

METHODOLOGY ARTICLE

Open Access

An algorithm to enumerate all possible protein conformations verifying a set of distance constraints

Andrea Cassioli⁴, Benjamin Bardiaux^{1,2}, Guillaume Bouvier^{1,2}, Antonio Mucherino⁶, Rafael Alves⁴, Leo Liberti^{4,5}, Michael Nilges^{1,2}, Carlile Lavor³ and Thérèse E Malliavin^{1,2*}

Abstract

Background: The determination of protein structures satisfying distance constraints is an important problem in structural biology. Whereas the most common method currently employed is simulated annealing, there have been other methods previously proposed in the literature. Most of them, however, are designed to find one solution only.

Results: In order to explore exhaustively the feasible conformational space, we propose here an interval Branch-and-Prune algorithm (*iBP*) to solve the Distance Geometry Problem (DGP) associated to protein structure determination. This algorithm is based on a discretization of the problem obtained by recursively constructing a search space having the structure of a tree, and by verifying whether the generated atomic positions are feasible or not by making use of pruning devices. The pruning devices used here are directly related to features of protein conformations.

Conclusions: We described the new algorithm *iBP* to generate protein conformations satisfying distance constraints, that would potentially allows a systematic exploration of the conformational space. The algorithm *iBP* has been applied on three α -helical peptides.

Keywords: Distance geometry, Branch-and-prune algorithm, Molecular conformation, Protein structure, Nuclear magnetic resonance

Background

Protein structure determination is crucial for understanding protein function, as it paves the way to the discovery of new chemical compounds and of new approaches to control the biological processes.

The problem of protein structure determination is certainly a problem with multiple solutions, as proteins are flexible polymers. As most of the experimental techniques of the structural biology obtain measurements averaged on an ensemble of protein conformations, the usual approaches for structure determination intend to find an average structure or a set of conformations describing fluctuations around an average structure. A path intending to get a complete coverage of the conformational

space, given a series of constraints, is usually not taken, although such an approach could provide precious information about the conformational equilibrium, which is essential in the function of many proteins, as the HIV protease [1].

An important class of experimental methods for protein structure determination is based on the measurement of inter-atomic distances and angles, such as Nuclear Magnetic Resonance (NMR) spectroscopy [2] and cross-linking coupled to mass spectrometry [3]. In NMR, distance intervals between hydrogens are determined from the measurement of nuclear Overhauser effects (NOE). The experimentally measured distances are then used as constraints for protein structure calculations. Pure *in silico* approaches have been also developed based on the use of inter-atomic distance constraints, such as homology modeling [4] or prediction of protein-protein complexes [5] and ligand poses [6].

*Correspondence: therese.malliavin@pasteur.fr

¹Institut Pasteur, Structural Bioinformatics Unit, 25, rue du Dr Roux, 75015 Paris, France

²CNRS UMR3528, 25, rue du Dr Roux, 75015 Paris, France

Full list of author information is available at the end of the article

The Distance Geometry Problem (DGP) [7,8] consists in identifying the sets of points which satisfy a set of constraints based on relative distances between some pairs of such points. The present work describes an algorithm developed to solve DGP in the context of protein structure determination: the points represent the protein atoms.

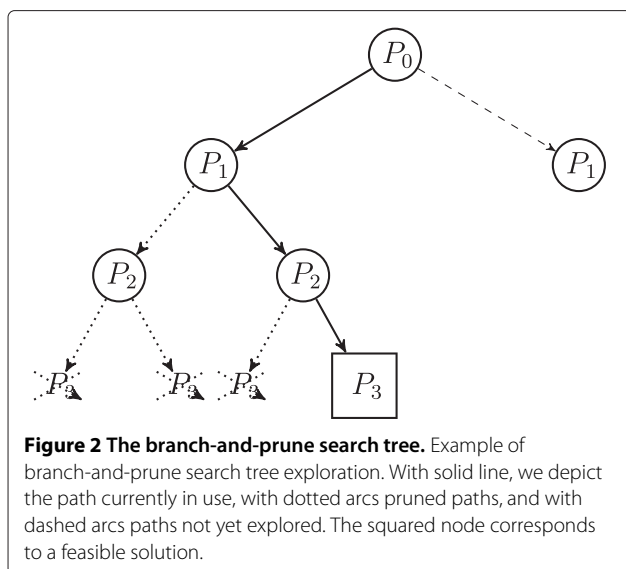
The DGP is a constraint satisfaction problem. Several approaches solve this problem by reformulating it [8] as a global optimization problem having a continuous search domain, and whose objective function is generally a penalty function designed to measure the violation of the distance constraints. Over the years, the solution of DGPs arising in structural biology have been typically attempted by Simulated Annealing (SA) approaches based on molecular dynamics [9]. Other proposed approaches are based on various optimization methods as in [10]. As all meta-heuristic approaches, SA may provide approximate solutions but does not deliver optimality certificates. In the case of protein structure determination, since the optimization problem is a reformulation of a constraint satisfaction problem, solutions given by SA can be successively verified by checking the violations of the distance constraints. However, additional solutions may exist but go undetected by SA. Thus, an algorithm for the systematic enumeration of the possible conformations of a given protein could find a widespread field of application. Branch-and-prune algorithms and similar were proposed in the general context of protein structure determination [11-16], (see also [8] and references therein). However, these studies primarily addressed the question of defining relative orientations of protein monomers in symmetric oligomers, not the determination of all possible conformation of a polypeptide chain with a very large number

of degrees of freedom from distance constraints. Systematic exploration was proved to be useful in the case of residual dipolar couplings (RDC) constraints [17], for exploring the sidechains conformations [18,19] and for assignment of NOEs, provided that the structure is known [20]. For the structure determination from RDCs, it has been shown [21] that when using RDCs but only sparse NOEs the problem can be solved in polynomial time. Such approaches have also been used for structure determination in X-ray crystallography for non-crystallographic symmetry by orienting and translating symmetric protein subunits [22]. To the best of our knowledge, in this paper we present the first application of a Branch-and-Prune algorithm to the problem of full protein structure determination based on unambiguous distance information.

Under certain conditions, DGPs can be discretized [23] (see below), which means that the search domain for the corresponding optimization problem can be reduced to a discrete set, which has the structure of a tree. The discretization makes the enumeration of the entire solution set of DGP instances possible. This is important when the experimental constraints do not specify the protein conformation uniquely, i.e., more than one conformation satisfies all constraints. For solving discretized DGP, we employ an *interval* branch-and-prune (*iBP*) algorithm [24], which is based on the idea of recursively exploring the tree while generating new candidate atomic positions (branching phase) and to verify the feasibility of such positions (pruning phase) (Figure 1). By making use of pruning devices, branches rooted at infeasible positions can be discarded from the tree, so that the search can be reduced to the feasible parts of the tree (Figure 2). Pruning devices can be conceived and integrated in *iBP* to improve

Algorithm 1: The <i>iBP</i> recursive algorithm.		
Input: atom index l , total number of atoms n , solution x		
1	if $l = n$ then	
	/* Solution found!	*/
2	return	
	/* Branching	*/
3	compute set P_l of possible position of atom l ;	
4	foreach $p \in P_l$ do	
	/* Check for infeasibility (Pruning)	*/
5	if p is feasible then	
	/* Value accepted	*/
6	$x^l \leftarrow p$;	
	/* Go to the next level	*/
7	<i>iBP</i> ($l + 1, n, x$);	
8	end	

Figure 1 The *iBP* recursive algorithm. Description of the *iBP* algorithm.



the performances of the pruning phase and thus of the algorithm.

In the present work, we first describe the branching phase and the pruning devices used to determine the solutions to the Distance Geometry problem. Then, an overall view of the method is given along with the use of the branching and pruning devices at different steps and the complexity of the algorithm is analyzed. We finally illustrate the algorithm application with three proteins for which α -helical regions are known along with few long-range NMR constraints (ie. constraints measured between residues i and j such that $|i - j| > 3$ in the protein sequence). The obtained conformations display good stereochemical quality parameters, and the conformational space explored is larger than the one sampled with traditional optimization methods such as simulated annealing.

Methods

In order to sample the conformational space of a protein, we use a Branch-and-Prune algorithm to build a tree in which each node represents a solution for one atomic position. We limit ourselves in the present work to the calculation of the backbone and $C\beta$ atomic coordinates.

The constraints used to generate atomic coordinates along the Branch-and-Prune algorithm are the following:

1. covalent distance constraints corresponding to bond lengths and bond angles, whose values are derived from high-resolution small molecule X-ray crystal structures [25];
2. NMR distance constraints;
3. van der Waals radii of atoms between non-bonded atom pairs (i, j) : a fraction of the sum of the van der

Waals radii of each atom provides a lower bound to the corresponding inter-atomic distances:

$$d_{ij} \geq \sigma (r_i^{vdw} + r_j^{vdw}), \quad (1)$$

where $\sigma \in [0, 1]$, and is typically around 0.85. The values for the radii are given in Table 1 [26,27]. These lower bounds apply only in the cases where no larger lower bound has been determined from NMR distance constraints;

4. distances derived from the backbone torsion angles ϕ and ψ ;
5. hydrogen bonds in α -helix;
6. amino-acid chirality;
7. α -helix geometry.

The atom coordinates are calculated, one by one, following the atom order P_{ato} described in Figure 3 and previously proposed in [24]. In this order, some atoms are repeated to insure that any entered atom is defined by distance constraints with respect to three preceding atoms in P_{ato} [24]. The carbonyl oxygens and the atoms $C\beta$, which were not present in the order P_{ato} , are calculated separately.

Then, the tree is built using a recursive procedure to create each node of the tree. This procedure is called branching phase. The created nodes are then submitted to the pruning devices in order to decide whether the node should be kept or removed. If the node is removed, the possible branches starting from this node are also pruned. A pruning device is responsible for checking whether a partial solution is feasible, i.e. to check whether a set of embedded atoms fulfill the constraints (1)-(7) described above.

In the following, we describe the branching phase and the pruning devices. Then, the complexity of the algorithm is described from a theoretical point of view, before presenting some application cases.

Branching devices

The tree parsed during iBP is formed by nodes, each corresponding to one set of atomic coordinates from the order P_{ato} (Figure 3) [24]. At each level of the tree, the atomic coordinates of the corresponding atom are calculated by making use of a recursive procedure, called branching phase. The current atom position is defined by distance constraints to three other atoms. These distances are obtained from the constraints (1-3) described above: (1) the covalent constraints, (2) the NMR distance constraints, (3) the van der Waals radii.

Table 1 Van der Waals radii (see [26] and [27])

atom	O	H	C	N
r^{vdw} (Å)	1.4	1.0	1.7	1.5

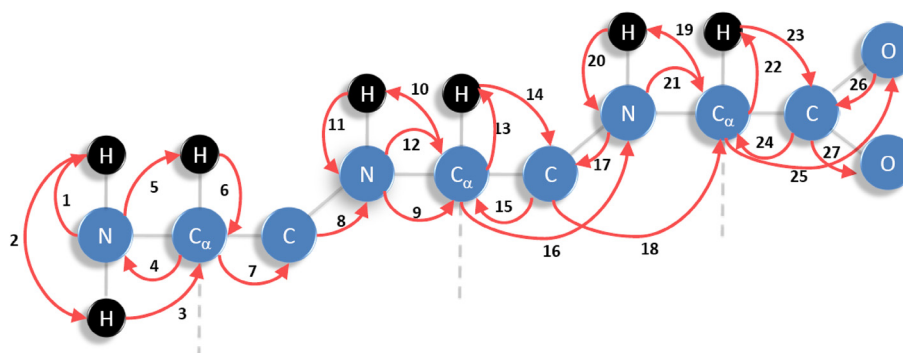


Figure 3 Order P_{ato} of the atoms parsed during the branch-and-prune algorithm.

If the distance constraints specify a unique value rather than an interval, this signifies that the distances to three immediate predecessors from the current vertex are known: these are the centers of the three spheres, and the distances are the radii of these spheres. The position of the current vertex/atom is thus defined by the intersection of three spheres, so there are at most two solutions for the current atom position: this is called a 2-branching situation (Figure 4).

When a distance is not uniquely defined, but rather defined by lower and upper bounds, i.e. $d_{ij} \in [l_{ij}, u_{ij}]$, this distance is uniformly discretized by sampling $b \geq 1$ values in $[l_{ij}, u_{ij}]$, as depicted in Figure 5.

$$\tilde{d}_i = \left\{ l_{i,i-3} + (t-1) \frac{(u_{i,i-3} - l_{i,i-3})}{b} : t = 1, \dots, b \right\}. \quad (2)$$

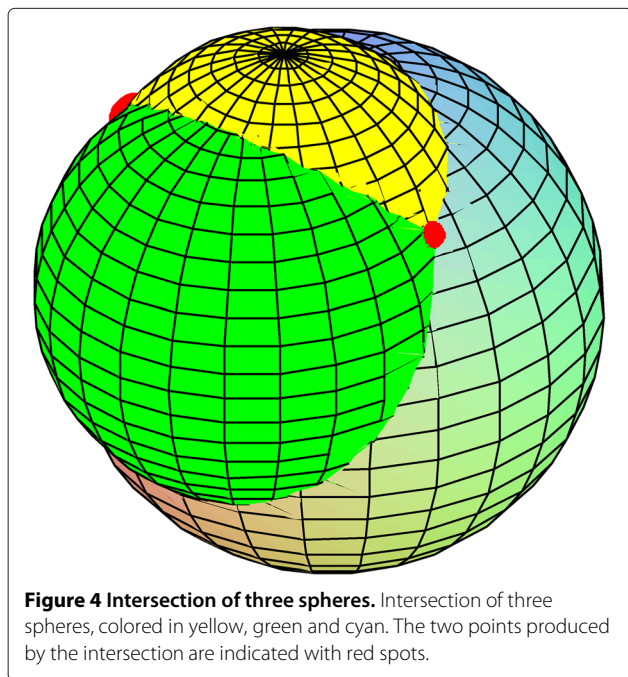


Figure 4 Intersection of three spheres. Intersection of three spheres, colored in yellow, green and cyan. The two points produced by the intersection are indicated with red spots.

In this case, we have a b-branching situation.

The algorithm used for calculating the atom coordinates is then applied to each set of \tilde{d}_i values sampled for the distance constraints. The choice of the *discretization factor* b is a crucial point: a small value might lead to an infeasible problem because we may not select any feasible distance; a larger value increases the computational burden. In general, the finer the discretization, the more accurate the computation is, but it is not trivial to figure out the optimal value for b . One way to choose b is to consider that the number of nodes in the search tree is bounded by

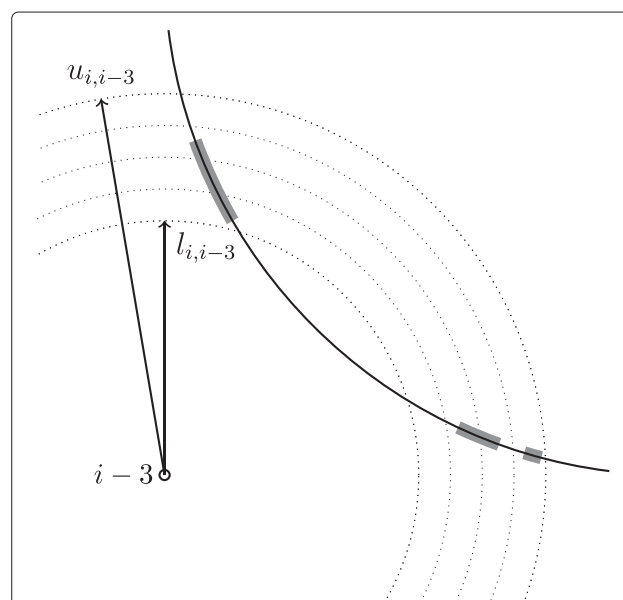


Figure 5 Discretization of the distance constraints. An example of discretization of the distance d_{ij-3} . The solid circle represents the result of the intersection of the spheres centered in $i-1, i-2$ with radii d_{ij-1}, d_{ij-2} , respectively. The distance d_{ij-3} is discretized accordingly to Equation 2 with $b = 5$: dotted circles represent the intersections of spheres centered in $i-3$ with radii in \tilde{d}_i with the plane containing the $i-3, i-2$ and $i-1$. Thick gray arcs represent the feasible regions for the atom i .

$3 + (2^l b^k)$, where l is the number of tree levels where we have a 2-branching situation, and k is the number of tree levels where we have a b-branching situation [28]. Appropriate values of b should result in a manageable number of nodes.

Given the position of the three previous atoms $k-3, k-2, k-1$ in the order P_{ato} and given the constraints to these atoms of the atom k to be embedded, the position of k is calculated by a recursive matrix multiplication by making use of the set of distances $d = \{d_{k,k-1}, d_{k,k-2}, d_{k,k-3}\}$ between the previous atoms and k . Although there are several methods to compute sphere intersections [29], in our experience, the best trade-off between efficiency and numerical stability is given by the use of recursion matrices [23], and of the two following angles: (i) the torsion angle ω_3 formed by atoms $\{k, k-1, k-2, k-3\}$ which depends on the distance between k and $k-3$, (ii) the angle θ_2 formed by atoms $\{k, k-1, k-2\}$.

The recursion is applied through the equation:

$$\begin{aligned} \begin{bmatrix} x_k \\ y_k \\ z_k \\ 1 \end{bmatrix} &= B_1 B_2 B_3 \dots B_k(d, \sigma) \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \\ &= Q_{k-1} B_k(d, \sigma) \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} = Q_k \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \end{aligned} \tag{3}$$

where:

$$B_k(d, \sigma) = \begin{bmatrix} -\cos \theta_2 & -\sigma \sin \theta_2 & 0 & -d_{k,k-1} \cos \theta_2 \\ \sigma \sin \theta_2 \cos \omega_3 & -\cos \theta_2 \cos \omega_3 & -\sin \omega_3 & \sigma d_{k,k-1} \sin \theta_2 \cos \omega_3 \\ \sigma \sin \theta_2 \sin \omega_3 & -\cos \theta_2 \sin \omega_3 & \cos \omega_3 & \sigma d_{k,k-1} \sin \theta_2 \sin \omega_3 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \tag{4}$$

and $\sigma \in \{+1, -1\}$. The series of recursion matrices is initialized as:

$$\begin{aligned} B_1 &= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, B_2 = \begin{bmatrix} -1 & 0 & 0 & -d_{2,1} \\ 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \\ B_3 &= \begin{bmatrix} -\cos \theta_3 & -\sin \theta_3 & 0 & -d_{3,2} \cos \theta_3 \\ \sin \theta_3 & -\cos \theta_3 & 0 & d_{3,2} \cos \theta_3 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}. \end{aligned} \tag{5}$$

$d_{2,1}$ being the distance between the first and the second atom, and $d_{3,2}$ the distance between the third and the second atom in the order P_{ato} .

The total number of B_k matrices to be calculated along the parsing of the tree is bounded by $2 |P_{\text{ato}}| b$, where $|P_{\text{ato}}|$ is the size of the ordered atom list P_{ato} . The product $Q_{k-1} B_k$ is calculated in two steps: (1) the fourth column of

Q_k , which gives us the coordinates of k , is computed; (2) only if k is not pruned, the three remaining columns are computed.

We must distinguish two cases when embedding an atom k . If it is the first appearance of k in P_{ato} , we use equation 3 to compute all possible embeddings of k for $\sigma \in \{+1, -1\}$ and the set of distances d . If it is not the first appearance of k in P_{ato} , we need to take into account the fact that numerical instabilities generate matrices which will lead to slightly different coordinates for k than those computed the first time. In order to decrease the impact of these numerical errors, we compute the set of distances d , the angles θ_2, ω_3 and for $\sigma \in \{+1, -1\}$ the corresponding matrices $B_k(d, +1), B_k(d, -1)$, which lead to two possible embeddings of k (Equation 3), as $k^+ = Q_{k-1} B_k(d, +1)$ and $k^- = Q_{k-1} B_k(d, -1)$. We choose the value of k that yields the updated coordinates of k being the closest to the previous coordinates of this atom.

Each carbonyl oxygen O^{i-1} is uniquely determined for residue i , once C^{i-1}, N^i and H^i have been embedded, since these atoms are all part of the peptide plane [30]. As is common practice (see, e.g., [31-33]), we fix here the torsion angle ω of the peptide plane to -180° or 0° . In a previous implementation [34], the positions of the carboxylic oxygens were not stored. Although this approach leads to memory savings, the availability of carboxylic oxygen positions can improve the definition of the α -helix secondary structure.

The positions of the carbonyl oxygens are thus now calculated in the following way. If $k = O^{i-1}$ is the carboxylic oxygen atom located at the vertex k , and $\{v_1, v_2, v_3\}$ are the vertices corresponding to atoms $\{C^{i-1}, N^i, H^i\}$, belonging on the same peptide plane π , we denote n_π the normal vector to π . The coordinates of k can then be computed by solving the following non-linear system:

$$\begin{cases} \|k - v_i\|^2 = d_{ki}^2, & i = 1, 2, 3 \\ n_\pi^T (v_1 - k) = 0 \end{cases} \tag{6}$$

where d_{ki} are the distances between atoms k and i . Using an approach similar to those employed in [35], we obtain the equivalent linear system:

$$\begin{cases} 2(v_2 - v_1)^T k = d_{k1}^2 - d_{k2}^2 - \|v_1\|^2 + \|v_2\|^2 \\ 2(v_3 - v_1)^T k = d_{k1}^2 - d_{k3}^2 - \|v_1\|^2 + \|v_3\|^2 \\ n_\pi^T (v_1 - k) = 0 \end{cases} \tag{7}$$

The parameter d_{k1} is the length of the bond connecting O^{i-1} and C^{i-1} , the parameters d_{k2} and d_{k3} are the distances between $k = O^{i-1}$ and N^i, H^i , calculated from bond angles and bond lengths between atoms of the peptide plane, and the angle ω of 180° in a *trans* peptide plane. The case of the *cis* peptide plane can be treated in the same way, modifying the value of ω to be 0° .

Following the idea proposed for carbonyl oxygens, the coordinates k of a $C\beta$ atom can be computed from previously calculated atoms, because the four distances of k to atoms $\{v_1 = C\alpha, v_2 = H\alpha, v_3 = N, v_4 = C\}$ are exactly known, and because these five atoms are not coplanar. The coordinates k are calculated by solving the linear system:

$$\begin{cases} 2(v_2 - v_1)^T k = d_{k1}^2 - d_{k2}^2 - \|v_1\|^2 + \|v_2\|^2 \\ 2(v_3 - v_1)^T k = d_{k1}^2 - d_{k3}^2 - \|v_1\|^2 + \|v_3\|^2 \\ 2(v_4 - v_1)^T k = d_{k1}^2 - d_{k4}^2 - \|v_1\|^2 + \|v_4\|^2 \end{cases} \quad (8)$$

The parameter d_{k1} is the length of the bond connecting $k = C\beta$ and $C\alpha$, the parameters d_{k2} , d_{k3} and d_{k4} are the distances between $k = C\beta$ and $H\alpha$, N , C , calculated from bond angles and bond lengths between these atoms.

Pruning devices

Once the set of possible coordinates of the atom k has been determined in the branching phase described above, pruning devices are used to check whether the coordinates of k are feasible. In some cases described below, the coordinates of k along with the coordinates of previously embedded atoms are checked together. If the check is negative, the solution obtained for k is discarded, which prunes all tree branches originating from the node k . In this section, we present the pruning devices used to accept or discard the coordinates of the atom k generated by the branching devices. The pruning device applies all these tests as soon as the involved atoms have been embedded.

Direct distance feasibility (DDF)

As the coordinates for an atom k are determined, we first check that all distances between k and the other embedded atoms respect the given lower and upper bounds arising from the constraints (1-3) listed in section "Solving the DGP with iBP ".

Torsion angle feasibility (TAF)

The values of the backbone torsion angles ϕ , ψ , are used as a pruning device, checking whether they are located in the permitted regions of the Ramachandran plot. The pruning device, first introduced in [34], is implemented in the following way. The torsion angle ξ_{ijkl} defined by a quadruple of atoms $\{i, j, k, l\}$ falls into a domain Ξ_{ijkl} , up to a certain tolerance $\epsilon_t > 0$. In general, Ξ_{ijkl} is the union of κ dis-jointed intervals, i.e.

$$\Xi_{ijkl} = \bigcup_{c=1}^{\kappa} \Xi_{ijkl}^c \quad (9)$$

From the bounds on a torsion angle ξ_{ijkl} it is possible to derive bounds on the distance d_{il} , noticing that

$$d_{il}(\xi_{ijkl}) = \sqrt{d_{ij}^2 + d_{jk}^2 - 2(\cos(\xi_{ijkl})\sqrt{ef} + bc)d_{ij}d_{jk}}, \quad (10)$$

where:

$$\begin{aligned} b &= \frac{1}{2} \frac{d_{lj}^2 + d_{jk}^2 - d_{lk}^2}{d_{ij}d_{kj}} \\ c &= \frac{1}{2} \frac{d_{ij}^2 + d_{jk}^2 - d_{ik}^2}{d_{ij}d_{jk}} \\ e &= 1 - b^2, f = 1 - c^2. \end{aligned} \quad (11)$$

Taking the maximum and minimum values of $d(\xi_{ijkl})$ for $\xi_{ijkl} \in \Xi_{ijkl}$, we obtain an interval $[l_{il}, u_{il}]$ for the distance d_{il} . The sign of the angle ξ_{ijkl} is used as an additional pruning criterion along with the d_{il} interval.

Dijkstra shortest-path (DSP)

As introduced in [23], we can exploit the fact that the distances are Euclidean to improve the iBP pruning capabilities. We extend and generalize the procedure presented in [36] in the following way. We introduce an auxiliary graph G^+ with the same topology as the graph connecting the atoms in the protein, but such that the weight of each edge (i, j) is the upper bound of the distance d_{ij} . For every pair of atoms i, j , the shortest-path between i, j in G^+ is a valid over-estimate of d_{ij} . Thus we used an all-to-all shortest-path algorithm, the Floyd-Warshall algorithm [37], to refine the upper bound for each pair of atoms.

The Dijkstra Shortest-Path pruning device uses the refined upper bounds of inter-atomic distances in the following way. According to Lemma 4 in [23], for an atom k and for each atom pair i, j such that $i < j < k$ in the order P_{ato} and for which d_{ik} is known, the embedding of k can be pruned if:

$$\|i - j\| - d_{ik} > u_{jk} \quad (12)$$

where u_{jk} is the upper bound of the atom pair (j, k) obtained using the Floyd-Warshall algorithm [37].

Chirality (CHI)

The pruning of atom coordinates through the amino-acid chirality is implemented through the so-called CORN rule of thumb: in amino acids, the groups COOH, R (sidechain), NH₂ and H are bonded to the chiral center $C\alpha$ carbon. Starting with the hydrogen atom away from the viewer, if these groups are arranged clockwise around the $C\alpha$ carbon, then the amino-acid is in the D-form. If these groups are arranged counter-clockwise, the amino-acid is in the L-form. The CORN rule was restated by imposing that the torsion angle defined by the atoms $C, C\beta, N, H\alpha$ of residue i for the D-form or $C, N, C\beta, H\alpha$ of residue i for the L-form, is positive.

α -helix secondary structure

We proposed the use of α helix information as a pruning device in the context of the iBP algorithm first in [34].

The α helix location can be determined from an analysis of the NMR chemical shifts by TALOS [38]. Four criteria are used to enforce the formation of an α helix: (i) the formation of backbone hydrogen bonds between amide hydrogens and carbonyl oxygens, (ii) the alignment of the amide and carbonyl functions checked by a qualitative condition on the energy of the hydrogen bond, (iii) the definition of backbone ϕ and ψ torsion angles already described in the Torsion Angle Feasibility, (iv) the definition of three additional angles θ , θ' and θ'' similar to the ones introduced by Grishaev et al. [39].

On a sequence of $m + 1$ contiguous residues $I_\alpha = \{i, i + 1, \dots, i + m\}$ forming an α helix, for any pair of residues $(i - 4, i)$ belonging to I_α , the lower and upper bounds on the distance between the carboxylic oxygen O^{i-4} and the amide hydrogen H^i should be compatible with the formation of an hydrogen bond. The upper and lower bounds are defined in an input parameter file of *iBP*, and were set to 1.9 and 3.0 Å in the present work.

The condition checking the alignment of atoms involved in the hydrogen bond is implemented by calculating a local energy information defined in the DSSP package [40]:

$$q_1 q_2 \left[\frac{1}{d_{O_{i-4}N_i}} + \frac{1}{d_{C_{i-4}H_i}} - \frac{1}{d_{O_{i-4}H_i}} - \frac{1}{d_{C_{i-4}N_i}} \right] \cdot f < -0.5, \quad (13)$$

with $q_1 = 0.42$, $q_2 = 0.2$ and $f = 332$, and d_{AB} correspond to the distance between atoms A and B .

The last criterion enforces the angles θ , θ' , θ'' to be respectively into the interval values $0/70^\circ$, $0/90^\circ$ and $110/180^\circ$.

Implementation details

In this section we provide an overview of the main implementation features. The *iBP* algorithm has been coded in C++ with extensive use of template meta-programming [41], STL [42,43], and BOOST (www.boost.org). Linear systems, as for instance (7), are solved using the LAPACK library [44].

Discretizable DGP instances were represented by simple weighted undirected graphs $G = (V, E, d)$, which were handled by the Boost Graph Library (BGL) [45]. The points in \mathbb{R}^3 were represented using the Boost Geometry Library (also known as Generic Geometry Library, GGL: www.boost.org).

Constraints on distances, angles or energy are typically expressed by enforcing a variable x to take values in a domain \mathcal{D} , which is generally the union of intervals and singletons:

$$\mathcal{D} = \left\{ \bigcup_{j=1}^m \bar{x}_j \right\} \cup \left\{ \bigcup_{i=1}^k [x'_i, x''_i] \right\}. \quad (14)$$

The Boost Interval Library (BIL – see [46,47]) was used to store such representation, and to perform basic operations for intervals and singletons. On top of the BIL, we define the type `domain` which contains a set of intervals and operations as intersection, scaling, etc. The BIL allows also to select the underlining data format for the interval (single/double precision real, integer).

Theory

In this section we give some details about the worst-case asymptotic complexity behavior of the *iBP* algorithm. The description given above includes many details which are useful for finding the structure of proteins but which somewhat complicate the precise mathematical treatment. We first give a very brief abstract description of the *iBP* and of the formal problem it solves, and then proceed to discuss its complexity.

Formally speaking, the DGP is the following decision problem: given an integer $K > 0$, a simple undirected graph $G = (V, E)$ and an edge weight function $d : E \rightarrow \mathbb{R}_+$, is there a realization $x : V \rightarrow \mathbb{R}^K$ such that for each $\{u, v\} \in E$ we have $\|x_u - x_v\|_2 = d_{uv}$? Note that we are writing x_u for $x(u)$ and d_{uv} for $d(u, v)$. We also remark that in the more “applied” interpretation given in the preceding section, the range of the edge function d is $\mathbb{I}\mathbb{R}_+$, i.e. the set of all non-negative closed real intervals, and $K = 3$. The DGP is **NP**-hard for any $K > 1$ and **NP**-complete for $K = 1$ [48]. Since we are interested in finding *all* solutions of the DGP rather than just one, we denote by X the set of all realizations of G .

Assumptions on the DGP input data

In fact, due to the fact that our data come from a protein structure setting, we can also make the following assumptions about G and d :

1. there is an order $1, 2, \dots, n$ on the vertices such that $1, 2, 3$ is a triangle in the graph G and, for each vertex $v > 3$, v is adjacent to $v - 1, v - 2, v - 3$;
2. the set of edges E can be partitioned in two subsets E_D and E_P , such that E_P consists of all edges $\{u, v\}$ with $v > 4$ and $|v - u| > 3$, and $E_D = E \setminus E_P$;
3. E_D can be further subdivided in E'_D and E''_D , so that E''_D consists of all edges $\{u, v\}$ with $|v - u| = 3$, and $E'_D = E_D \setminus E''_D$;
4. the distance function d is such that: (a) d_{uv} is a scalar for each $\{u, v\} \in E'_D$; (b) d_{uv} consists of a discrete set of b scalars for each $\{u, v\} \in E''_D$; (c) d_{uv} is a general interval for all $\{u, v\} \in E_P$.

We remark that the above definitions can be appropriately extended to Euclidean spaces of any dimension $K > 0$, not just $K = 3$. We call E_D the *discretization edges* and E_P the *pruning edges*. Discretization edges ensure that the

graph G is rigid, which implies that there are finitely many realizations of G in \mathbb{R}^K . Pruning edges make some of those realizations infeasible, and thereby make the solution set X smaller. A few remarks are in order:

- we consider that distances which are known because of covalent bond relations are sufficiently precise to be represented by a scalar;
- we consider that distances which are known from NOESY (or other) experiments can be represented by intervals;
- we assume that a limited number of the intervals can be discretized into sets containing a finite number b of values within the intervals;
- the edges in E'_D represent atom pairs of the form $\{v, v-1\}$ or $\{v, v-2\}$ for any $v > 2$: these are involved in covalent bonds;
- the edges in E''_D represent atom pairs which are assigned a certain number b of possible values (optionally $b = 1$ for certain pairs);
- the edges in E_P represent atom pairs for which the distance might be a general interval.

We remark that the order on V was initially intended to follow the protein backbone [49], but new orders which better exploit the hydrogen atoms in or close to the backbone have been defined in [50,51]: these are the orders on which the above assumptions are based.

The DGP with the restrictions above, but where all intervals are replaced by scalars, is called DISCRETIZABLE MOLECULAR DGP (DMDGP). Both the DMDGP and its generalization to any K (denoted by K DMDGP) are NP-hard [52,53]. The problem defined above, involving intervals, obviously contains the DMDGP as a sub-case and is hence also NP-hard by inclusion.

When all distances are precise

We first focus on the simplest case, where all intervals are replaced by scalar values. Then $d : E \rightarrow \mathbb{R}_+$, and $b = 1$. In this simplified setting, the *i*BP is simply called BP [52], and the order on V is called a *contiguous trilateration order* [54] or a *DMDGP order* [55].

The BP can be defined as a recursive procedure: assuming we already found a realization x_1, \dots, x_{v-1} for the vertices $1, \dots, v-1$, and that we mean to find a consistent realization x_v for v , the discretization edges E_D guarantee that there will be at most two positions for x_v compatible with the distances restricted to E_D [49]. This can be intuitively understood in \mathbb{R}^3 by considering the intersection of three spheres centered at $x_{v-1}, x_{v-2}, x_{v-3}$ with radii $d_{v,v-1}, d_{v,v-2}, d_{v,v-3}$: the first two spheres either do not meet or their intersection is in general a circle, and the intersection of the third sphere with this circle is either

empty or consists in general of two points [56]. We can now consider the distances defined on pruning edges in E_P , linking v to its preceding vertices in order to accept or reject these two points. For each accepted point we recursively call BP with v replaced by $v+1$, for all $v < n$. When $v = n$ we have a valid realization of the graph: we save it in X , and proceed to complete the recursive search. This yields a search tree which is explored depth-first. The recursion starts after placing the initial triangle 1,2,3 (either arbitrarily or by using BP restricted to subspaces), so this tree starts branching at level 4. It can be proved that, at completion, X contains all incongruent (modulo translations and rotations) realizations of G .

In the case where $E_P = \emptyset$, the search tree is a complete binary tree with 2^{n-3} nodes at the n -th (and last) level: in other words, its depth is n and its width is 2^{n-3} . This is the worst case, since the BP must explore all of the nodes in the tree, and proves that the BP (and hence the *i*BP, since it generalizes the BP) is an exponential-time algorithm in n .

When $E_P \neq \emptyset$, it was shown that X almost always contains a number of solutions which is either zero or a power of two [55]; this discovery led to a set of results where the BP search tree width can be kept polynomial in n during the search [53]. Since the exponential behavior is only due to the tree width, this yields a set of cases where the BP is actually fixed-parameter tractable (FPT). Throughout all our experiments with protein data we were always able to fix the parameter controlling the exponential growth of the tree width to a universal constant, which makes BP “polynomial on proteins” (this is an informal statement — the precise statement is given in [53]).

Intervals and discrete distance sets

The theory supporting the case where d might map edges to discrete sets of distance values or intervals, which is the case treated in this paper, is not so clearly understood yet. As it generalizes the simpler case sketched above, in a certain sense it inherits its properties, but this is an oversimplification: for instance, if all intervals are $[0, \infty]$, it is obvious that the problem is easy independently of the graph topology, since every realization is valid.

Some bounds on the cardinality of X in the presence of discrete sets and intervals are given in [55]. Our understanding is that if the intervals are small enough, the theory which led to fixed-parameter tractability goes through with few changes, but we have no way so far of establishing an aprioristic maximum width for the intervals. If the intervals are very large the problem might become tractable, as mentioned above, for the purposes of finding at least one solution. The *i*BP would still behave exponentially, however.

Results and discussion

We applied the presented algorithm to three examples of proteins displaying α helical secondary structures. Before presenting the obtained results, we emphasize that the method proposed here has a completely different philosophy than classical optimization approaches commonly used in the field of NMR structure determination. In the present approach, each constraint is treated in the strict sense, that is, no violation, however small, is tolerated. This is why we consistently use the word *constraint* in the paper. This is what potentially allows us to systematically explore the entire search space. However, the use of the procedure demands that the data have been pre-processed accordingly, and all geometric inconsistencies that exist in three-dimensional space have been removed.

For the proteins studied here, if one includes the ensemble of NMR interval distance constraints stored in the .mr file at the Protein Data Bank (PDB) [57] as well as all pruning devices described above, all solutions are pruned out, indicating that no solution to the distance geometry problem exists with the deposited data. This is not really surprising, since the optimization algorithms generally used in NMR structure determination are based on optimization of a target function or hybrid energy rather than on strict constraint satisfaction. That is, there is always a phase where the algorithm tries to find a trade-off when inconsistencies exist between constraints. The optimization thus produces solutions in which chemical and NMR constraints are optimized, but in which small violations are always present. These inconsistencies are present in any structure determination, in particular because distance constraints are imprecise, due to experimental limitations.

Since the data in the PDB for the examples presented here were not pre-processed the way our algorithm requires, we decided to use a subset of the stored data sets: the definition of α -helix regions and a few long-range distance constraints arbitrary selected from the set of NMR constraints for structures with more than one α -helix. In order to further reduce the risk of all solutions being pruned, we used tolerance values for atomic positions and angles between atoms (Table 2).

The three examples we chose to illustrate the algorithm display an increasing structural complexity: (i) a single α helix, corresponding to the structure of peptide CM15 determined in micelles (PDB id: 2JMY [58]), (ii) an α helical hairpin (PDB id: 2KXA [59]), (iii) the insecticidal toxin TAITX-1a, formed as a bundle of four α helices, restrained by three disulphide bridges (PDB id: 2KSL). The main characteristics of the studied proteins are given in Table 2. All three examples were originally determined by NMR spectroscopy, and the corresponding constraint lists are available from the PDB. The analysis by PROCHECK

[60] of the Ramachandran diagram of these three PDB structures shows that more than 85% of the residues are located in the core region. For 2KXA and 2KSL, more than 95% of the residues are located in the core and allowed region, whereas in 2JMY, 7% of the residues are located in the generously allowed region. For 2KXA, one PRO residue was replaced by an ALA, as the PRO cycle has not yet been included in the current version of the *iBP* algorithm.

We generated conformations using the branching phase and the pruning devices described above. The long-range constraints added for the calculations of 2KXA and 2KSL, are:

- (i) for 2KXA, one constraint between $H\alpha$ hydrogen and carbonyl oxygen of Ala-5 and Met-17, enforcing the pairing of the two α -helices,
- (ii) for 2KSL, three constraints between Carbons β of Cys-7 and Cys-37, of Cys-23 and Cys-33 and of Cys-26 and Cys-46, corresponding to the formation of the three disulphide bridges.

For all calculations, except the one of 2JMY with the α helix defined along the whole sequence, the obtained conformations were filtered according to the coordinate root mean-squared deviation (RMSD: 1.5 Å) with respect to the previously obtained conformation in the *iBP* procedure. Enforcing an RMSD value larger than 1.5 Å between two successively stored conformations, avoids an over-sampling of the conformational space. Each calculation was stopped after storing 10000 filtered conformations. For our three examples, five calculations were performed in total: three on 2JMY with different definitions of the α helix (residues 1-15, 3-13 and 5-11), and one each for 2KXA and 2KSL. For the first calculation on 2JMY, one conformation was obtained and saved. The second and third calculations on 2JMY were quite short, of the order of minutes (Table 2), which is due to the small size of the corresponding tree. For the 2KXA and 2KSL calculations, 10000 conformations were obtained in about 30 mins of calculation. Large total numbers of conformations were generated: this number increases from ~634,000 (2JMY_1) up to ~3,400,000 (2KXA) with the size of the considered problem, depending on the number of residues and on the number of constraints. Despite 2KSL being the largest example, the second smallest number of conformations was generated, which is the sign of a severe pruning arising from a rather restricted conformational space.

The reliability of the obtained conformations was checked in three ways. First, the whole set of NMR constraints deposited along with the PDB entries and involving backbone hydrogens, were probed on the conformations. Second, the quality of the obtained

Table 2 Analysis of conformations obtained by the branch-and-pruning algorithm on the three proteins targets: 2JMY, 2KXA and 2KSL

Proteins	2JMY	2JMY_1	2JMY_2	2KXA	2KSL
Number of residues	15	15	15	24	51
Number of vertices	107	107	107	170	359
Definition of α helices	1-15	3-13	5-11	1-11, 13-23	4-11, 13-27, 29-36, 41-50
Position tolerance (Å)	0.2	0.2	0.2	0.2	0.2
Angle tolerance (°)	2	2	2	4	4
<i>b</i> value	4	4	4	8	4
Number of long-range constraints	0	0	0	1	3
Number of saved conformations	1	10000	10000	10000	10000
Number of generated conformations	1	633,937	928,399	3,380,964	491,498
CPU time	-	1 min	1 min	25 min	31 min
Number of violated constraints (> 1 Å)	0	4.0 ± 2.1	11.6 ± 3.6	9.6 ± 2.9	12.8 ± 1.1
Maximum violation (Å)	0	3.3 ± 1.4	4.8 ± 0.7	3.7 ± 1.0	8.1 ± 0.6
Minimum RMSD from PDB structure (Å)	1.4	1.3	2.1	1.1	3.0
RMSD from PDB structure for minimum violated conformations (Å)	1.4	2.9	2.8	1.3	3.5
PROCHECK					
core residues (%)	100	65.7 ± 25.9	49.2 ± 7.6	60.4 ± 8.1	76.9 ± 2.4
allowed residues (%)	0	17.9 ± 9.7	40.9 ± 8.3	39.6 ± 8.0	21.3 ± 2.8
gen.allow. residues (%)	0	3.6 ± 4.8	9.9 ± 7.2	0.0 ± 0.0	1.9 ± 1.7
disall. residues (%)	0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0

2JMY_1 and 2JMY_2 correspond to the target 2JMY with shorter definitions of α helices. The total number of generated conformations is given, along with the number conformations filtered according to RMSD values.

conformations was checked using PROCHECK [60] analysis of the Ramachandran plot. Third, the obtained conformations were clustered with an unsupervised clustering method, namely the self-organizing map or SOM [61-64], in order to investigate the properties of sampled conformations.

The agreement of the obtained conformations with the backbone NMR constraints deposited with the PDB structures was checked by calculating the distances between the backbone hydrogens in each obtained conformation. The distances larger than the upper bound of the constraint correspond to violations of this constraint. The

mean number of violated constraints along with the mean value of the difference to the upper bound for these constraints were calculated on all conformations (Table 2). For the 2JMY calculation with the 1-15 α helix definition, no violation of the NMR constraints could be observed. As expected, when the α helix definition is reduced (2JMY_1 and 2JMY_2), the average number of violations increases as well as the average maximum violation. Not surprisingly, the most violated constraints involve residues located at the N and C terminal parts of the α -helix, TRP-2, PHE-5, LYS-3, LYS-6 and VAL-11, VAL-14, LEU-15 for 2JMY_1 and 2JMY_2. The largest violations and number of violations are of the same order or value for 2KXA than for 2JMY_1 and 2JMY_2. In contrast, the largest violations and number of violations are observed for 2KSL and involve residues CYS-33, GLU-34, PHE-38, TYR-43. Such over-restraining of NMR structures have been put in evidence in the past, through molecular dynamics simulations [65] and analysis of the structure quality [66].

The average number of violations is similar for 2JMY_2, 2KXA and 2KSL, but the average maximum violation for 2KSL is twice as large as that for 2JMY_2 and 2KXA. This might be due to the very restrained conformations of 2KSL, which contain three disulphide bridges. Due to this restrained conformation, the NMR constraint list is probably more prone to contain inconsistencies, and large mechanical strain can be stored in the structure if one uses an optimization procedure such as simulated annealing. In contrast, no mechanical strain whatsoever is generated by the *iBP* algorithm, and the obtained conformations might have a stronger tendency to deviate from the PDB conformations.

For each example, the obtained conformations were compared to the first conformation deposited in the PDB. Minimum RMSD values in the range 1.1-2.1 Å were obtained for all targets, except 2KSL for which the minimum RMSD value was 3.0 Å. Thus the Branch-and-Prune algorithm was able to capture conformations close to the PDB conformations, the larger value obtained for 2KSL arising from the larger mechanical strain quoted above.

For each calculation, the conformation displaying the smallest number of NMR constraint violations was compared to the first conformation deposited in the PDB. The RMSD values are smaller than 1.5 Å for 2JMY and 2KXA. This shows that, in the context of the *iBP* algorithm, the measured NMR constraints also push the structure toward the PDB structure. For 2JMY_1 and 2JMY_2, the RMSD value increases since the definition of the α helical region is shorter. For 2KSL, the conformation displaying the smallest number of constraint violations, displays an RMSD of 3.5 Å with the PDB first conformation, which agrees with the maximum number of violations observed for this protein and with the minimum RMSD with the PDB structure analyzed above.

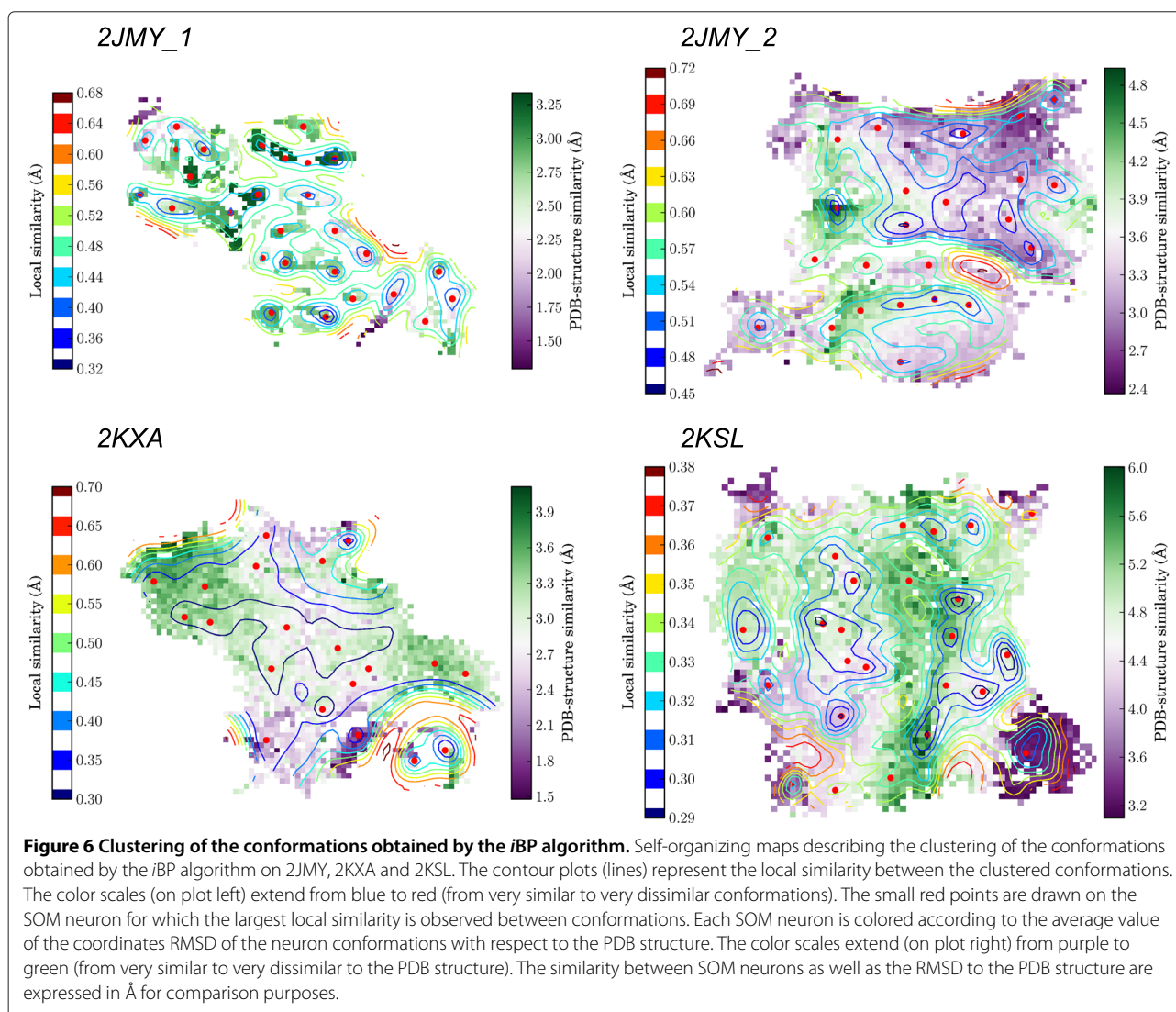
From the PROCHECK [60] analysis, the percentage of residues located in core and allowed Ramachandran regions, is larger than 95% for all targets except 2JMY_1, 2JMY_2, for which the percentages are about 80% due to the reduced definition of the α helix. For all targets, the percentage of residues in disallowed regions is equal to zero. The relatively important percentage of residues located in the allowed region may arise from the systematic exploration performed by the Branch-and-Prune algorithm, the strict nature of the constraints, and the nature of the pruning devices.

In order to further probe the robustness of the proposed algorithm, *iBP* calculations on 2KXA and 2KSL have been performed, using input data degraded in the following way: (i) the length of each α helix has been reduced by 1 residues at each extremity, (ii) the lower and upper bounds of the long-range distance constraints have been increased by 0.5 Å. The introduction of this noise into the α helical and long-range constraints makes the *iBP* solution moving apart from the PDB structure, as the minimum RMSD to PDB structure changes from 1.1 to 2 Å for 2KXA, and from 3.0 to 4.3 Å for 2KSL. Nevertheless, the quality of the Ramachandran diagram remains satisfying, with 93.3% and 95.4% of the residues located in the core and allowed regions of the Ramachandran plot for 2KXA and 2KSL.

The conformations were clustered using a self-organizing map (SOM) approach [62,63], on which the coordinate RMSD values between the conformers obtained by Branch-and-Prune and the corresponding PDB structure, were projected on the SOMs (Figure 6). These RMSD values lay in the 1.3-3.2 Å range for 2JMY_1, in the 2.4-4.9 Å range for 2JMY_2, in the 1.5-4.0 Å range for 2KXA, and in the 3.2-6.0 Å for 2KSL.

In the SOMs for the four calculations (Figure 6), the RMSD values are colored according to their RMSD from the PDB entry, violet color indicating values smaller than the median value of the sampled RMSD value, green color indicating RMSD values larger than this median value. For 2JMY_1, 2KXA and 2KSL, a larger number of neurons of the SOMs belongs to the second group, which is the sign of an enhanced sampling of the conformational space with respect to the region sampled by simulated annealing. For 2JMY_2, the inverse picture is observed, which may arise from the more limited conformational space available to be sampled for a unique α -helix.

In 2KSL and 2KXA SOMs, the protein conformations corresponding to the region displaying the smallest coordinate RMSD values with respect to the PDB structure, were extracted (Figure 7). These sets of conformers are similar to the superimposed conformations obtained in a usual NMR calculation.



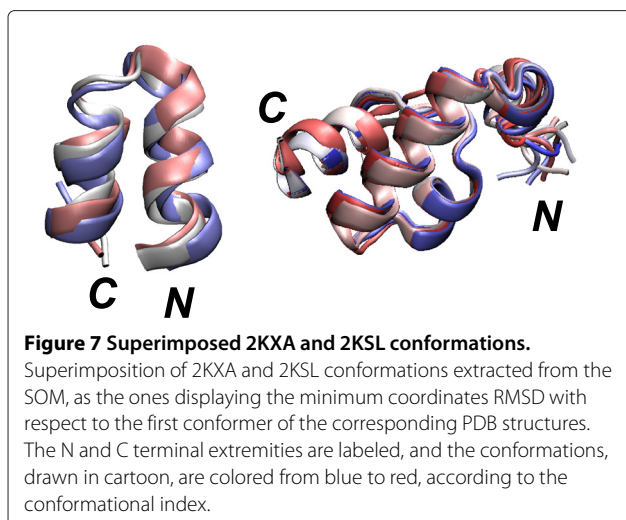
Conclusions

We proposed here a Branch-and-Prune algorithm (*iBP*) to solve the Distance Geometry Problem, in order to sample exhaustively the conformational space of the backbone of α -helical proteins. The *iBP* algorithm bears a very slight reminiscence to variable target function approaches for example implemented in DISMAN [67], due to the sequential nature of introducing constraints and non-bonded interactions. However, the precise way of introducing the constraints and non-bonded interactions differs significantly, and DISMAN does not systematically search space but is an optimization approach.

We introduced new pruning devices integrated in the *iBP* algorithm for DGP with intervals and we tested our *iBP* implementation on the backbones of α -helical proteins. Several pruning devices have been designed to enforce amino-acid chirality, α -helix geometry and van

der Waals steric hindrance. The algorithm allowed to efficiently reconstruct backbone conformations of three α -helical peptides, of various sizes, and for which the structure were previously solved by NMR. The obtained solutions satisfy most of the NMR constraints involving backbone hydrogen bonds, and display very acceptable Ramachandran statistics. The present work represents a first successful step on the way to reconstruct protein structures using a branch-and-prune algorithm applied to the Distance Geometry problem.

Applications where this approach could have significant advantages are cases where there are few distances defining the tertiary structure of a protein, where it is important to characterize the space of all solutions. It might also be useful as part iterative automated assignment algorithms such as ARIA [68], CYANA [69] or UNIO [70], where in a first iteration all solutions compatible with a



few unambiguous long-range constraints could be generated to reduce the ambiguity of the remaining constraints. Another application of the approach proposed here would be to provide input molecular conformations to model the structure of multi-subunit complexes into an electron microscopy density map [71].

Some limitations of the current version of *iBP* prevent for the moment its use with real nuclear Overhauser effect (NOE) data. These limitations are the use of unambiguous distance constraints, the non-inclusion of protein side-chains, the loss of information intervals and the appropriate weighting of the various constraints in order to overcome the inconsistencies contained among the whole constraint set. Protein side-chains can be added to the protein backbone afterward. The discretization of circle arcs could be tackled using algebraic geometry and geometric algebra approaches [72]. The Bayesian approach [73] developed for the objective weighting of various NMR constraints according to the data quality could be used to alleviate the inconsistency problems. The use of unambiguous distance constraints is probably the most unavoidable aspect of the current set-up of the algorithm.

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

AC, TM, MN, LL and AM designed the work. AC implemented the algorithm. TM, GB and BB performed and analyzed the application cases. AC, BB, AM, LL, CL, RA, MN and TM wrote the manuscript. All authors read and approved the final manuscript.

Acknowledgments

TM, MN, BB thank the Institut Pasteur and the CNRS for support. This work was funded by the European Union (FP7-IDEAS-ERC 294809 to MN), the "investissement d'avenir" program (grant bip:bip to MN and LL), and the Brazilian Research Agencies FAPESP and CNPq (to CL and RA).

Author details

¹Institut Pasteur, Structural Bioinformatics Unit, 25, rue du Dr Roux, 75015 Paris, France. ²CNRS UMR3528, 25, rue du Dr Roux, 75015 Paris, France. ³University of Campinas (IMECC-UNICAMP), 13083-859 Campinas-SP, Brasil. ⁴LIX, Ecole Polytechnique, 91128 Palaiseau, France. ⁵IBM TJ Watson Research Center, 10598 NY Yorktown Heights, USA. ⁶Université de Rennes-I, Rennes, France.

Received: 9 July 2014 Accepted: 5 January 2015

Published online: 28 January 2015

References

- Huang X, Britto M, Kear-Scott J, Boone C, Rocca J, Simmerling C, et al. The role of select subtype polymorphisms on HIV-1 protease conformational sampling and dynamics. *J Biol Chem* 2014;289:17203–14.
- Kanelis V, Forman-Kay J, Kay L. Multidimensional NMR methods for protein structure determination. *IUBMB Life* 2001;52:291–302.
- Sinz A. Chemical cross-linking and mass spectrometry for mapping three-dimensional structures of proteins and protein complexes. *J Mass Spectrometry* 2003;38:1225–37.
- Marti-Renom M, Stuart A, Fiser A, Sánchez R, Melo F, Sali A. Comparative protein structure modeling of genes and genomes. *Annu Rev Biophys Biomol Struct* 2000;29:291–325.
- Vajda S, Kozakov D. Convergence and combination of methods in protein-protein docking. *Curr Opin Struct Biol* 2009;19:164–70.
- Bello M, Martínez-Archundia M, Correa-Basurto J. Automated docking for novel drug discovery. *Expert Opin Drug Discovery* 2013;8:821–34.
- Crippen G, Havel T. Distance geometry and molecular conformation. New York: Wiley; 1988.
- Liberti L, Lavor C, Maculan N, Mucherino A. Euclidean distance geometry and applications. *SIAM Rev* 2014;56:3–69.
- Nilges M, Gronenborn A, Brünger A, Clore G. Determination of three-dimensional structures of proteins by simulated annealing with interproton distance restraints. Application to crambin, potato carboxypeptidase inhibitor and barley serine proteinase inhibitor 2. *Protein Eng* 1988;2:27–38.
- Alipanahi B, Krislock N, Ghodsi A, Wolkowicz H, Donaldson L, Li M. Determining protein structures from NOESY distance constraints by semidefinite programming. *J Comput Biol* 2013;20:296–310.
- Wang C, Lozano-Pérez T, Tidor B. AmbiPack: a systematic algorithm for packing of macromolecular structures with ambiguous distance constraints. *Proteins* 1998;32:26–42.
- Potluri S, Yan A, Chou J, Donald B, Bailey-Kellogg C. Structure determination of symmetric homo-oligomers by a complete search of symmetry configuration space using nmr restraints and Van Der Waals packing. *Proteins* 2006;65:203–19.
- Potluri S, Yan A, Donald B, Bailey-Kellogg C. A complete algorithm to resolve ambiguity for intersubunit NOE assignment in structure determination of symmetric homo-oligomers. *Protein Sci* 2007;16:69–81.
- Martin J, Yan A, Bailey-Kellogg C, Zhou P, Donald B. A geometric arrangement algorithm for structure determination of symmetric protein homo-oligomers from NOEs and RDCs. *J Comput Biol* 2011;18:1507–23.
- Martin J, Yan A, Bailey-Kellogg C, Zhou P, Donald B. A graphical method for analyzing distance restraints using residual dipolar couplings for structure determination of symmetric protein homo-oligomers. *Protein Sci* 2011;20:970–85.
- Reardon P, Sage H, Dennison S, Martin J, Donald B, Alam S, et al. Structure of an HIV-1-neutralizing antibody target, the lipid-bound gp41 envelope membrane proximal region trimer. *Proc Natl Acad Sci USA* 2014;111:1391–6.
- Zeng J, Boyles J, Tripathy C, Wang L, Yan A, Zhou P, et al. High-resolution protein structure determination starting with a global fold calculated from exact solutions to the rdc equations. *J Biomol NMR* 2009;45:265–81.
- Gordon D, Hom G, Mayo S, Pierce N. Exact rotamer optimization for protein design. *J Comput Chem* 2003;24:232–43.
- Kingsford C, Chazelle B, Singh M. Solving and analyzing side-chain positioning problems using linear and integer programming. *Bioinformatics* 2005;21:1028–36.
- Wang L, Donald B. An efficient and accurate algorithm for assigning nuclear Overhauser effect restraints using a rotamer library ensemble and residual dipolar couplings. The IEEE computational systems

- bioinformatics conference (CSB). Stanford, CA: The Institute of Electrical and Electronics Engineers, Inc.; 8-12 Aug 2005, pp. 189–202.
21. Wang L, Mettu R, Donald B. A polynomial-time algorithm for de novo protein backbone structure determination from NMR data. *J Comput Biol* 2006;13:1276–88.
 22. O'Neil R, Lilien R, Donald B, Stroud R, Anderson A. Phylogenetic classification of protozoa based on the structure of the linker domain in the bifunctional enzyme, dihydrofolate reductase-thymidylate synthase. *J Biol Chem* 2003;278:52980–7.
 23. Lavor C, Liberti L, Maculan N, Mucherino A. The discretizable molecular distance geometry problem. *Comput Optimization App* 2012;52:115–46.
 24. Lavor C, Liberti L, Mucherino A. The interval Branch-and-Prune Algorithm for the Discretizable Molecular Distance Geometry Problem with Inexact Distances. *J Global Optimization* 2013;56:855–71.
 25. Engh RA, Huber R. Accurate bond and angle parameters for x-ray protein structure refinement. *Acta Crystallogr Sect A: Found Crystallogr* 1991;47(4):392–400.
 26. Rocchia W, Alexov E, Honig B. Extending the applicability of the nonlinear poisson-boltzmann equation: Multiple dielectric constants and multivalent ions. *J Phys Chem B* 2001;105(28):6507–14.
 27. Honig B, Nicholls A. Classical electrostatics in biology and chemistry. *Science* 1995 May 26;268(5214):1144–9.
 28. Liberti L, Masson B, Lee J, Lavor C, Mucherino A. On the number of realizations of certain Henneberg graphs arising in protein conformation. *Discrete Appl Mathematics* 2014;165:213–32.
 29. Coope I. Reliable computation of the points of intersection of n spheres in \mathbb{R}^n . *ANZIAM Journal* 2000;42:461–77.
 30. Berg J, Tymoczko J, Stryer L. *Biochemistry: International Edition*. New York: WH Freeman & Co; 2006.
 31. Güntert P, Mumenthaler C, Wüthrich K. Torsion angle dynamics for NMR structure calculation with the new program DYANA. *J Mol Biol* 1997;273: 283–98.
 32. Güntert P, Wüthrich K. Sampling of conformation space in torsion angle dynamics calculations. *Comp Phys Commun* 2001;138:155–69.
 33. López-Méndez B, Güntert P. Automated protein structure determination from NMR spectra. *J Am Chem Soc* 2006;128:13112–22.
 34. Mucherino A, Lavor C, Malliavin T, Liberti L, Nilges M, Maculan N. Influence of pruning devices on the solution of molecular distance geometry problems. In: Pardalos, P., Rebennack, S. (eds.) *Lecture Notes in Computer Science* 6630. Germany: Springer; 2011. p. 206–17.
 35. Dong Q, Wu Z. A geometric build-up algorithm for solving the molecular distance geometry problem with sparse distance data. *J Global Optimization* 2003;26(3):321–33.
 36. Lavor C, Liberti L, Mucherino A, Maculan N. On a discretizable subclass of instances of the molecular distance geometry problem. In: *Proceedings of the 2009 ACM Symposium on Applied Computing*. ACM Press; 2009. p. 804–5.
 37. Floyd RW. Algorithm 97: shortest path. *Commun ACM* 1962;5(6):345.
 38. Shen Y, Delaglio F, Cornilescu G, Bax A. TALOS+: a hybrid method for predicting protein backbone torsion angles from NMR chemical shifts. *J Biomol NMR* 2009;44:213–23.
 39. Grishaev A, Bax A. An empirical backbone-backbone hydrogen-bonding potential in proteins and its applications to NMR structure refinement and validation. *J Am Chem Soc* 2004;126(23):7281–92.
 40. Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* 1983;22(12):2577–637.
 41. Abrahams D, Gurtovoy A. *C++ Template Metaprogramming: Concepts, Tools, and Techniques from Boost and Beyond*. Boston, Massachusetts: Addison-Wesley Professional; 2004.
 42. Austern MH. *Generic Programming and the STL: Using and Extending the C++ Standard Template Library*. Boston, Massachusetts: Addison-Wesley Longman Publishing Co., Inc.; 1998.
 43. Josuttis N. *The C++ Standard Library: a Tutorial and Reference*. Boston, Massachusetts: Addison-Wesley Professional; 1999.
 44. Anderson E, Bai Z, Bischof C, Blackford S, Demmel J, Dongarra J, et al. *LAPACK Users' Guide*, 3rd edn. Philadelphia, PA: Society for Industrial and Applied Mathematics; 1999.
 45. Lee L-Q, Lumsdaine A. *The Boost Graph Library: User Guide and Reference Manual*. Boston, Massachusetts: Addison-Wesley Professional; 2002.
 46. Brönnimann H, Melquiond G, Pion S. The design of the boost interval arithmetic library. *Theor Comput Sci* 2006;351(1):111–8.
 47. Brönnimann H, Melquiond G, Pion S. The boost interval arithmetic library. In: *Proceedings of the 5th Conference on Real Numbers and Computers*. Lyon, France; 2003. p. 65–80. <http://www.lri.fr/~melquion/doc/03-rnc5-article.ps.gz>.
 48. Saxe J. Embeddability of weighted graphs in k -space is strongly NP-hard. *Proc 17th Allerton Conference Commun Control Comput*. Monticello, Illinois; 1979:480–9.
 49. Liberti L, Lavor C, Maculan N. A branch-and-prune algorithm for the molecular distance geometry problem. *Int Trans Operational Res* 2008;15: 1–7.
 50. Lavor C, Mucherino A, Liberti L, Maculan N. On the computation of protein backbones by using artificial backbones of hydrogens. *J Global Optimization* 2011;50:329–44.
 51. Costa V, Mucherino A, Lavor C, Cassioli A, Carvalho L, Maculan N. Discretization orders for protein side chains. *J Global Optimization* 2014;60:333–49.
 52. Lavor C, Liberti L, Maculan N, Mucherino A. The discretizable molecular distance geometry problem. *Comput Optimization App* 2012; 52:115–46.
 53. Liberti L, Lavor C, Mucherino A. The discretizable molecular distance geometry problem seems easier on proteins. In: Mucherino A, Lavor C, Liberti L, Maculan N (eds.) *Distance Geometry: Theory, Methods, and Applications*. New York: Springer; 2013.
 54. Cassioli A, Günlük O, Lavor C, Liberti L. Discretization vertex orders for distance geometry. *Discrete Applied Mathematics* (accepted). in press.
 55. Liberti L, Masson B, Lavor C, Lee J, Mucherino A. On the number of realizations of certain Henneberg graphs arising in protein conformation. *Discrete Appl Mathematics* 2014;165:213–32.
 56. Lavor C, Lee J, Lee-St. John A, Liberti L, Mucherino A, Sviridenko M. Discretization orders for distance geometry problems. *Optimization Lett* 2012;6:783–96.
 57. Berman H, Kleywegt G, Nakamura H, Markley J. The future of the Protein Data Bank. *Biopolymers* 2013;99:218–22.
 58. Respondek M, Madl T, Göbl C, Golser R, Zangger K. Mapping the orientation of helices in micelle-bound peptides by paramagnetic relaxation waves. *J Am Chem Soc* 2007;129:5228–34.
 59. Lorieau J, Louis J, Bax A. The complete influenza hemagglutinin fusion domain adopts a tight helical hairpin arrangement at the lipid: water interface. *Proc Nat Acad Sci USA* 2010;107:11341–6.
 60. Laskowski R, MacArthur M, Moss D, Thornton J. PROCHECK: a program to check the stereochemical quality of protein structure. *J Appl Crystallogr* 1993;26:283–91.
 61. Bouvier G, Desdouts N, Ferber M, Blondel A, Nilges M. An automatic tool to analyze and cluster macromolecular conformations based on Self-Organizing Maps. *Bioinformatics* 2015. in press.
 62. Miri L, Bouvier G, Kettani A, Mikou A, Wakrim L, Nilges M, et al. Stabilization of the integrase-DNA complex by Mg²⁺ ions and prediction of key residues for binding HIV-1 integrase inhibitors. *Proteins* 2014;82: 466–78.
 63. Bouvier G, Duclert-Savatier N, Desdouts N, Meziane-Cherif D, Blondel A, Courvalin P, et al. Functional motions modulating VanA ligand binding unraveled by self-organizing maps. *J Chem Inf Model* 2014;54:289–301.
 64. Kohonen T. *Self-organizing Maps*. Heidelberg, Germany: Springer; 2001.
 65. Fan H, Mark A. Relative stability of protein structures determined by X-ray crystallography or NMR spectroscopy: a molecular dynamics simulation study. *Proteins* 2003;53:111–20.
 66. Nabuurs S, Spronk C, Vuister G, Vriend G. Traditional biomolecular structure determination by NMR spectroscopy allows for major errors. *PLoS Computational Biol* 2006;2:9.
 67. Braun W, Gö N. Calculation of Protein Conformations by Proton-Proton Distance Constraints: A New Efficient Algorithm. *J Mol Biol* 1985;186: 611–26.
 68. Rieping W, Habeck M, Bardiaux B, Bernard A, Malliavin T, Nilges M. ARIA2: automated NOE assignment and data integration in NMR structure calculation. *Bioinformatics* 2007;23:381–2.
 69. Güntert P. Automated NMR structure calculation with CYANA. *Methods Mol Biol* 2004;278:353–78.

70. Guerry P, Herrmann T. Comprehensive automation for NMR structure determination of proteins. *Methods Mol Biol* 2012;831:429–51.
71. Lasker K, Sali A, Wolfson H. Determining macromolecular assembly structures by molecular docking and fitting into a electron density map. *Proteins* 2010;78:3205–11.
72. Lavor C, Alves R, Figueiredo W, Petraglia A, Maculan N. Clifford Algebra and the discretizable molecular distance geometry problem. *Adv Appl Clifford Algebras* 2015. in press.
73. Bernard A, Vranken W, Bardiaux B, Nilges M, Malliavin T. Bayesian estimation of NMR restraint potential and weight: a validation on a representative set of protein structures. *Proteins* 2008;79:1525–37.

**Submit your next manuscript to BioMed Central
and take full advantage of:**

- **Convenient online submission**
- **Thorough peer review**
- **No space constraints or color figure charges**
- **Immediate publication on acceptance**
- **Inclusion in PubMed, CAS, Scopus and Google Scholar**
- **Research which is freely available for redistribution**

Submit your manuscript at
www.biomedcentral.com/submit

