees during the search, their effect



Protein	Processors				[14]
	1	2	4	8	1
IPG1	12.51	6.21	3.5	2.11	48.7
1KVG	13.03	7.56	3.99	2.73	24.9
1LE3	16.36	10.63	5.66	2.51	8.2
1EDP	11.91	6.4	3.42	1.19	67.8
1EON	239.9	160	136.2	118.5	598.4
1EON 3 (Levels)	2249	1367	141.72	82	598.4

Fig. 5. Parallel Execution Time without rescheduling (in seconds)

on a local portion of the search tree can result in a different pruning compared to the same subtree during a sequential search. This effect can have appreciated side-effects: for instance, the superlinear speed-up obtained in the cases of 1EDP and 1EON (3 levels).

Let us observe that, although the sequential speed up for some proteins is of the order of 40 (e.g., for 1EDP—see Table 1) and the parallel speed up is for some proteins of the order of 10 with eight processors (e.g., for 1EDP), the combined speed up is only of the order of 60 (instead of 400). The reason is the overheads associated to the management of parallel tasks—e.g., task exchange, communication costs, and the cost of resetting the constraint store when a new task is acquired. The latter cost, that seems to be the most significant one, could be removed by introducing a more ad-hoc handling of constraints—i.e., bypassing some of the constraint-handling functionalities of SICStus.

Note that clients share their intermediate results by placing them in shared memory, where it becomes accessible to every other client. In particular, the agents can share their current best solution, effectively parallelizing the branch-and-bound process.

Let us assume that the tasks T_1, \dots, T_n are explored in that order in the equivalent sequential algorithm. In the parallel case, every task can take advantage of an earlier-found new bound for the search. For example, T_i and T_j , with $i < j$, are explored in parallel, and a new best solution is found in T_j . This allows a client to prune T_i as well, resulting in a speedup in the search. On the other hand, though, the symmetric case can arise: T_j can be explored extensively because of a late discover of a bound in T_i . In the sequential case T_j would have been pruned from the very beginning, because T_i would have been completed before even starting the exploration of T_j . Thus, the partitioning can dramatically affect the search (see, e.g., protein 1EON). We plan to investigate in the future the design of problem-dependent partitionings, in order to take advantage of the parallel sharing of information.

Let us conclude by observing that, in our experiments, the rescheduling procedure has been rarely effective in improving performance. This is due to the poor interaction between the fairly “sequential” way of rescheduling tasks and the more sophisticated exploration imposed by our heuristic strategies. Work is in progress to adapt rescheduling to better match our search strategies.

6.6 Scalability on Macro Blocks

When dealing with large proteins, it is common to encounter situations where the conformations of various subsequences are already known (e.g., by homology). Thus, the problem of predicting the structure of the whole protein is conceptually equivalent to predicting the placement of few rigid macro blocks, that are linked together by some coils. These ideas could be exploited by our prediction tool.

We defined the following example, to demonstrate that the current model can be feasibly used to attack larger proteins. We use the sequence XYZ, where X and Z are known protein sequences and Y is a fixed-length linking coil of non-interacting amino acids. The idea is to predict the structure of the sequence XYZ, using as input some strong constraints (i.e. Euclidean distances between every pair of amino acids) derived from X and Z.

In our tests, the execution times grow according to the length of the coil Y. As expected, for proteins of this size but with partially known structure, the times are significantly lower than a prediction that uses only the secondary structure information. This allows our system to handle larger protein complexes.

7 Conclusion and Future Work

In this paper, we provided a formalization of the protein folding problem on face-centered cubic lattice structures. The formalization has been transformed in a constraint system, and solved using constraint solving over finite domains. We analyzed different ways of organizing the constraint structure, and different heuristics and search strategies to solve them. We presented and tested a new search heuristic (Bounded Block Fails), well suited for this problem. We also provided a way to parallelize the process of exploring the search space, allowing concurrent constraint solvers to cooperate in the search of an optimal folding.

The results collected from the different approaches to the problem (sequential search strategies, parallel implementations, different implementations of solvers) converge on the need of a new dedicated and efficient solver. The analysis in Section 5 suggests that the only way to significantly improve our framework is to access the search tree at a lower level, and to implement new heuristics more suitable to the problem. Since this is not allowed by the current implementations of SICStus and ECLiPSe, we plan to develop our own ad-hoc constraint solver. The new solver will be dedicated to problem on lattices (and thus more efficient) and will implement ad-hoc search strategies. The new solver, applied to the protein folding problem of fcc, is expected to allow us to tackle larger problems—the final goal is to manage proteins composed of 500 amino acids. As shown in Table 1, currently we can solve (without using scalability—Section 6.6) proteins of length in the order of 80 amino acids.

We also plan to continue the development of the parallel solution; in particular, we intend to develop rescheduling strategies that better match the heuristics employed by the sequential system, allowing for a more effective load balancing and scalability.

Acknowledgments We thank Federico Fogolari, for his guidance and comments. The research has been partially supported by NSF grants HRD-0420407 and EIA-0220590, by the MIUR Project *Sybilla*, and by GNCS 2005.

References

1. R. Agarwala et al. Local rules for protein folding on a triangular lattice and generalized hydrophobicity in the HP model. *J. of Computational Biology*, pages 275–296, 1997.
2. K. R. Apt. *Principles of constraint programming*. Cambridge University press, 2003.

3. R. Backofen. The protein structure prediction problem: A constraint optimization approach using a new lower bound. *Constraints*, 6(2–3):223–255, 2001.
4. R. Backofen and S. Will. Fast, constraint-based threading of HP sequences to hydrophobic cores. *Int. Conf. on Principle and Practice of Constraint Programming*, 494–508, 2001. Springer Verlag.
5. R. Backofen and S. Will. A Constraint-Based Approach to Structure Prediction for Simplified Protein Models that Outperforms Other Existing Methods. *ICLP*, 2003, Springer Verlag.
6. H. M. Berman et al. The protein data bank. *Nucleic Acids Research*, 28:235–242, 2000. <http://www.rcsb.org/pdb/>.
7. M. Berrera, H. Molinari, and F. Fogolari. Amino acid empirical contact energy definitions for fold recognition in the space of contact maps. *BMC Bioinformatics*, 4(8), 2003.
8. R. Bonneau and D. Baker. Ab initio protein structure prediction: progress and prospects. *Annu. Rev. Biophys. Biomol. Struct.*, 30:173–89, 2001.
9. M. Carlsson, G. Ottosson, and B. Carlson. An open-ended finite domain constraint solver. *PLILP*, Springer Verlag, 1997.
10. A. M. Cheadle et al. ECLiPSe: An Introduction. Technical Report IC-Parc 03–1, IC-Parc, 2003.
11. D. Clark, J. Shirazi, and C. Rawlings. Protein topology prediction through constraint-based search and the evaluation of topological folding rules. *Protein Engineering*, 4:752–760, 1991.
12. P. Clote and R. Backofen. *Computational Molecular Biology: An Introduction*. John Wiley & Sons, 2001.
13. P. Crescenzi et al. On the complexity of protein folding. In *Proc. of STOC*, pages 597–603, 1998.
14. A. Dal Palù, A. Dovier, and F. Fogolari. Constraint logic programming approach to protein structure prediction. *BMC Bioinformatics*, 5(186), 2004.
15. A. Dal Palù, A. Dovier, and F. Fogolari. Protein folding in $CLP(\mathcal{FD})$ with empirical contact energies. In *Recent Advances in Constraints*, Springer Verlag, 2004.
16. G. Gupta, E. Pontelli, M. Carlsson, M. Hermegildo, K. Ali. Parallel Execution of Prolog: a Survey. *ACM TOPLAS*, 23(4):472–602, 2001.
17. F. Fogolari et al. Modeling of polypeptide chains as C- α chains, C- α chains with C- β , and C- α chains with ellipsoidal lateral chains. *Biophysical Journal*, 70:1183–1197, 1996.
18. S. Forman. *Torsion Angle Selection and Emergent Non-local Secondary Structure in Protein Structure Prediction*. PhD thesis, U. of Iowa, 2001.
19. W. Hart and S. Israil. Fast protein folding in the hydrophobic-hydrophilic model within three-eighths of optimal. *J. of Computational Biology*, pages 53–96, 1996.
20. S. Miyazawa and R. L. Jernigan. Residue-residue potentials with a favorable contact pair term and an unfavorable high packing density term, for simulation and threading. *J. of Molecular Biology*, 256(3):623–644, 1996.
21. A. Newman. A New Algorithm for Protein Folding in the HP Model. In *Symposium on Discrete Algorithms*. Springer Verlag, 2002.
22. Committee on Mathematical Challenges from Computational Chemistry. National Research Council, 1995.
23. G. Raghunathan and R. L. Jernigan. Ideal architecture of residue packing and its observation in protein structures. *Protein Science*, 6:2072–2083, 1997.
24. B. Rost and C. Sander. Prediction of protein secondary structure at better than 70% accuracy. *J. Mol. Biol.*, 232:584–599, 1993.
25. J. Skolnick and A. Kolinski. Computational studies of protein folding. *Computing in Science and Engineering*, 3(5):40–50, 2001.
26. J. Skolnick and A. Kolinski. Reduced models of proteins and their applications. *Polymer*, 45:511–524, 2004.
27. Swedish Institute for Computer Science. Sicstus Prolog. www.sics.se/sicstus/.
28. L. Toma and S. Toma. Folding simulation of protein models on the structure-based cubo-octahedral lattice with contact interactions algorithm. *Protein Science*, 8:196–202, 1999.
29. M.L.G. William and D. Harvey. Limited discrepancy search. *IJCAI*, pages 607–615, Morgan Kaufmann, 1995.