

Document downloaded from:

<https://riunet.upv.es/handle/10251/233320>

This paper must be cited as:

Inostroza-Ponta, M.; Dorn, M.; Escobar, I.; De Lima Correa, L.; Rosas-Olivos, Erika Susana; Hidalgo, N.; Marin, M. (2020). Exploring the high selectivity of 3-D protein structures using distributed memetic algorithms. *Journal of Computational Science*. 41. <https://doi.org/10.1016/j.jocs.2020.101087>



The final publication is available at

<https://doi.org/10.1016/j.jocs.2020.101087>

Copyright Elsevier

Additional Information

# Exploring the high selectivity of 3-D protein structures using distributed memetic algorithms

Mario Inostroza-Ponta<sup>a,\*</sup>, Márcio Dorn<sup>d</sup>, Iván Escobar<sup>a</sup>, Leonardo de Lima Correa<sup>d</sup>, Erika Rosas<sup>c</sup>, Nicolás Hidalgo<sup>b</sup>, Mauricio Marin<sup>a</sup>

<sup>a</sup>*Universidad de Santiago de Chile, Estación Central, Santiago, Chile.*

<sup>b</sup>*Universidad Diego Portales*

<sup>c</sup>*Universidad Técnica Federico Santa María*

<sup>d</sup>*Institute of Informatics, Federal University of Rio Grande do Sul  
Av. Bento Gonçalves 9500, 91501-970, Porto Alegre, RS, Brazil.*

---

## Abstract

This paper addresses the problem of predicting the tertiary structure of a protein given its amino acid sequence, which has been reported to belong to the NP-Complete class of problems. We design an ad-hoc distributed memetic algorithm (DMA) and evaluate several algorithm configurations in terms of different distributed population structures, ad-hoc local search strategies and the combination of two energy functions. The algorithm uses an asynchronous hierarchical population of agents that exchange solutions along the execution of the algorithm. Extensive computational experiments were carried out in order to test: (1) the impact of the communication on different population structures, (2) the combination of the energy functions used for fitness calculations, (3) the scalability of the algorithm for structures with a larger number of agents, (4) the performance of the different approaches proposed for local search and diversity calculations, (5) the biological significance of the predicted structures and (6) to compare the best performing configuration of the DMA with other algorithms from the literature. The algorithm was tested on 20 sequences of different size, and the analysis was performed regarding both computational quality and biological significance of the predicted structures. Results show that the combination of energy functions and the proposed Distributed Memetic Algorithm allows the prediction of structures that are similar to the experimental ones. Performance analysis shows that increasing parallelism improves the execution times, without worsening the quality of the solutions.

*Keywords:* protein structure prediction, memetic algorithm, distributed memetic algorithms, structural bioinformatics.

---

\*Corresponding author

*Email address:* [mario.inostroza@usach.cl](mailto:mario.inostroza@usach.cl) (Mario Inostroza-Ponta)

## 1. Introduction

It is well known that proteins are present in all living systems, and are responsible for performing a large set of functions in the cell. Knowing the structure of a protein allows the investigation of the biological processes in which the protein could be involved more directly, with higher resolution and finer detail. The sequence-protein-structure paradigm (also known as the “lock-and-key” hypothesis) says that the protein can achieve its biological function only by reaching a given tertiary structure determined by its amino acid sequence. The Protein Structure Prediction (PSP) problem is one of the most important problems in Structural Bioinformatics. Each protein, or polypeptide, is defined by an unique sequence of chained amino acids that under some physiological conditions reach a particular 3-D shape. *In silico* prediction of the three-dimensional (3-D) structure of proteins has been the ‘*holy grail*’ of computational biologists for many years [41].

The PSP problem has challenged Biochemists, Biologists, Computer Scientists, and Mathematicians, comprising the domain of high-dimensional continuous optimization problems [9]. Nowadays, a wide range of optimization algorithms and metaheuristics are being applied in an attempt to find approximate solutions to predict the 3-D structure of protein only from its linear sequence of amino acid residues [10]. Despite the number and the variety of methods developed for structure prediction over the past 40 years, no truly accurate methods exist to predict the correct structure of polypeptides. Most of the existing challenging optimization problems, such as the PSP, cannot be optimally solved by any known computational method due to the high dimensionality and complexity of the search space [1].

In this paper we propose a knowledge-based *Distributed Memetic Algorithm* (DMA) [30] to deal with the 3D PSP problem. The proposed algorithm takes advantage of the experimental knowledge stored in the *Protein Data Bank* (PDB). The main goal is to decrease the PSP complexity, reducing the size of the search space and increasing the method effectiveness. To incorporate the structural information of protein conformations, we use the *Angle Probability List* (APL) [5, 4]. APL analyses the conformational preferences of amino acids in proteins according to its secondary structure. The DMA proposed in this work also uses a *Local Search* (LS) technique to intensify the search around the most favourable regions.

Parallelism is intrinsic for memetic algorithms. Harris et al. [19] and Blocho et al. [3] use memetic algorithms with a parallel local search. The first solves the quadratic assignment problem using a ternary tree of 13 agents. The second solves the vehicle routing problem with time windows using a ring structure of  $N$  agents. This approach significantly improves execution times. However, the synchronization among the agents makes it difficult to scale over 32 agents [3]. Our proposal also explores different structures for cooperation between agents. We used tree-based structures to calibrate knowledge sharing and analyze the scalability of the approach.

This paper is organized as follows. Section 2 presents a brief review about

46 the problem of PSP and the current techniques that exist to use distributed  
47 metaheuristics. Section 3 describes in detail the proposed method and their  
48 variants. In Section 4 we present the computational experiments and discussion  
49 of results. Finally, the last section present the conclusions and points out some  
50 future works.

## 51 2. Background

### 52 2.1. The 3D Protein Structure Prediction Problem

53 The prediction of protein structure is one of the most important and chal-  
54 lenging problems in structural bioinformatics. Traditionally, protein's structure  
55 can be obtained through X-ray crystallography and multi-dimensional magnetic  
56 resonance in the laboratory. Although this method is expensive, it provides an  
57 accurate protein structure. However, the high costs of the method, its extreme  
58 difficulty, and the time required to carry out the study cannot cope with the  
59 increasing number of proteins data sequences. In this scenario, computational  
60 methods are a more suitable solution to deal with this problem due to its low  
61 costs and fast results [27]. These methods have the potential to correlate and  
62 predict the structure of a protein from its primary sequence to overcome the  
63 difficulties associated with the experimental approaches [21].

64 In an attempt to solve this problem many methods and algorithms have been  
65 proposed and they can be classify over four categories: (a) methods without  
66 database information; (b) methods with database information; (c) fold recogni-  
67 tion methods, and (d) comparative modelling methods.

68 Methods without database information are founded on thermodynamics and  
69 based on the fact that the structure of a protein corresponds to the global  
70 minimum energy conformation. Thus, a stable structure of a protein always  
71 consumes minimum energy among all the conformations of the protein. These  
72 methods are supported by two main components: a mathematical model for the  
73 protein-energy function and an optimization algorithm to find a global minimum  
74 of the potential energy [42] [7] [16] [15] [35].

75 The methods with database information follow a set of general rules ex-  
76 tracted from protein databases to establish a starting protein structure. These  
77 methods compare fragments of the protein with fragments of the structure of  
78 reference proteins [49] [37] [45] [11] [28] [44] [8].

79 Fold recognition methods assume that the structure is more likely to be  
80 preserved in the protein evolution rather than the sequence. i.e., proteins with  
81 no apparent sequence similarity could present similar folds. These methods'  
82 goal is to predict the protein structure by fitting a protein sequence correctly  
83 against a structural model [2] [22] [43] [17].

84 Finally, the comparative modelling employs a target sequence of amino acid  
85 is aligned against the amino acid sequence of another protein with known struc-  
86 ture (template protein) and stored in the protein data bank (pdb) [29] [38] [13] [43].

87 Predicting the 3D structure of a protein molecule *in silico* is a NP-hard  
88 problem. Metaheuristic optimization algorithms have potential ability to solve

89 complex problems efficiently providing near optimal solutions. Our proposal  
90 follows this optimization strategy as well as many of the before cited solutions.

## 91 2.2. Distributed Memetic Algorithms

92 Memetic algorithms have been distributed in literature given their collabo-  
93 rative nature and the complexity of the real-world problems they are used to  
94 solve. In Gong et al. [18] the authors review the models of the distributed  
95 evolutionary algorithms. Different authors have proposed different memetic al-  
96 gorithms that use distributed models that include the well-know island model  
97 [14] [33] [39], the cellular model [26] [34], and the hierarchical architecture [19].  
98 In this work we study the hierarchical model of organisation among the agents,  
99 which has been previously shown good results for our particular problem in [8].

100 The different cooperation topology inside the components of each model has  
101 influence on the algorithm search capabilities in terms of increasing exploration  
102 or exploitation of the search space. Cooperation is to define: what is to be  
103 transferred, what is the cooperation topology, how often the processes commu-  
104 nicate, and what is the strategy to handle immigrant and emigrants. Nalepa  
105 et al. [33] [32] studied the impact of this cooperation for the VRPTW and the  
106 PDTW problem. The authors conclude that if the convergence time is to be  
107 as minimum as possible, frequent cooperation intervals are better. Moreover,  
108 if the running time is a restriction, then rarer co-operation intervals benefit a  
109 broader exploration of the search space. Our proposal also explores different  
110 structures for cooperation between agents. We used tree-based structures in order  
111 to calibrate knowledge sharing and analyse the scalability of the approach.

## 112 3. DMA-3DPSP: A Distributed Memetic Algorithm for the PSP prob- 113 lem

114 We propose a Distributed Memetic Algorithm to tackle the 3-D Protein  
115 Structure Prediction problem (DMA- 3DPSP). The distributed model follows a  
116 hierarchical population structure based on the *Island Model* [18], with a set of  
117 agents working independently in a particular search space and cooperating by  
118 exchanging information according to a hierarchical population structure. The  
119 hierarchical structure used will be determined from the computational exper-  
120 iments. In the DMA-3DPSP, agents share the information about local solutions  
121 found and preserve the most promissory ones. The goal of this distributed strat-  
122 egy is to improve both the execution times of the search process and the quality  
123 of the solutions obtained.

124 Algorithm 1 shows the pseudocode for each agent of the DMA-3DPSP. Agents  
125 are organised under a hierarchical structure or *n-ary* tree. The structure used  
126 defines the concept of sub-populations composed by a *leader agent* and *sup-*  
127 *porter agents*. *Leaders* are responsible for receiving the best solutions generated  
128 by their *supporters* and also to spread good quality solutions among their *sup-*  
129 *porters*. Information updates are performed asynchronously each time an agent  
130 finds a better solution. Each agent keeps their best solutions on a local memory

131 called *pockets* (line 5) that keeps track of the best solutions found so far by  
 132 the agent. In each generation, the agent selects one of its pocket solutions and  
 133 performs a *crossover* with the solution received from its leader. Then, on the  
 134 *offspring* generated it applies a local search algorithm to improve the quality of  
 135 the solution. After that, if a new best solution is found, the agent sends this  
 136 solution to its leader, and it starts a new generation. In the following sections  
 137 we described the main components of the proposed algorithm.

### 138 3.1. Agent Structure

139 Figure 1 shows a scheme of the main components of the agent. Each agent  
 140 has two structures to store and to manage the progress of the search: (1) the  
 141 *current*, that is a temporary solution on which evolutionary operators are being  
 142 applied, and (2) a set of *pockets*, a local memory structure that stores the best  
 143 solutions found by the agent and received by its supporters. The pockets enable  
 144 the agent to perform a better exploration of the search space by focusing on the  
 145 sectors where these solutions are located; their number is variable, and it is a  
 146 parameter of the algorithm.

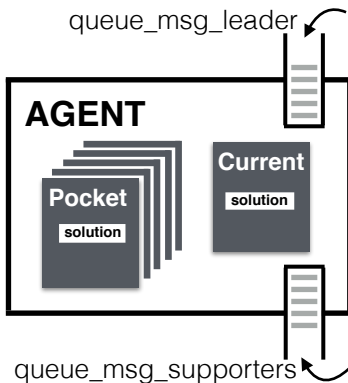


Figure 1: Structure of each agent of DMA-3DPSP. Agents work in the *current* solution, save the best solutions found in the *pockets* and send them to *leader* and supporters through message queues.

147 Each *solution* (*current* and *pockets*) describes a possible 3-D structure for  
 148 the peptide and is composed by: (1) a vector of size  $2 \times n$  (for a sequence of  $n$   
 149 amino acids) with the value of the  $\phi$  and  $\psi$  torsion angles for each residue of the  
 150 sequences and (2) the free energy of the structure given by the energy fitness  
 151 function. For communication, the agents use a distributed message queuing  
 152 system. Each agent has two queues, one dedicated to receiving the solutions  
 153 from its leader and the other to receive them for its supporters (sub-agents).  
 154 This asynchronous model allows the agents to send solutions to other agents  
 155 without interrupting the search process.

---

**Algorithm 1** Distributed Memetic Algorithm for the PSP problem - DMA-3DPSP

---

```
1: Input: seq  $\leftarrow$  Aminoacid sequence, ss  $\leftarrow$  Secondary structure, histograms  
    $\leftarrow$  APL Histograms  
2: procedure DMA-3DPSP  
3:   for all Agent do                                      $\triangleright$  Initialise agent population  
4:     pockets  $\leftarrow$   $\emptyset$   
5:     pockets_leader  $\leftarrow$  None  
6:     APL  $\leftarrow$  LoadHistograms(seq, histograms)  
7:     current  $\leftarrow$  initSolution(seq, APL)  
8:     updatePockets(current)                                $\triangleright$  Perform search  
9:     while stop_criteria has not been reached do  
10:      if Agent is root then  
11:        current  $\leftarrow$  GetRandom(pockets)  
12:      else  
13:        current  $\leftarrow$  Crossover(ss, GetRandom(pockets),  
   Roulette(pockets_leader))  
14:      end if  
15:      current  $\leftarrow$  SimulatedAnnealing(current, seq, APL)  
16:      updatePockets(current)                                $\triangleright$  Perform communication  
17:      while queue_supporters  $\neq \emptyset$  do  
18:        updatePockets(queue_supporters.pop())  
19:      end while  
20:      while queue_leader  $\neq \emptyset$  do  
21:        pockets_leader  $\leftarrow$  queue_leader.pop()  
22:      end while  
23:      if pockets were updated then  
24:        for each Supporter do  
25:          send pockets to Supporter  
26:        end for  
27:        if pockets[0] was updated then  
28:          send pockets[0] to Leader  
29:        end if  
30:      end if  
31:      restartControl()  
32:    end while                                              $\triangleright$  End search process  
33:    if Agent is root then  
34:      return pockets[0]  
35:    end if  
36:  end for  
37: end procedure
```

---

156 *3.2. Population Structure and Communication*

157 Agents in our approach are organised in a *n-ary* tree, where the father  
158 of an agent is called its *leader* and the children of the agent are part of its  
159 *supporters*. Each group of a leader and its supporters defines a subpopulation  
160 of the system. *Leaders* keep in their *pockets* the best solutions found in the  
161 subpopulation, and feed the search process of their *supporters* with these good  
162 solutions. Supporters, on the other hand, are responsible for sending their  
163 best solutions to their leaders as soon as they found a new one. Agents work  
164 independently and communicate with them only to share solutions and major  
165 events. In order to avoid interruption and the loss of solutions, a message  
166 queueing system is deployed at each agent. It enables the agent to send solutions  
167 without interrupt the receiver’s search process nor wait for an acknowledgment.

168 We evaluate three population structures: (1) a binary tree, (2) a ternary tree,  
169 and (3) a 12-ary tree, all of them composed by 13 agents. In the 12-ary tree, the  
170 structure has one level, following a master-slave approach where twelve of the  
171 thirteen agents have the same leader agent. The ternary tree builds a structure  
172 with two levels and the binary tree a three level structure. A structure with more  
173 agents should reach a better exploration of the search space and therefore obtain  
174 better results regarding processing time and solution quality, however, it also  
175 impacts communication. Communication used for sharing solutions between  
176 the agents is bi-directional: each time an agent has found a new best solution  
177 it sends the solution to its leader and also to its supporters, then a structure  
178 composed with a larger number of agents requires a more significant amount of  
179 communication. It is necessary to explore the balance between the convergence  
180 time and the exploration of the search space to reach high-quality solutions.  
181 Section 4.1 presents a complete description of each one of these structures.

182 *3.3. Fitness Function*

183 In order to evaluate the quality of the predicted structures and to guide  
184 the search process of the algorithm, a fitness function is required. Theoretically,  
185 the native structure of a protein corresponds to the global minimum of  
186 its free energy [6]. We propose a fitness function composed by two different energy  
187 approaches: (1) the *Rosetta all-atom high-resolution* energy function, using  
188 “*Talaris*” weights for the energy terms, and (2) the Solvent Accessible Surface  
189 Area calculation (SASA) [36] with an atomic radius of 1.5Å, both of them implemented  
190 in the *PyRosetta toolkit*. The former seeks to improve the stability of the  
191 structure, looking for solutions with low values of energy among all interactions  
192 and that fulfil the restrictions of atoms clashes, and the latter aims to improve  
193 the packing of the 3-D structure by reducing the area exposed to solvent.

For the combination of the energy functions terms *Talaris* and *SASA*, we used the fitness function described in equation 1:

$$fitness(S) = (1 - \alpha) \times Talaris + \alpha \times SASA \quad (1)$$

194 where  $\alpha$  determines the weight of the energy functions. In the computational  
195 experiments we will show the results of testing the fitness function using

196  $\alpha \in \{0.0; 0.1; \dots; 0.9; 1.0\}$ , where an  $\alpha = 0.2$  means that the fitness function is  
197 composed by 80% of Talaris and 20% of *SASA*.

### 198 3.4. Crossover Operator

199 The DMA-3DPSP combines population operators and a *Simulated Annealing*  
200 local search strategy. We used the same operators as those used in [20], but we  
201 adapt them to work asynchronously. At the beginning of each generation the  
202 algorithm generates new *offspring* using a *crossover* operator. The *crossover*  
203 process first selects two solutions as parents, one of them randomly selected from  
204 the pockets of the agent, and the other one selected from its leader’s pockets  
205 through a roulette wheel selection, so the best solutions are more likely to be  
206 selected. The new *offspring* is generated by taking the angles of each residue  
207 of the sequence from the residues of the parents, with a probability of 0.6 for  
208 the parent taken of the agent’s pockets and 0.4 for the one selected from its  
209 leader’s pockets. In order to preserve the formation of secondary structures,  
210 we implement a *crossover* operator, where the amino acid chain is separated  
211 into segments based on its secondary structure (received as a parameter of the  
212 algorithm). Sequential amino acids that share the same secondary structure  
213 will be treated as a whole section for effects of *crossover*, and the angles for  
214 the residues of the section will be taken from the same parent. The only agent  
215 that does not perform crossover is the leader of the population or root agent.  
216 Instead, this agent only replaces its current with one of its pocket solutions  
217 randomly chosen.

### 218 3.5. Local Search Algorithm

219 The new solution generated in the *crossover* will go through the local search  
220 procedure. For each residue of the sequence, the algorithm decides whether it  
221 goes to local search with a probability of 0.9. A residue selected for local search  
222 first determines whether to perform a jump to another sector in the search space  
223 based on the preferences of the *APL* [5, 4]. Then, the jump is performed with  
224 a probability of 0.3, and it will add a value  $< jump$ , which will be decreased in  
225 each generation at a rate of 0.85. After that, a *Simulated Annealing* algorithm  
226 is performed on the angles  $\phi$  and  $\psi$  of the residue. For each angle, the algorithm  
227 makes small modifications such that it is always within the same angle cell.  
228 Each change will be followed by an evaluation of the globally free energy of  
229 the structure. Then, the solution is accepted if it improves the fitness value,  
230 otherwise the algorithm applies the Boltzmann probability function to decide if  
231 the new solution is accepted. The temperature function controls the acceptance  
232 function and the number of iteration of the *Simulated Annealing* algorithm:

$$T(iteration) = T_{initial} \times \alpha^{iteration} \quad (2)$$

233 where  $T_{initial}$  is the initial temperature and  $\alpha$  is the cooling factor. We use  
234  $T_{initial} = 4000$  for the root agent and  $T_{initial} = 2000$  for the rest of the agents,  
235  $\alpha = 0.9$  and a limit of  $T(iteration) < 0.1$ , using the results of the parametrisation  
236 tests performed in [20].

237 *3.6. Pockets Updating Policy*

238 After memetic operators are applied, the agent evaluates whether it should  
 239 update its *pockets* with the new solution. To determine if a new solution is  
 240 worth to keep in the *pockets*, the agent compares its free energy (through the  
 241  $E(sol)$  function) and its diversity (through the  $Div(sol_a, sol_b)$  function) with  
 242 each one of its *pockets*. To be saved in the *pockets*, the new solution must fulfill  
 243 at least one of the following conditions:

- 244 •  $\exists j / (Div(current, pocket_j) < Div_{min}) \wedge (E(current) < E(pocket_j))$   
 245 So *current* will replace  $pocket_i$  such that  $Div(current, pocket_i)$  is mini-  
 246 mum,  $\forall i$  where  $E(current) < E(pocket_i)$ .
- 247 •  $(Div(current, pocket_j) \geq Div_{min}, \forall j) \wedge (pockets \neq full)$   
 248 So *current* will replace the first empty pocket.
- 249 •  $(Div(current, pocket_j) \geq Div_{min}, \forall j) \wedge (pockets = full) \wedge (\exists k \text{ such that}$   
 250  $E(current) < E(pocket_k))$   
 251 So *current* will replace  $pocket_i$  such that  $E(pocket_i)$  is maximum,  $\forall i$ .

252 The agent also uses these rules to determine whether it should save the  
 253 solutions put on its queue from its supporters. Then, if there was an update  
 254 in the pockets, the agent proceeds to the communication process. The best  
 255 solution is sent to the agent's leader and the list of *pockets* is forwarded to each  
 256 supporter to perform the crossovers of the next generations. Therefore, good  
 257 solutions go towards the top of the hierarchical structure.

258 *3.7. Local Search Approaches*

259 We evaluate a set of modifications on the original Simulating Annealing  
 260 local search implementation in order to improve the quality of its results. In  
 261 the following we give a detail description of each one of these approaches:

- 262 • **LS1:** in the original local search algorithm, before the execution of the  
 263 Simulating Annealing on the angle cell, the residues could perform a  
 264 change of neighborhood. Then, with a probability of 0.3, the angles  $\phi$   
 265 and  $\psi$  of the residue will relocate to a different position in the search  
 266 space, based on the preferences for these angles on the APL. This was  
 267 meant to improve the exploration of the search space, but there is the  
 268 risk of losing a good solution when the changed is performed. We add an  
 269 additional restriction: the residue will only maintain the new position if  
 270 it generates an improvement in the fitness value.
- 271 • **LS2:** the original local search performs the Simulating Annealing on the  
 272 dihedral angles independently. Therefore, the algorithm first applies small  
 273 modifications in the  $\phi$  angle while  $\psi$  remains the same, and then it does  
 274 the same for  $\psi$  with  $\phi$  remaining static. This implementation limits the  
 275 exploration around the angle cell. In order to avoid the above, this local  
 276 search modifies both angles in each iteration of the Simulated Annealing.

- 277 • **LS3**: in the original local search algorithm, the angle cell that limits the  
278 range of the small random modifications is determined at the start of the  
279 Simulated Annealing, and it remains static throughout the iterations of  
280 the algorithm. The angle cell with the size of one degree is determined  
281 by taking the two integer values that surround each dihedral angle as  
282 the limit, and every modification must be kept in this area. We modify  
283 this approach by implementing a dynamic angle cell that is recalculated  
284 every time the algorithm applies an adjustment in the dihedral angles,  
285 allowing the algorithm to focus the search towards the sectors that show  
286 improvements in the results. The new cell will be the area delimited by a  
287 range of  $\pm 0.5^\circ$  from the new  $\phi$  and  $\psi$  angles.
- 288 • **LS4**: this local search has two phases: (1) during the first iterations, the  
289 algorithm uses the **APL** to determine the new positions for the  $\phi$  and  $\psi$   
290 angles. The size of the modification is delimited by a range not greater  
291 than *jump*, defined before, and it is only performed if the new position  
292 improves the global free energy or it is accepted by the acceptance function  
293 of the Simulated Annealing. When *jump* value has become less than one  
294 degree, (2) for the rest of the iterations the algorithm will use the original  
295 approach applying small random modifications in the angle cell. This local  
296 search algorithm allows searching good sections of the search space in the  
297 first iterations and focus the search around them in the last iterations.

### 298 3.8. Diversity Approaches

299 We implement two different techniques to determine the diversity among the  
300 solutions in the pockets of the agents:

- **Div1**: the first approach calculates the average difference among the di-  
hedral angles from the residues of the two solutions, according to the  
equation [3](#)

$$Div(a, b) = \frac{\sum_{i=1}^n (\phi_i^a - \phi_i^b) + \sum_{i=1}^n (\psi_i^a - \psi_i^b)}{n} \quad (3)$$

301 where  $\phi_i^a$  and  $\psi_i^a$  represent the dihedral angles of the residue  $i$  in the  
302 structure  $a$ , and  $\phi_i^b$  and  $\psi_i^b$  do the same for the structure  $b$ .

- **Div2**: the second approach uses the calculation of the RMSD (CA atoms)  
of the structures as a measure of the difference between the solutions.  
The RMSD is calculated by superimposing the atoms of the structures and  
measuring their average distance, as it is shown in the equation [4](#)

$$RMSD(a, b) = \sqrt{\frac{1}{n} \sum_{i=1}^n \|a_i - b_i\|^2} \quad (4)$$

303 where  $a_i$  and  $b_i$  are the vector of positions for the atom  $i$  of the structures  
304  $a$  and  $b$  when they are optimally superimposed.

305 *3.9. Restart Procedure*

The search space is continuous, which makes easy for the algorithm usually to fall on local optima, so the algorithms implements a restart procedure to escape from these situations. If during a  $g_{limit}$  number of generations the best solution of the leader agent of the population does not change, the algorithm restarts the whole population. To do it so, it deletes the solutions stored in the pockets of every agent, keeping only the best solution found so far by the entire population. The value of  $g_{limit}$  is a parameter of the algorithm, but it dynamically changes during the execution according to the number of generations that took to reach the last best solution. This gives time to the algorithm to converge in advanced stages of the evolution process, where usually is more difficult to reach better solutions. The  $g_{limit}$  for the  $n$  restart is calculated following the equation:

$$f(n) = \begin{cases} g_{min} & n = 1 \\ g_{min} + (g_{current} - g_{last}) - f(n - 1) & n > 1 \end{cases} \quad (5)$$

306 where  $g_{min}$  is given as parameter for the first  $g_{limit}$ ,  $g_{current}$  is the current  
 307 generation in which the algorithm is working and  $g_{last}$  is the generation in which  
 308 the last restart was made. The root agent is in charge of the restart procedure,  
 309 evaluating the number of generations since the last time its best pocket was  
 310 updated, and communicating and synchronizing the rest of agents when it is  
 necessary to make a restart.

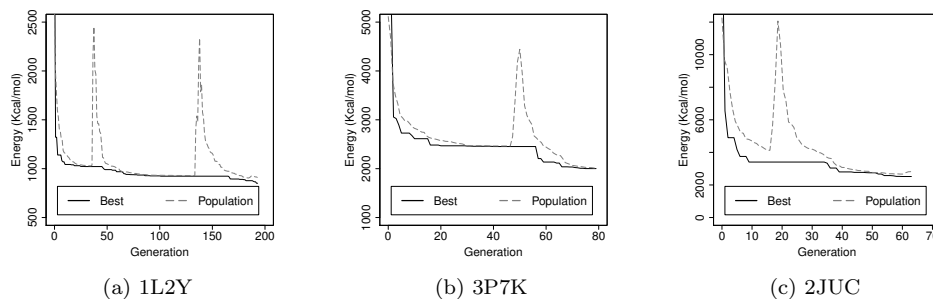


Figure 2: Analysis of energy convergence of the agent pockets of the DMA-3DPSP on three different protein sequences. The solid line shows the energy value of the best solution found so far by the entire population, and the dashed line shows the average energy of the best solutions stored in each agent of the population.

311 In Figure 2 we show the convergence of the free energy on the DMA-3DPSP for  
 312 three different sequences. When a restart is performed, the energy of the pockets  
 313 of the population (dashed line) increases rapidly, because the relocation of the  
 314 solutions in the search space, but after these solution converge it is possible to  
 315 see how the relocation allows the algorithm to escape from local minima.  
 316

317 **4. Computational Experiments**

318 We performed several computational experiments to evaluate the perfor-  
 319 mance of the alternative configurations of the proposed DMA. The tests aim to  
 320 evaluate: (1) the impact of the communication on different population struc-  
 321 tures, (2) the combination of the energy functions used for fitness calculations,  
 322 (3) the scalability of the algorithm for structures with a greater number of  
 323 agents, (4) the performance of the different approaches proposed for local search  
 324 and diversity calculations, (5) the biological significance of the predicted struc-  
 325 tures and (6) compare the best performing configuration of the DMA with other  
 326 algorithms from the literature.

327 We select 20 sequences from known peptide structures as cases of study  
 328 for our experiments. Table 1 shows a list of the chosen peptides, its number  
 329 of amino acids and secondary structures. The selection was made aiming at  
 330 using sequences that have been used on other works for comparison purposes  
 331 [24, 31, 8]. Also we selected those that have different secondary structures to  
 332 see the performance of the proposal on different conditions. We run a small  
 333 test to determine the time required to execute the first 100,000 evaluations of  
 334 the energy function. Based on those results and the size of the sequences, we  
 335 classified them into three groups: (1) *short* sequences, with a size lower than  
 336 40 residues; (2) *medium* sequences, with a size between 40 and 50 residues; and  
 337 (3) *long* sequences, with a size over 50 residues.

PDB ID	Size	Folding pattern	100K Test (sec)
1K43	14	2 $\beta$ -sheet	35,04
1L2Y	20	1 $\alpha$ -helix, 1 $3_{10}$ helix	55,75
2MTW	20	1 $\alpha$ -helix	66,51
*1WQC	26	2 $\alpha$ -helix	90,31
*1ACW	29	1 $\alpha$ -helix, 2 $\beta$ -sheet	63,19
1DFN	30	3 $\beta$ -sheet	75,81
1Q2K	31	1 $\alpha$ -helix, 2 $\beta$ -sheet	105,94
2F4K	35	3 $\alpha$ -helix	125,74
2MR9	44	3 $\alpha$ -helix	152,79
2P81	44	2 $\alpha$ -helix	241,76
*3P7K	45	1 $\alpha$ -helix	229,34
1AB1	46	2 $\alpha$ -helix, 2 $\beta$ -sheet	157,69
*3V1A	48	2 $\alpha$ -helix	225,18
2P6J	52	3 $\alpha$ -helix	324,41
1ENH	54	3 $\alpha$ -helix	272,15
*2JUC	59	3 $\alpha$ -helix, 1 $3_{10}$ helix	278,51
1ROP	63	2 $\alpha$ -helix	300,25
*2P5K	64	3 $\alpha$ -helix, 2 $\beta$ -sheet	330,16
1AIL	73	3 $\alpha$ -helix	332,03
2PMR	87	3 $\alpha$ -helix	402,44

Table 1: List of protein sequences used in computational tests. We classified them in three groups based on its size and a short test of 100,000 evaluations of the energy function. Sequences marked with \* were used for tests to define the best DMA implementation.

338 We also selected 2 sequences of each group (marked with a \* in Table 1)   
 339 that will be used in the tests to determine the best DMA implementation.

340 To validate the results we execute each computational experiment 10 times.   
 341 Experiments were executed in two different environments: (1) a 4-processor   
 342 2.6GHz Intel Xeon CPU with 8 cores (32 cores) and 32GB RAM, and (2) a   
 343 cluster composed by 7 machines (112 cores in total): two 4-processor Intel Xeon   
 344 E-5405 with 2 cores ( $2 \times 8$  cores), three 8-processor Amd Opteron 6128 with   
 345 2 cores ( $3 \times 16$  cores) and two 12-processor Amd Opteron 6168 with 2 cores   
 346 ( $2 \times 24$  cores).

#### 347 4.1. Topology and Energy Configuration

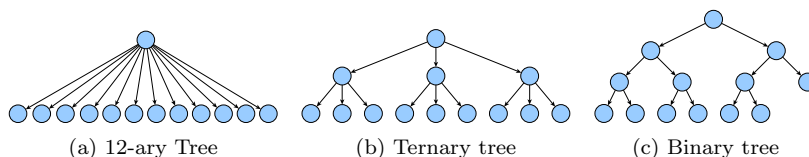


Figure 3: Tree structures evaluated in the tests. All the structures had 13 agents, which are distributed in one, two and three levels on the tree. The binary tree was left incomplete to keep the same number of agents among the structures.

348 The first tests were designed to evaluate the structure of agents and the   
 349 combination of energy functions terms to determine the best configuration for   
 350 the next experiments. We evaluate three hierarchical structures (Figure 3) to   
 351 distribute the population of agents: a two-level ternary tree, an incomplete   
 352 three-level binary tree and a 12-ary tree. All of them were adapted to have 13   
 353 agents to have comparable results.

354 We use the sequences marked with \* in Table 1 to test the tree structures   
 355 and energy combinations. For each sequence, we execute ten runs using each   
 356 one of the three topologies and one of the eleven tested  $\alpha$  to combine the energy   
 357 functions, which gives a total of 1980 experiments. For *short* sequences we use   
 358  $1 \times 10^6$  evaluations of the energy function as a stop criteria,  $2 \times 10^6$  for *medium*   
 359 sequences and  $3 \times 10^6$  for *long* sequences.

#### 360 4.2. Tree structures assessment

361 Figure 4 shows a comparison between the results obtained by the three pop-   
 362 ulation structures evaluated regarding quality of the predicted structures (RMSD   
 363 (CA atoms), free energy) and performance of the execution (total time, num-   
 364 ber of transactions). We consider the mean of the results from each energy   
 365 configuration tested using each structure. As can be seen, the quality of the   
 366 results are similar between different structures. For both RMSD and free energy,   
 367 the average and standard deviation remain in similar ranges for the structures   
 368 evaluated. It reflects that there is no a great impact on the results obtained   
 369 for different distributions of the same number of agents. Similarly occurs for

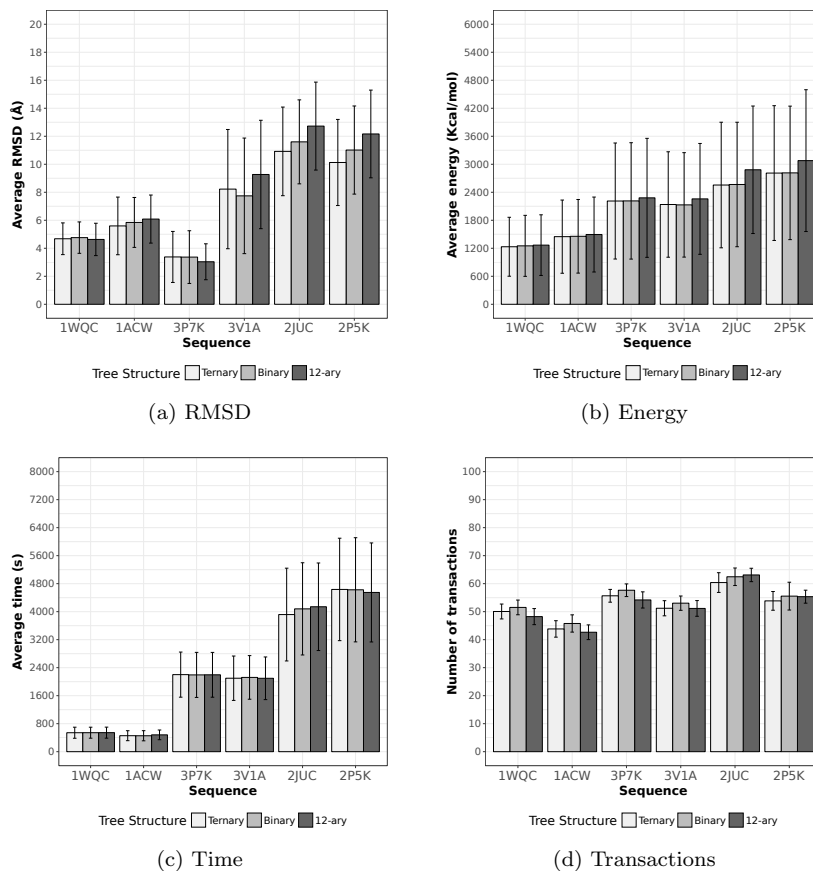


Figure 4: Results of the experiments using different hierarchical structures. For each peptide and tree structure we show the average of the results from all the energy configurations defined by  $\alpha \in \{0.0; 0.1; \dots; 0.9; 1.0\}$ . We show the distribution for the RMSD (CA atoms) vs experimental structures (a), the free energy of the predicted structures (b), the total time of the experiments (c) and average number of transactions between the 13 agents (d).

370 the performance of the executions, where the runtime and the number of trans-  
 371 actions also do not seem to be profoundly affected by the tree structure used  
 372 in the experiment. Extended information about the results can be checked in  
 373 Table 1 and 2 of Supplementary Material.

374 In Table 2 we show the  $p$ -value obtained from the comparison between the  
 375 results of the different topologies using the *Wilcoxon signed-rank test*. As it  
 376 is shown, there are almost no significant differences between the results of the  
 377 tree structures evaluated. Except for the 12-ary tree that shows evidence of a  
 378 slightly lower number of transactions, none of the structures prove to be superior  
 379 respect to the free energy of the solutions, RMSD (CA atoms) calculation with  
 380 experimental structures or executions times. We decided to use the two-level

Sequence	Comparison	p-value				
		RMSD	Energy	Time	Generations	Transactions
1WQC	Ter-Bin	0.478	0.847	0.797	0.818	0.148
	Ter-12ary	0.898	0.748	0.748	0.450	<b>0.047</b>
	Bin-12ary	0.606	0.748	0.606	0.792	<b>0.006</b>
1ACW	Ter-Bin	0.365	1.000	0.949	0.818	0.101
	Ter-12ary	0.171	0.797	0.116	0.921	0.401
	Bin-12ary	0.606	0.797	0.0879	0.742	<b>0.008</b>
3P7K	Ter-Bin	0.562	0.949	0.847	0.645	<b>0.040</b>
	Ter-12ary	0.101	0.748	0.847	0.669	0.139
	Bin-12ary	0.401	0.748	0.847	0.293	<b>0.004</b>
3V1A	Ter-Bin	0.797	0.898	0.949	1.000	0.076
	Ter-12ary	0.365	0.748	0.699	0.693	1.000
	Bin-12ary	0.193	0.699	0.606	0.767	0.076
2JUC	Ter-Bin	0.438	0.949	0.193	0.250	0.088
	Ter-12ary	0.0879	0.606	<b>0.040</b>	0.131	<b>0.013</b>
	Bin-12ary	0.193	0.652	0.847	0.038	0.606
2P5K	Ter-Bin	0.300	0.898	0.748	0.532	0.168
	Ter-12ary	<b>0.0334</b>	0.652	0.562	0.598	0.133
	Bin-12ary	0.193	0.606	0.898	0.224	0.694

Table 2: Comparison between the results from the three hierarchical structures through the *Wilcoxon signed-rank test*. We show the *p-values* obtained from the test for comparing the differences in RMSD (CA atoms), free energy, runtime, total generations and number of transactions. In **bold** we show the *p-values* smaller than 0.05.

381 ternary tree for the next set of testing.

382 *4.3. Energy Function Combination*

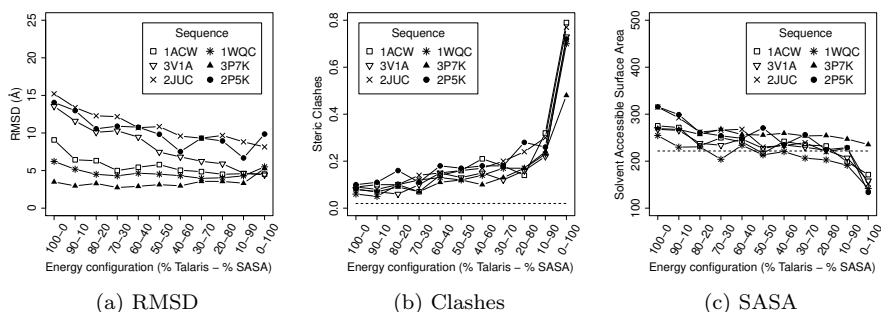


Figure 5: Results of the experiments for different combinations of weights for **Talaris** and **SASA** functions using a ternary tree structure. For each target peptide and energy configuration we show the best RMSD obtained, and the steric clashes and solvent accessible surface area (SASA) of the best solution, calculated using *Gaia Web Tool*.

383 Figure 5 shows a comparison between the results from the 11 weighted combinations of the Talaris and SASA energy functions using a ternary tree. Figure 384 5a shows the RMSD values of the best solution found by each energy configuration for the six sequences tested, while Figures 5b and 5c show the packing 386 quality of the best solutions from the test of *Gaia Web Tool* [25]. The last 387 analysis also extends to the binary and 12-ary trees (see Figures 1 and 2 from 388 the Supplementary Material). 389

390 In Figure 5a shows that by increasing the weighting of the SASA function 391 over Talaris function, the solutions tend to be slightly better regarding RMSD 392 values. It can be explained by the packing properties of the SASA function, 393 leading the structure to a form in which fewer atoms are exposed to the solvent 394 environment. However, solutions also start to show a high amount of steric 395 clashes in their structures as shown in Figure 5b. Talaris function seeks for 396 structures with low levels of energy among all their atom interactions. It means 397 that if Talaris weight is lower than SASA weight, the algorithm will give priority 398 to compact solutions even if these are highly unstable. Structures in nature 399 rarely have clashes in their structure, therefore it is important to maintain this 400 value as lower as possible to generate biologically significant solutions.

401 An extended analysis of the combination of the energy terms in a ternary 402 tree structure can be reviewed in [12]. Based on the above, we decided to use 403 for the next test a value of  $\alpha = 0.4$ , which means that the fitness function will 404 be made of 60% of Talaris energy function and 40% of SASA.

#### 405 4.4. Scalability

406 We perform scalability tests to evaluate the performance of the distributed 407 implementation of the proposed memetic algorithm. These tests aim to assess 408 the performance of the execution and quality of the solutions when a structure 409 with a high number of agents is used. We evaluate two hierarchical structures, a 410 binary and a ternary tree. The ternary tree was tested using two and three levels, 411 for a total of 13 and 40 agents, while the binary tree was tested using three and 412 five levels. We execute 10 runs for each one of these structures using the same 413 testing environment: 20 peptides from Table 1,  $\alpha = 0.4$  to weighing the energy 414 functions based on equation 1 and a stop criteria of  $5 \times 10^6$  evaluations of the 415 energy function among all the agents of the structure. We expect that, since all 416 algorithms evaluate the same number of solutions throughout the execution, the 417 results obtained should be similar, but the execution times should be reduced 418 proportionally to the number of agents used by each structure.

419 Figure 6 shows the average results for 6 of the 20 sequences tested using 420 each one of the structures evaluated. The following conclusion also applies to 421 the others sequences evaluated in the scalability tests. Tables 3 and 4 of Supple- 422 mentary Material have detailed information of all the peptide target sequences 423 evaluated in these tests.

424 As we can observe in Figure 6c, the execution time drops almost proportion- 425 ally to the number of agents used. The three-levels ternary tree, using around 426 three times as many agents, reached the  $5 \times 10^6$  evaluations of the energy func- 427 tion in about a third of the time used for the two-levels ternary tree, as same

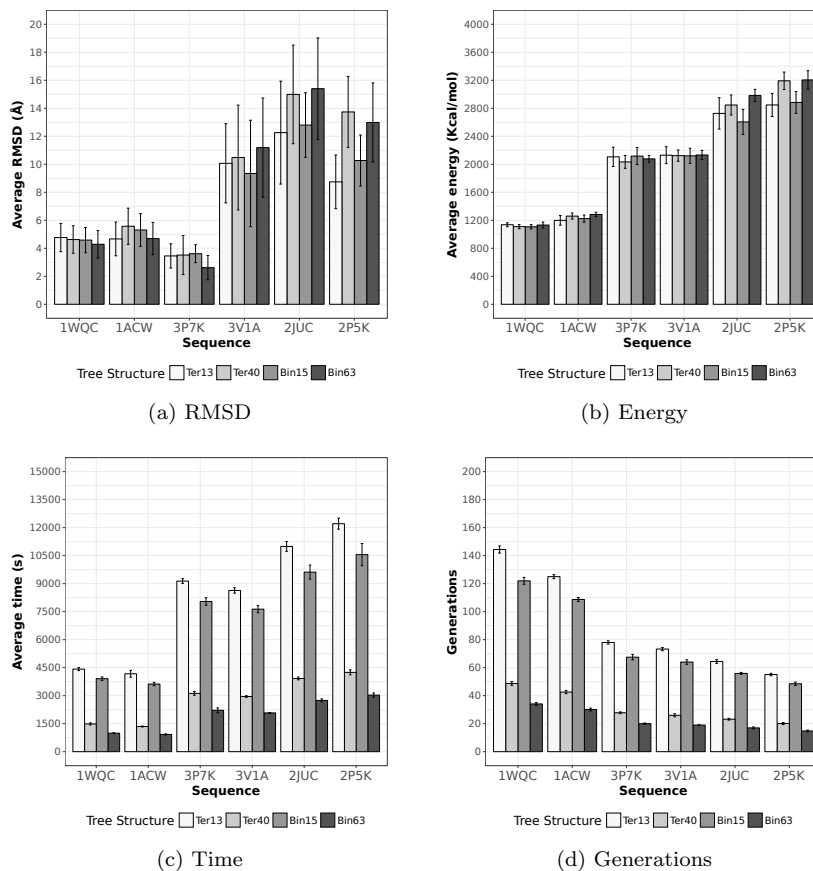


Figure 6: Results of the scalability tests using a ternary and binary trees. The ternary tree was tested using two (13 agents) and three (40 agents) levels and the binary tree using three (15 agents) and five (63 agents). For each peptide and tree structure we show the average results of 10 experiments in terms of RMSD vs experimental structures (a), free energy of the predicted structures (b), total time of the experiments (c) and average generations executed (d).

428 as the five-levels binary tree, which with 4 times more agents reached the end  
 429 criteria in about a quarter of the time used for the three-level binary tree. These  
 430 results apply to all the sequences, being slightly better in the shorter sequences.  
 431 On the other hand, the number of generations reached for the structures was  
 432 reduced proportionally to the increase in the number of agents. This outcome is  
 433 expected, as structures with more agents perform a higher number of evaluations  
 434 of the energy function in each generation.

435 Regarding the quality of the results, Figures 6a and 6b show that the so-  
 436 lutions predicted by all the structures remain in similar ranges, as it would be  
 437 expected for using the same stopping criteria. It can be seen that for *short* and

438 *medium* sequences the differences between all the structures is small, in partic-  
439 ular regarding free energy. For *long* sequences the solutions prove to be slightly  
440 better for both RMSD and free energy in the structures with a fewer number of  
441 agents. However, this difference could be more associated with the efficiency  
442 of the smaller trees in the convergence to better solutions, than the efficiency  
443 of the distributed implementation. This explains why the difference is reduced  
444 in smaller sequences, where all the trees have enough evaluations of the energy  
445 function to reach optimal values. Even considering the latter, the distributed  
446 implementation shows to be highly scalable, which allows it to reduce the execu-  
447 tion times almost proportionally to the number of agents used and to maintain  
448 to a large degree the effectiveness of the algorithm.

#### 449 4.5. Local Search and Diversity Assessment

450 In Section 3.7 we proposed several modifications of the local search algorithm  
451 to improve its efficiency and the quality of the predicted solutions. Two different  
452 approaches were proposed to calculate the diversity of the solutions generated  
453 throughout the algorithm. We tested all these approaches to determine the  
454 best implementation of the DMA-3DPSP. Our goal is to verify which are the  
455 best local search strategies to be employed in 3D protein structure prediction  
456 methods based on memetic algorithms. We seek to keep the diversity of the  
457 population as a way to better explore the protein conformational search space.  
458 The comparison is made in terms of solution quality and execution time. It is  
459 important to mention that we are identifying the algorithm configuration that  
460 is able to find the best solution using a fixed number of evaluations. In order to  
461 deal with the problem, larger runs will be later performed and compared with  
462 the state of the art algorithms.

463 For the testing environment we use the six target sequences marked with \*  
464 in Table 1, a two-level ternary tree structure with 13 agents (Figure 3b) and an  
465 energy fitness function (Equation 1) with  $\alpha = 0.4$  for weighting 60% of Talaris  
466 function and 40% of SASA function. Each experiment was executed ten times  
467 using an end criteria of  $1 \times 10^6$  evaluations of the energy function.

468 The first set of tests compare the original local search algorithm with the  
469 first three modifications proposed in Section 3.7. Unlike LS4, which redefines the  
470 simulating annealing used in the local search, the first three approaches apply  
471 specific changes in the local search algorithm that could generate improvements  
472 in its results. Table 3 shows the results of the tests for the original local search  
473 and the first three new approaches.

474 From the results, all the proposed modifications shown slightly improve the  
475 energy convergence, especially for *long* target sequences (2JUC and 2P5K). The  
476 improvement varies between the different sequences, and none of the three im-  
477 plementations show to be particularly superior to the others. Regarding RMSD,  
478 in four of the six sequences tested (1WQC, 3V1A, 2JUC and 2P5K), all the new  
479 implementations reduced the average RMSD values of the original local search.  
480 The target sequence 3P7K was the only exception where the original local search  
481 shows better solutions regarding RMSD.

Sequence	Structure	Quality of solutions		Execution	
		RMSD (Å)	Energy (Kcal/mol)	Total time (s)	Generations
1WQC	LS-Original	4.91 ( $\pm 1.14$ )	1150.7 ( $\pm 48.2$ )	956.9 ( $\pm 19.0$ )	38.8 ( $\pm 1.0$ )
	LS1	4.75 ( $\pm 1.03$ )	1159.0 ( $\pm 48.1$ )	993.7 ( $\pm 19.2$ )	75.9 ( $\pm 1.5$ )
	LS2	4.79 ( $\pm 0.84$ )	1131.6 ( $\pm 41.1$ )	1055.6 ( $\pm 21.8$ )	156.6 ( $\pm 2.5$ )
	LS3	4.38 ( $\pm 0.90$ )	1108.4 ( $\pm 31.8$ )	975.8 ( $\pm 20.6$ )	20.4 ( $\pm 0.5$ )
1ACW	LS-Original	5.23 ( $\pm 1.94$ )	1298.3 ( $\pm 76.8$ )	871.0 ( $\pm 31.7$ )	34.5 ( $\pm 0.7$ )
	LS1	4.81 ( $\pm 1.45$ )	1260.7 ( $\pm 57.5$ )	925.0 ( $\pm 19.5$ )	68.3 ( $\pm 1.3$ )
	LS2	5.31 ( $\pm 1.31$ )	1251.4 ( $\pm 36.9$ )	995.3 ( $\pm 20.5$ )	138.2 ( $\pm 3.7$ )
	LS3	4.76 ( $\pm 1.46$ )	1297.9 ( $\pm 42.9$ )	887.4 ( $\pm 27.2$ )	18.0 ( $\pm 0.5$ )
3P7K	LS-Original	2.83 ( $\pm 1.16$ )	2173.8 ( $\pm 125.7$ )	2001.1 ( $\pm 51.9$ )	22.0 ( $\pm 0.5$ )
	LS1	2.88 ( $\pm 1.11$ )	2147.3 ( $\pm 170.6$ )	1981.9 ( $\pm 23.6$ )	41.5 ( $\pm 0.8$ )
	LS2	3.54 ( $\pm 0.98$ )	2096.8 ( $\pm 125.9$ )	2063.0 ( $\pm 30.8$ )	84.5 ( $\pm 1.8$ )
	LS3	3.18 ( $\pm 1.39$ )	2053.7 ( $\pm 75.9$ )	2076.5 ( $\pm 61.8$ )	12.1 ( $\pm 0.3$ )
3V1A	LS-Original	12.62 ( $\pm 3.11$ )	2292.5 ( $\pm 110.5$ )	1909.5 ( $\pm 75.8$ )	21.3 ( $\pm 0.7$ )
	LS1	8.53 ( $\pm 5.25$ )	2144.4 ( $\pm 188.4$ )	1890.7 ( $\pm 46.5$ )	39.5 ( $\pm 0.8$ )
	LS2	10.36 ( $\pm 3.66$ )	2118.8 ( $\pm 159.8$ )	1957.3 ( $\pm 23.5$ )	78.3 ( $\pm 1.3$ )
	LS3	10.32 ( $\pm 4.00$ )	2107.1 ( $\pm 85.2$ )	2013.2 ( $\pm 81.0$ )	11.6 ( $\pm 0.5$ )
2JUC	LS-Original	15.05 ( $\pm 2.18$ )	3260.8 ( $\pm 169.1$ )	2494.8 ( $\pm 89.3$ )	19.2 ( $\pm 0.8$ )
	LS1	13.26 ( $\pm 3.40$ )	2805.1 ( $\pm 188.2$ )	2419.8 ( $\pm 58.2$ )	35.3 ( $\pm 0.9$ )
	LS2	11.30 ( $\pm 2.53$ )	2755.9 ( $\pm 233.7$ )	2491.6 ( $\pm 75.9$ )	69.4 ( $\pm 1.4$ )
	LS3	12.36 ( $\pm 2.00$ )	2671.1 ( $\pm 78.7$ )	2629.3 ( $\pm 93.3$ )	10.6 ( $\pm 0.5$ )
2P5K	LS-Original	13.49 ( $\pm 4.33$ )	3497.6 ( $\pm 320.7$ )	2754.3 ( $\pm 92.7$ )	16.8 ( $\pm 0.6$ )
	LS1	12.15 ( $\pm 2.24$ )	3084.6 ( $\pm 113.8$ )	2595.4 ( $\pm 53.9$ )	30.1 ( $\pm 0.7$ )
	LS2	8.96 ( $\pm 2.33$ )	3005.9 ( $\pm 227.4$ )	2713.4 ( $\pm 48.4$ )	60.6 ( $\pm 1.3$ )
	LS3	12.01 ( $\pm 2.28$ )	3078.2 ( $\pm 108.4$ )	2813.0 ( $\pm 147.6$ )	9.3 ( $\pm 0.5$ )

Table 3: Results of the experiments for different local search approaches in six peptide target sequences. LS-Original uses the original local search implementation described in Section 3.5, while LS1, LS2 and LS3 apply the modifications in the local search described in Section 3.7. For each amino acid and local search implementation, we show the average results of 10 experiments in terms of quality of the solutions (RMSD and free energy) and performance of the execution (runtime and generations performed).

482 Regarding the execution performance, the new implementations do not seem  
483 to generate great differences in the runtimes. Only for *short* target peptide  
484 sequences, we can see a slight increase in the total time by the LS1 and LS2  
485 implementations. On the other hand, the number of generations reached by the  
486 algorithm changes with the modifications in the local search, where LS1 and  
487 especially LS2 perform a larger number of generations and LS3 considerably  
488 less. These results are justified by each implementation. LS1 favours better  
489 solutions in the neighbourhood changes, which naturally reduces the number  
490 of iterations of simulated annealing by having more optimal solutions. LS2  
491 combines the small modifications on angles phi and psi in one cycle of iterations,  
492 so the simulated annealing naturally consumes less time and allows to reach a  
493 larger number of generations. LS3 improves the search range of the angles, which

494 generates a more efficient simulated annealing that performs more number of  
iterations per generation.

Sequence	Structure	Quality of solutions		Execution	
		RMSD (Å)	Energy (Kcal/mol)	Total time (s)	Generations
1WQC	Div1	3.83 ( $\pm 1.15$ )	1060.5 ( $\pm 34.6$ )	1004.8 ( $\pm 17.2$ )	26.6 ( $\pm 0.5$ )
	Div2-2.0	3.86 ( $\pm 0.75$ )	1061.0 ( $\pm 19.3$ )	1070.4 ( $\pm 27.5$ )	23.6 ( $\pm 0.7$ )
	Div2-4.0	4.09 ( $\pm 1.04$ )	1065.0 ( $\pm 42.7$ )	1037.4 ( $\pm 28.5$ )	25.1 ( $\pm 0.9$ )
1ACW	Div1	4.68 ( $\pm 1.54$ )	1161.2 ( $\pm 47.2$ )	934.6 ( $\pm 26.0$ )	23.8 ( $\pm 0.6$ )
	Div2-2.0	5.04 ( $\pm 0.91$ )	1211.5 ( $\pm 43.1$ )	961.8 ( $\pm 30.5$ )	20.4 ( $\pm 0.5$ )
	Div2-4.0	4.26 ( $\pm 1.64$ )	1206.1 ( $\pm 58.2$ )	971.0 ( $\pm 28.4$ )	21.2 ( $\pm 0.6$ )
3P7K	Div1	2.99 ( $\pm 1.33$ )	2009.2 ( $\pm 71.9$ )	2048.6 ( $\pm 50.5$ )	19.0 ( $\pm 0.5$ )
	Div2-2.0	3.04 ( $\pm 0.84$ )	1988.6 ( $\pm 52.5$ )	2171.7 ( $\pm 92.3$ )	17.8 ( $\pm 0.9$ )
	Div2-4.0	3.86 ( $\pm 1.36$ )	2038.7 ( $\pm 141.7$ )	2064.1 ( $\pm 51.0$ )	18.6 ( $\pm 0.5$ )
3V1A	Div1	9.21 ( $\pm 3.66$ )	2080.6 ( $\pm 69.4$ )	1916.0 ( $\pm 52.7$ )	17.6 ( $\pm 0.5$ )
	Div2-2.0	8.18 ( $\pm 4.48$ )	2059.2 ( $\pm 113.0$ )	2029.5 ( $\pm 36.5$ )	16.0 ( $\pm 0.0$ )
	Div2-4.0	8.84 ( $\pm 4.26$ )	2035.3 ( $\pm 86.1$ )	2063.0 ( $\pm 73.4$ )	17.0 ( $\pm 0.7$ )
2JUC	Div1	13.48 ( $\pm 2.54$ )	2642.1 ( $\pm 87.9$ )	2503.4 ( $\pm 96.1$ )	16.0 ( $\pm 0.5$ )
	Div2-2.0	14.42 ( $\pm 2.55$ )	2762.4 ( $\pm 143.0$ )	2776.3 ( $\pm 79.1$ )	14.9 ( $\pm 0.3$ )
	Div2-4.0	12.71 ( $\pm 2.54$ )	2737.1 ( $\pm 115.2$ )	2787.6 ( $\pm 83.9$ )	15.1 ( $\pm 0.6$ )
2P5K	Div1	13.29 ( $\pm 3.95$ )	2929.8 ( $\pm 151.5$ )	2756.5 ( $\pm 87.7$ )	14.6 ( $\pm 0.5$ )
	Div2-2.0	12.88 ( $\pm 2.96$ )	3039.4 ( $\pm 106.3$ )	2967.2 ( $\pm 164.6$ )	13.8 ( $\pm 0.6$ )
	Div2-4.0	13.74 ( $\pm 1.29$ )	3041.9 ( $\pm 134.5$ )	2938.6 ( $\pm 174.2$ )	13.6 ( $\pm 0.5$ )

Table 4: Results of the experiments for different diversity approaches in six target sequences. Div1 and Div2 follow the diversity implementations described in Section 3.8. Div2 was tested using a RMSD cutoff of 2.0Å and 4.0Å. For each amino acid and diversity implementation, we show the average results of 10 experiments in terms of quality of the solutions (RMSD and free energy) and performance of the execution (runtime and generations performed). For the RMSD and GDT measurements we show in bold the best results between the two implementations

495  
496 We also tested the two diversity approaches proposed in Section 3.8: Div1  
497 that compares the average difference of the angles of two solutions, and Div2  
498 that calculates the RMSD between their structures. We tested Div2 using 2.0  
499 and 4.0 as RMSD cutoff to determine if a solution is diverse, while Div1 limit  
500 was calculated based on the size of the peptide (six times). Table 4 shows the  
501 results of the tests using the two diversity implementations. Both, the quality  
502 of the solutions and the execution performance reached similar ranges and do  
503 not show significant differences between the three implementations evaluated.

504 The first group of experiments shows that all of the local search approaches  
505 evaluated improve the predicted solutions. These three approaches are not mu-  
506 tually exclusive, therefore, for the last set of tests, we compare a new implemen-  
507 tation of the local search that combines the approaches of LS1, LS2, and LS3  
508 with the original local search implementation and with LS4, the last approach  
509 proposed in Section 3.7. LS4 defines a new local search algorithm, combining  
510 the change of neighborhood with the simulating annealing process, therefore  
511 approaches as LS1 does not work in this implementation.

Sequence	Structure	Quality of solutions		Execution	
		RMSD (Å)	Energy (Kcal/mol)	Total time (s)	Generations
1WQC	LS-Original	4.33 ( $\pm 1.10$ )	1180.2 ( $\pm 59.1$ )	1040.3 ( $\pm 26.1$ )	23.1 ( $\pm 0.3$ )
	LS1-2-3	3.86 ( $\pm 0.75$ )	1061.0 ( $\pm 19.3$ )	1070.4 ( $\pm 27.5$ )	23.6 ( $\pm 0.7$ )
	LS4	3.55 ( $\pm 0.53$ )	1009.6 ( $\pm 26.4$ )	1106.5 ( $\pm 25.1$ )	25.4 ( $\pm 0.5$ )
1ACW	LS-Original	5.70 ( $\pm 0.61$ )	1348.3 ( $\pm 58.2$ )	966.7 ( $\pm 27.5$ )	20.7 ( $\pm 0.5$ )
	LS1-2-3	5.04 ( $\pm 0.91$ )	1211.5 ( $\pm 43.1$ )	961.8 ( $\pm 30.5$ )	20.4 ( $\pm 0.5$ )
	LS4	4.43 ( $\pm 1.17$ )	1102.3 ( $\pm 42.7$ )	1065.6 ( $\pm 28.2$ )	22.6 ( $\pm 0.5$ )
3P7K	LS-Original	2.70 ( $\pm 1.06$ )	2302.3 ( $\pm 166.9$ )	2176.5 ( $\pm 79.2$ )	18.0 ( $\pm 0.7$ )
	LS1-2-3	3.04 ( $\pm 0.84$ )	1988.6 ( $\pm 52.5$ )	2171.7 ( $\pm 92.3$ )	17.8 ( $\pm 0.9$ )
	LS4	2.96 ( $\pm 1.00$ )	1879.9 ( $\pm 36.3$ )	2187.8 ( $\pm 59.6$ )	19.7 ( $\pm 0.7$ )
3V1A	LS-Original	11.73 ( $\pm 1.97$ )	2475.5 ( $\pm 225.6$ )	2078.0 ( $\pm 63.0$ )	16.4 ( $\pm 0.5$ )
	LS1-2-3	8.18 ( $\pm 4.48$ )	2059.2 ( $\pm 113.0$ )	2029.5 ( $\pm 36.5$ )	16.0 ( $\pm 0.0$ )
	LS4	4.88 ( $\pm 2.16$ )	1829.2 ( $\pm 47.8$ )	2102.0 ( $\pm 78.1$ )	17.7 ( $\pm 0.5$ )
2JUC	LS-Original	13.08 ( $\pm 2.23$ )	3621.7 ( $\pm 422.0$ )	2851.2 ( $\pm 132.2$ )	15.5 ( $\pm 0.7$ )
	LS1-2-3	14.42 ( $\pm 2.55$ )	2762.4 ( $\pm 143.0$ )	2776.3 ( $\pm 79.1$ )	14.9 ( $\pm 0.3$ )
	LS4	9.94 ( $\pm 1.82$ )	2214.1 ( $\pm 65.8$ )	2655.9 ( $\pm 77.5$ )	15.9 ( $\pm 0.6$ )
2P5K	LS-Original	13.22 ( $\pm 2.94$ )	3624.9 ( $\pm 305.9$ )	2987.5 ( $\pm 72.8$ )	13.9 ( $\pm 0.3$ )
	LS1-2-3	12.88 ( $\pm 2.96$ )	3039.4 ( $\pm 106.3$ )	2967.2 ( $\pm 164.6$ )	13.8 ( $\pm 0.6$ )
	LS4	9.73 ( $\pm 1.98$ )	2492.8 ( $\pm 147.9$ )	2966.3 ( $\pm 80.3$ )	14.8 ( $\pm 0.4$ )

Table 5: Results of the experiments for different local search approaches in six peptides. LS-Original uses the original local search implementation described in Section 3.5. LS1-2-3 apply the modifications LS1, LS2 and LS3 described in Section 3.7 at the same time in the local search. LS4 use the last local search implementation described in Section 3.7. For each amino acid and local search implementation, we show the average results of 10 experiments in terms of quality of the solutions (RMSD and free energy) and performance of the execution (runtime and generations performed).

512 Table 5 shows the results of the three local search implementations tested.  
513 LS1-2-3 and LS4 reached energy values significantly better than the original  
514 local search implementation. Regarding RMSD, LS1-2-3 obtained better solutions  
515 than the original implementation in four of the six target peptide sequences  
516 (1WQC, 1ACW, 3V1A, and 2P5K), while LS4 got better solutions than the other  
517 two implementations in all the peptides tested, except for the sequence 3P7K  
518 where the original algorithm was slightly better. Again, the case for the 3P7K  
519 sequence can be explained by an imprecision in the energy function, since the  
520 energy values of LS1-2-3 and LS4 were lower than those obtained by the original  
521 implementation. The performance of the execution remained in similar ranges  
522 for all implementations; therefore there does not seem to be a drawback in using  
523 the newly proposed implementations.

#### 524 4.6. Best implementation of the DMA-3DPSP

525 Using the above results, we measure the performance and effectiveness of  
526 the best performing configuration: LS4 local search, and Div2 with a cutoff of  
527 2.0.

528 Table 6 shows the results of the last group of experiments. The DMA-Original  
 529 corresponds to the original simulating annealing algorithm described in Section  
 530 3.5 using Div1 as diversity strategy in updating pockets. DMA-Final uses the  
 531 latest local search implementation, the LS4 strategy updated with LS2 and  
 532 LS3 approaches in its second phase and it uses Div2 for diversity calculations.  
 533 The different methods were tested using the two-level ternary tree topology for  
 534 agents distribution, an energy fitness function weighting 60% of Talaris and 40%  
 535 of SASA, and  $5 \times 10^6$  evaluations of the energy functions as the stop criteria.

Results in Table 6 were assessed in terms of structural quality of the solutions  
 by comparing the three-dimensional structures with the experimental ones using  
 two measurements: RMSD and the *Global Distance Test* total score (GDT\_TS [48]),  
 metric that gives a more accurate comparison by measuring the percentage of  
 atoms of the structures within the distance cutoff using different superpositions.  
 Equation 6 describes the calculation of the GDT total score. Small values of  
 RMSD and large values of GDT means good quality solutions.

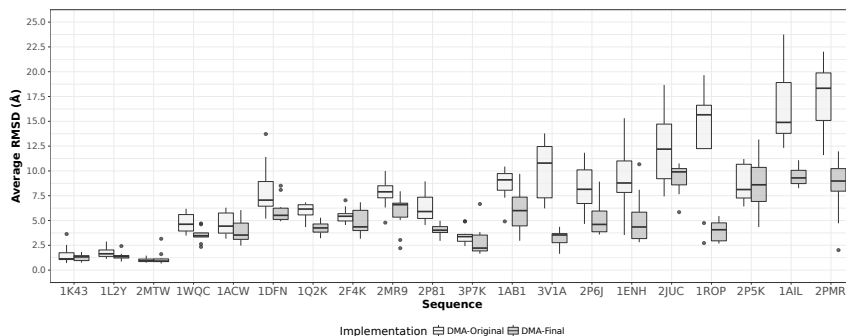
$$GDT_{TS} = \frac{GDT_{P1} + GDT_{P2} + GDT_{P4} + GDT_{P8}}{4} \quad (6)$$

536 where  $GDT_{P_n}$  is the percentage of the residues under the cutoff  $\leq n\text{\AA}$ . We  
 537 also show the average values of free energy, total execution time and number of  
 538 generations reached by each implementation.

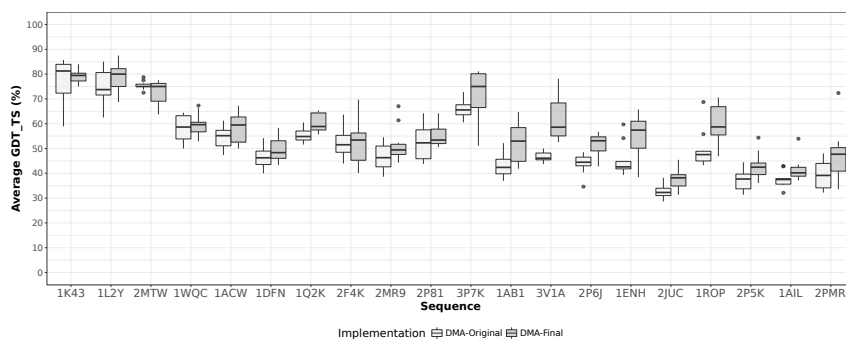
539 Figure 7 shows a comparison between the distribution of the results from  
 540 both implementations regarding RMSD and GDT. We observed that the solutions  
 541 predicted by the DMA-Final show improvements in both metrics over the original  
 542 implementation. In 18 out of 20 peptides evaluated the DMA-Final solutions  
 543 obtained lower average RMSD values, especially for *long* sequences. Also, in 19  
 544 out of 20 target peptides, the best solutions found were predicted by the DMA-  
 545 Final implementation. The same applies concerning the GDT\_TS percentages. In  
 546 18 out of 20 peptides the average GDT\_TS of the DMA-Final solutions were higher  
 547 than the ones predicted by DMA-Original, and in 17 out of 20 of these sequences  
 548 the best solutions in terms of GDT\_TS were found by DMA-Final.

#### 549 4.7. Biological Significance Assessment

550 Even with good RMSD and GDT scores, the solutions predicted by the imple-  
 551 mentations must prove to have biological significance. We examine the tertiary  
 552 structure of the best solutions found by both implementations using PyMOL [40].  
 553 Cartoon representation of the solutions with the lowest RMSD obtained by each  
 554 DMA-3DPSP implementation are shown in Figure 8. We observed that the  $\alpha$ -  
 555 helix secondary structures are formed and positioned correctly in almost all  
 556 structures. The  $\beta$ -sheet secondary structures, on the other hand, although well  
 557 positioned prove to be difficult to form correctly in several of the predicted struc-  
 558 tures (1ACW, 1DFN, 1Q2K and 2P5K are good examples). This can be explained  
 559 because the search space of these structures tends to overlap much more with  
 560 one of turns and coils than the  $\alpha$ -helix structures, which makes it much harder  
 561 to predict even with experimental knowledge.



(a) RMSD



(b) GDT\_TS

Figure 7: Comparison of the quality of the solutions from the original and the latest implementation. For each peptide we show the distribution of the results in terms of RMSD values and GDT percentage vs experimental structures.

562 We also performed quality assessment of the best solutions regarding RMSD  
 563 using *Gaia Web Tool* [25]. *Gaia* measures the packing quality of a given protein  
 564 structure using several metrics, between them, the number of steric clashes, the  
 565 percentage of unsatisfied hydrogen bonds, total void volume and the scaling of  
 566 the accessible surface area with protein length.

567 Table 7 shows the results of the *Gaia* reports for the predicted structures  
 568 with the lowest RMSD of both DMA implementations. It can be observed that, in  
 569 all the peptide sequences evaluated in the tests, the DMA-Final implementation  
 570 reduced both the number of steric clashes and the area exposed to the solvent  
 571 of the predicted structures over the original DMA implementation. It proves that  
 572 the solutions of the new implementation not only have better RMSD and GDT but  
 573 also their structures are more biologically significant. There is no much impact  
 574 in the percentage of unsatisfied hydrogen bonds, and these remain in good  
 575 margin for both implementations, but there was an excess of void volume in  
 576 several of the structures predicted by the DMA-Final implementation. It is also

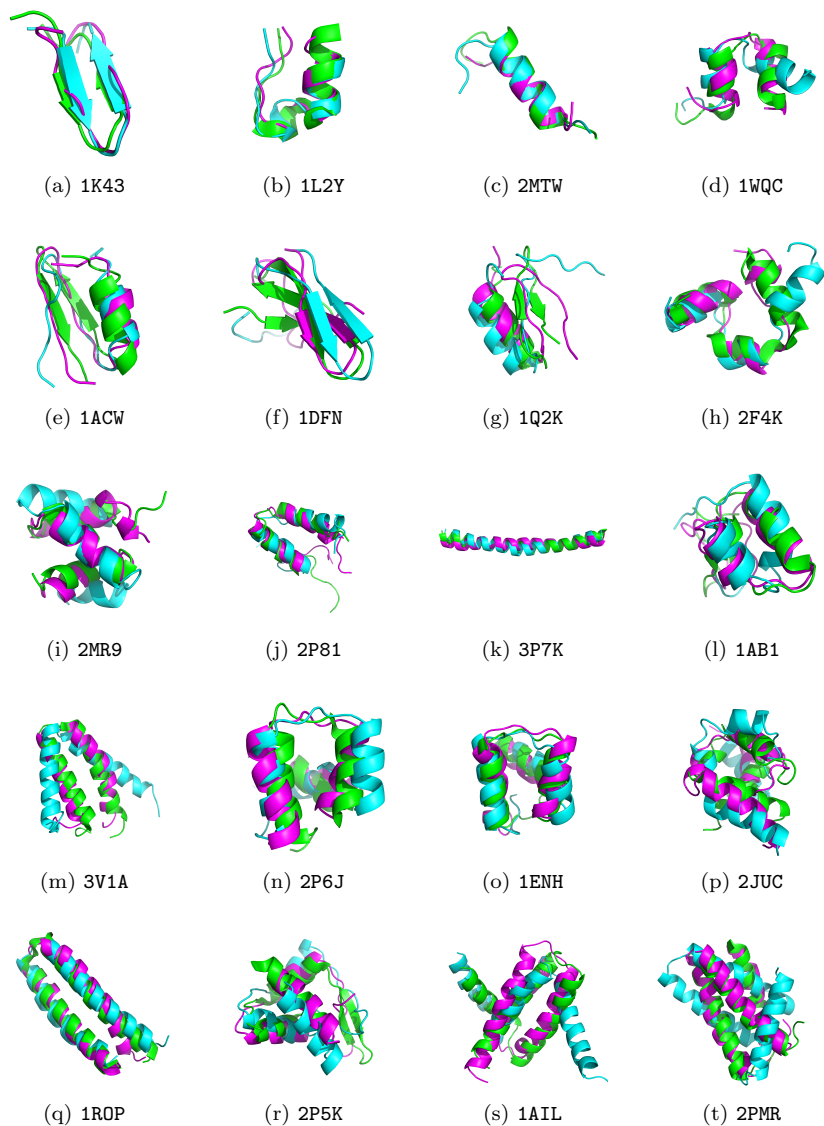


Figure 8: Cartoon representation of the best solution found for each peptide tested in the final experiments. We show in green the experimental protein, in cyan the structure predicted by the original DMA-3DPSP and in magenta the structure predicted by the final update of the DMA-3DPSP. The  $C_{\alpha}$  atoms of the experimental and predicted structures are fitted.

577 noticeable that the results of both implementations show a much higher level  
 578 of steric clashes than the target expected by Gaia. In experimental structures,  
 579 the number of clashes rarely exceeds 1% of total contacts, so there is still room

580 to improve the results of the DMA-3DPSP.

#### 581 4.8. Comparison with State of the Art Methods

582 We also test the best approach using peptides that have been tested in the  
583 CASP-2016 experiments [31] and compare our results with some of the “reference  
584 methods” in the PSP area. We chose five of the protein targets from CASP 11  
585 experiments [24] with sizes from 70 to 128 residues and different folding patterns.  
586 Since the sequences evaluated are larger and structurally more complex than our  
587 previous case studies, we set  $1 \times 10^7$  evaluations of the energy functions as the  
588 stopping criteria of the DMA-3DPSP. We performed 10 runs of the algorithm  
589 using a two-levels ternary tree structure to distribute the agents and an energy  
590 configuration weighting 60% of Talaris function and 40% fo SASA function.

591 Table 9 shows the results of the DMA-3DPSP for the five protein targets from  
592 CASP 11 experiments. We show the results obtained by two of the methods  
593 with the best performance over the last CASP experiments, Zhang-Server [46],  
594 which combine the *ab initio* strategy of QUARK [45] and the threading modeling  
595 of I-TASSER [47], and Robetta [23], the Rosetta PSP web server of Baker Labs.

596 We see that the results of the DMA-3DPSP are still far from the top methods of  
597 the last CASP experiments in both RMSD and GDT values. Even the best solutions  
598 found in all the sequences are still not close to the Zhang-Server and Robetta  
599 results. The proposed strategy requires a large amount of time and resource for  
600 the predictions of sequences over 100 residues, therefore the  $1 \times 10^7$  evaluations  
601 of the energy function proves to be insufficient to predict the sequences selected  
602 in the tests. Further testing is still needed to determine the maximum scope of  
603 the proposed strategy.

## 604 5. Conclusion

605 In this work, we presented a Distributed Memetic Algorithm (DMA-3DPSP) to  
606 deal with the Tertiary Protein Structure Prediction Problem (3D-PSP problem).  
607 The 3D-PSP problem is one of the hardest problem to solve and it requires the  
608 developed of smart and efficient techniques. In this work we evaluated different  
609 configurations of a Distributed Memetic Algorithm in order to measure which  
610 one shows a better contribution to improve the quality of solutions.

611 The algorithm uses an asynchronous hierarchical population of agents that  
612 shares the solution’s information along with the execution of the algorithm.  
613 We also used a combination of well-known energy scoring function as a fitness  
614 value for the optimisation process. The DMA-3DPSP algorithm uses information  
615 extracted from the PDB to guide the search process. We evaluated three different  
616 hierarchical population structures and several alternatives to local search and  
617 population operators.

618 We perform a thorough testing process in order to correctly identify the  
619 best performing configuration in terms of computational results and biological  
620 significance. All the approaches developed were tested in order to elucidate dif-  
621 ferent ways to improve the current protein structure prediction methods based

622 on memetic algorithms. We investigated different strategies to be considered  
623 in the design of prediction methods for the 3D structure of proteins, show-  
624 ing which mechanisms bring advantages in terms of computational time and  
625 biological significance. Results show that the selection of the combination of  
626 energies allows the generation of solutions that have a good trade-off between  
627 similarity with the experimental protein and biological significance. The use of  
628 a distributed strategy improved the time of the solutions as expected, and they  
629 show slight improvements in terms of solution quality. In our experiments,  
630 communications strategies, and solution sharing in the distributed algorithm,  
631 did not show a significant improvement effect, so this is an area that needs to  
632 be further explored.

633 We were able to identify a configuration for the DMA-3DPSP that can find  
634 solutions with good computational quality, and also with good biological signif-  
635 icance, which is a must in bioinformatics. Further work will be focused on the  
636 improvement of the implementation to reach similar quality solutions on larger  
637 sequences, and also in the incorporation of new machine learning techniques to  
638 extract new information from the PDB. Also, the incorporation of different local  
639 search strategies that can be combined during the evolution process, in order to  
640 exploit better the search space.

#### 641 **Acknowledgments**

642 This work was partially funded by grants STIC-AMSUD (88887.135130 /2017-  
643 01 CAPES-Brazil and 17-STIC-05 Conicyt-Chile), FAPERGS ( 16/2551-0000520-  
644 6; 19/2551-0001906-8), MCT/CNPq (311611/2018-4), and DICYT-VRIDEI USACH  
645 (061919IP).

#### 646 **References**

- 647 [1] R. M. Abaskharon and F. Gai. Meandering down the energy landscape  
648 of protein folding: Are we there yet? *Biophysical Journal*, 110(9):1924 –  
649 1932, 2016.
- 650 [2] Nikolai Alexandrov and Ilya Shindyalov. PDP: protein domain parser.  
651 *Bioinformatics*, 19(3):429–430, 02 2003.
- 652 [3] M. Blocho and Z.J. Czech. A parallel memetic algorithm for the vehicle  
653 routing problem with time windows. In *2013 Eighth International Confer-*  
654 *ence on P2P, Parallel, Grid, Cloud and Internet Computing (3PGCIC)*,  
655 pages 144–151, Oct 2013.
- 656 [4] B Borguesan, M.B. e Silva, B. Grisci, M. Inostroza-Ponta, and M. Dorn.  
657 Apl: An angle probability list to improve knowledge-based metaheuris-  
658 tics for the three-dimensional protein structure prediction. *Comput. Biol.*  
659 *Chem.*, 59:142–157, 2015.

- 660 [5] B. Borguesan, M. Inostroza-Ponta, and M. Dorn. NIAS-Server: Neighbors  
661 influence of amino acids and secondary structures in proteins. *Journal of*  
662 *Computational Biology*, 24(3):255–265, 2017.
- 663 [6] Joseph D Bryngelson, Jose Nelson Onuchic, Nicholas D Socci, and Peter G  
664 Wolynes. Funnels, pathways, and the energy landscape of protein folding:  
665 a synthesis. *Proteins: Struct., Funct., Bioinf.*, 21(3):167–195, 1995.
- 666 [7] David Case, Thomas Cheatham, Tom Darden, Holger Gohlke, Ray Luo,  
667 Kenneth Merz, Alexey Onufriev, Carlos Simmerling, Bing Wang, and  
668 Robert Woods. The amber biomolecular simulation programs. *Journal*  
669 *of computational chemistry*, 26:1668–88, 12 2005.
- 670 [8] Leonardo de Lima Correa, Bruno Borguesan, Camilo Farfán, Mario  
671 Inostroza-Ponta, and Márcio Dorn. A memetic algorithm for 3d protein  
672 structure prediction problem. *IEEE/ACM Trans. Comput. Biology Bioin-*  
673 *form.*, 15(3):690–704, 2018.
- 674 [9] K. A. Dill, S. B. Ozkan, T. R. Weikl, J. D. Chodera, and V. A. Voelz.  
675 The protein folding problem: when will it be solved? *Current Opinion in*  
676 *Structural Biology*, 17 3:342–6, 2007.
- 677 [10] M Dorn, MB e Silva, L Buriol, and L Lamb. Three-dimensional protein  
678 structure prediction: methods and computational strategies. *Comput. Biol.*  
679 *Chem.*, 53:251–276, 2014.
- 680 [11] Márcio Dorn and Osmar Norberto de Souza. Cref: A central-residue-  
681 fragment-based method for predicting approximate 3-d polypeptides struc-  
682 tures. In *Proceedings of the 2008 ACM Symposium on Applied Computing*,  
683 SAC '08, pages 1261–1267, New York, NY, USA, 2008. ACM.
- 684 [12] I. Escobar, N. Hidalgo, M. Inostroza-Ponta, M. Marín, E. Rosas, and  
685 M. Dorn. Evaluation of a combined energy fitness function for a distributed  
686 memetic algorithm to tackle the 3d protein structure prediction problem.  
687 In *2016 35th International Conference of the Chilean Computer Science*  
688 *Society (SCCC)*, pages 1–10, Oct 2016.
- 689 [13] Narayanan Eswar, David Eramian, Ben Webb, Min-Yi Shen, and Andrej  
690 Sali. *Protein Structure Modeling with MODELLER*, pages 145–159. Hu-  
691 mana Press, Totowa, NJ, 2008.
- 692 [14] Antonio J. Fernández and Álvaro Gutiérrez-Rodríguez. On distributed  
693 user-centric memetic algorithms. *Soft Comput.*, 23(12):4019–4039, 2019.
- 694 [15] Rasmus Fonseca, Martin Paluszewski, and Pawel Winter. Protein struc-  
695 ture prediction using bee colony optimization metaheuristic. *Journal of*  
696 *Mathematical Modelling and Algorithms*, 9(2):181–194, 2010.

- 697 [16] Nick Gibbs, Anthony R. Clarke, and Richard B. Sessions. Ab initio pro-  
698 tein structure prediction using physicochemical potentials and a simpli-  
699 fied off-lattice model. *Proteins: Structure, Function, and Bioinformatics*,  
700 43(2):186–202, 2001.
- 701 [17] Pawel Gniewek, Andrzej Kolinski, Andrzej Kloczkowski, and Dominik  
702 Gront. Bioshell-threading: versatile monte carlo package for protein 3d  
703 threading. *BMC bioinformatics*, 15(1):22, 2014.
- 704 [18] Yue-Jiao Gong, Wei-Neng Chen, Zhi-Hui Zhan, Jun Zhang, Yun Li, Qingfu  
705 Zhang, and Jing-Jing Li. Distributed evolutionary algorithms and their  
706 models: A survey of the state-of-the-art. *Applied Soft Computing*, 34:286  
707 – 300, 2015.
- 708 [19] M. Harris, R. Berretta, M. Inostroza-Ponta, and P. Moscato. A memetic al-  
709 gorithm for the quadratic assignment problem with parallel local search. In  
710 *IEEE Congress on Evolutionary Computation, CEC 2015, Sendai, Japan,*  
711 *May 25-28, 2015*, pages 838–845, 2015.
- 712 [20] Mario Inostroza-Ponta, Camilo Farfán, and Marcio Dorn. A memetic algo-  
713 rithm for protein structure prediction based on conformational preferences  
714 of aminoacid residues. In *Genetic and Evolutionary Computation Confer-*  
715 *ence, GECCO 2015, Madrid, Spain, July 11-15, 2015, Companion Material*  
716 *Proceedings*, pages 1403–1404, 2015.
- 717 [21] Nanda Dulal Jana, Swagatam Das, and Jaya Sil. *Backgrounds on Protein*  
718 *Structure Prediction and Metaheuristics*, pages 1–28. Springer International  
719 Publishing, Cham, 2018.
- 720 [22] Morten Källberg, Gohar Margaryan, Sheng Wang, Jianzhu Ma, and Jinbo  
721 Xu. Raptorx server: a resource for template-based protein structure mod-  
722 eling. In *Protein Structure Prediction*, pages 17–27. Springer, 2014.
- 723 [23] David E Kim, Dylan Chivian, and David Baker. Protein structure predic-  
724 tion and analysis using the rosetta server. *Nucleic acids research*, 32(suppl  
725 2):W526–W531, 2004.
- 726 [24] Lisa N Kinch, Wenlin Li, R Dustin Schaeffer, Roland L Dunbrack, Bohdan  
727 Monastyrskyy, Andriy Kryshtafovych, and Nick V Grishin. Casp 11 target  
728 classification. *Proteins: Structure, Function, and Bioinformatics*, 2016.
- 729 [25] Pradeep Kota, Feng Ding, Srinivas Ramachandran, and Nikolay V.  
730 Dokholyan. Gaia: automated quality assessment of protein structure mod-  
731 els. *Bioinformatics (Oxford, England)*, 27(16):2209–2215, Aug 2011.
- 732 [26] Nuno Leite, Carlos M. Fernandes, Fernando Melicio, and Agostinho C.  
733 Rosa. A cellular memetic algorithm for the examination timetabling prob-  
734 lem. *Computers & OR*, 94:118–138, 2018.

- 735 [27] Dapeng Li, Tonghua Li, Peisheng Cong, Wenwei Xiong, and Jiangming  
736 Sun. A novel structural position-specific scoring matrix for the prediction  
737 of protein secondary structures. *Bioinformatics*, 28(1):32–39, January 2012.
- 738 [28] Cheng-Jian Lin and Shih-Chieh Su. Protein 3d hp model folding simula-  
739 tion using a hybrid of genetic algorithm and particle swarm optimization.  
740 *International Journal of Fuzzy Systems*, 13:140–147, 06 2011.
- 741 [29] Dimitrios P. Lyras and Dirk Metzler. Reformalign: improved multiple  
742 sequence alignments using a profile-based meta-alignment approach. *BMC*  
743 *Bioinformatics*, 15(1):265, 2014.
- 744 [30] P. Moscato. On evolution, search, optimization, genetic algorithms and  
745 martial arts: Towards memetic algorithms. Technical Report Caltech Con-  
746 current Computation Program, Report. 826, CALTECH, Pasadena, Cali-  
747 fornia, USA, 1989.
- 748 [31] John Moult, Krzysztof Fidelis, Andriy Kryshchak, Torsten Schwede,  
749 and Anna Tramontano. Critical assessment of methods of protein structure  
750 prediction: Progress and new directions in round xi. *Proteins: Structure,*  
751 *Function, and Bioinformatics*, 84(S1):4–14, 2016.
- 752 [32] Jakub Nalepa and Miroslaw Blocho. Co-operation in the parallel memetic  
753 algorithm. *International Journal of Parallel Programming*, 43(5):812–839,  
754 2015.
- 755 [33] Jakub Nalepa and Miroslaw Blocho. Adaptive cooperation in parallel  
756 memetic algorithms for rich vehicle routing problems. *IJGUC*, 9(2):179–  
757 192, 2018.
- 758 [34] Quang Huy Nguyen, Yew-Soon Ong, Meng-Hiot Lim, and Natalio Krasno-  
759 gor. Adaptive cellular memetic algorithms. *Evolutionary Computation*,  
760 17(2):231–256, 2009.
- 761 [35] M. A. Rashid, M. A. H. Newton, M. T. Hoque, and A. Sattar. A local search  
762 embedded genetic algorithm for simplified protein structure prediction. In  
763 *2013 IEEE Congress on Evolutionary Computation*, pages 1091–1098, June  
764 2013.
- 765 [36] Timothy J Richmond. Solvent accessible surface area and excluded volume  
766 in proteins: Analytical equations for overlapping spheres and implications  
767 for the hydrophobic effect. *J. Mol. Biol.*, 178(1):63–89, 1984.
- 768 [37] Carol A. Rohl, Charlie E.M. Strauss, Kira M.S. Misura, and David Baker.  
769 Protein structure prediction using rosetta. In *Numerical Computer Meth-*  
770 *ods, Part D*, volume 383 of *Methods in Enzymology*, pages 66 – 93. Academic  
771 Press, 2004.

- 772 [38] Ruslan Sadreyev and Nick Grishin. Compass: a tool for comparison of mul-  
773 tiple protein alignments with assessment of statistical significance. *Journal*  
774 *of molecular biology*, 326(1):317–336, 2003.
- 775 [39] Claudio Sanhueza, Francia Jimenez, Regina Berretta, and Pablo Moscato.  
776 Pasmogap: A parallel asynchronous memetic algorithm for solving the  
777 multi-objective quadratic assignment problem. In *2017 IEEE Congress on*  
778 *Evolutionary Computation, CEC 2017, Donostia, San Sebastián, Spain,*  
779 *June 5-8, 2017*, pages 1103–1110. IEEE, 2017.
- 780 [40] Schrödinger, LLC. The PyMOL molecular graphics system, version 1.7,  
781 Nov 2015.
- 782 [41] N. Siew and D. Fischer. Convergent evolution of protein structure predic-  
783 tion and computer chess tournaments: Casp, kasparov, and cafasp. *IBM*  
784 *Systems Journal*, 40(2):410–425, 2001.
- 785 [42] J. E. Smith. *The Co-Evolution of Memetic Algorithms for Protein Structure*  
786 *Prediction*, pages 105–128. Springer Berlin Heidelberg, Berlin, Heidelberg,  
787 2005.
- 788 [43] Johannes Söding, Andreas Biegert, and Andrei N Lupas. The hhpred in-  
789 teractive server for protein homology detection and structure prediction.  
790 *Nucleic acids research*, 33(suppl.2):W244–W248, 2005.
- 791 [44] Ching-Wai Tan and David Jones. Using neural networks and evolutionary  
792 information in decoy discrimination for protein tertiary structure predic-  
793 tion. *BMC bioinformatics*, 9:94, 02 2008.
- 794 [45] Dong Xu and Yang Zhang. Ab initio protein structure assembly using  
795 continuous structure fragments and optimized knowledge-based force field.  
796 *Proteins: Structure, Function, and Bioinformatics*, 80(7):1715–1735, 2012.
- 797 [46] Wenxuan Zhang, Jianyi Yang, Baoji He, Sara Elizabeth Walker, Hongjiu  
798 Zhang, Brandon Govindarajoo, Jouko Virtanen, Zhidong Xue, Hong-Bin  
799 Shen, and Yang Zhang. Integration of quark and i-tasser for ab initio  
800 protein structure prediction in casp11. *Proteins: Structure, Function, and*  
801 *Bioinformatics*, 2015.
- 802 [47] Yang Zhang. I-tasser server for protein 3d structure prediction. *BMC*  
803 *bioinformatics*, 9(1):1, 2008.
- 804 [48] Yang Zhang and Jeffrey Skolnick. Scoring function for automated assess-  
805 ment of protein structure template quality. *Proteins: Structure, Function,*  
806 *and Bioinformatics*, 57(4):702–710, 2004.
- 807 [49] Wei Zheng, Chengxin Zhang, Eric W. Bell, and Yang Zhang. I-tasser  
808 gateway: A protein structure and function prediction server powered by  
809 xsede. *Future Generation Computer Systems*, 99:73 – 85, 2019.

Sequence	Structure	RMSD (Å)		GDT_TS (%)		Energy (Kcal/mol)	Total time (s)	Generations
		Best	Average	Best	Average			
1K43	DMA-Original	<b>0.71</b>	1.52 (±0.93)	<b>85.71</b>	77.32 (±8.83)	655.8 (±35.7)	2554.1 (±77.9)	286.6 (±6.0)
	DMA-Final	0.76	<b>1.28 (±0.36)</b>	83.93	<b>79.11 (±2.92)</b>	576.9 (±17.3)	2835.6 (±54.3)	374.6 (±9.1)
	Rosetta	*0.64	*0.93 (±0.15)	85.71	*80.54 (±2.58)	116.43 (±73.2)	-	-
1L2Y	DMA-Original	1.14	1.75 (±0.54)	85.00	75.12 (±6.57)	846.2 (±56.6)	3486.4 (±160.2)	191.6 (±3.7)
	DMA-Final	<b>0.86</b>	<b>1.39 (±0.44)</b>	<b>87.50</b>	<b>78.75 (±6.51)</b>	709.6 (±27.4)	3833.0 (±51.1)	250.9 (±2.9)
	Rosetta	*0.62	1.48 (±0.45)	*96.25	*80.75 (±7.52)	-26.04 (±5.51)	-	-
2MTW	DMA-Original	0.75	<b>1.02 (±0.23)</b>	<b>78.75</b>	<b>75.38 (±1.77)</b>	914.8 (±35.6)	3813.4 (±78.5)	190.8 (±2.9)
	DMA-Final	<b>0.67</b>	1.19 (±0.74)	77.50	72.62 (±4.73)	786.4 (±19.1)	4138.1 (±57.5)	244.7 (±3.3)
	Rosetta	3.9	4.4 (±0.36)	62.5	59.75 (±2.61)	-18.28 (±2.27)	-	-
1WQC	DMA-Original	3.48	4.77 (±1.01)	64.42	58.17 (±5.48)	1136.2 (±28.3)	4423.5 (±81.3)	144.4 (±2.5)
	DMA-Final	<b>2.34</b>	<b>3.54 (±0.74)</b>	<b>67.31</b>	<b>59.62 (±4.51)</b>	917.6 (±24.4)	4904.3 (±65.3)	188.7 (±3.7)
	Rosetta	*1.82	*2.19 (±0.26)	*76.92	*72.31 (±2.85)	-27.6 (±6.58)	-	-
1ACW	DMA-Original	3.17	4.67 (±1.21)	61.21	54.40 (±4.31)	1198.7 (±69.1)	4173.0 (±182.1)	125.0 (±1.4)
	DMA-Final	<b>2.47</b>	<b>3.87 (±1.16)</b>	<b>67.24</b>	<b>58.02 (±6.04)</b>	993.5 (±28.9)	4696.9 (±94.3)	163.4 (±2.2)
	Rosetta	*1.66	*2.52 (±1.15)	*77.59	*70.95 (±6.94)	-26.0 (±3.62)	-	-
1DFN	DMA-Original	5.19	8.03 (±2.71)	54.17	46.58 (±4.66)	1417.2 (±59.5)	4527.4 (±142.2)	120.2 (±2.1)
	DMA-Final	<b>4.93</b>	<b>6.02 (±1.28)</b>	<b>58.33</b>	<b>49.92 (±5.31)</b>	1126.9 (±50.7)	5168.7 (±91.1)	158.7 (±1.5)
	Rosetta	*3.63	*5.36 (±0.79)	49.17	44.67 (±2.67)	79.04 (±233.0)	-	-
1Q2K	DMA-Original	4.34	5.89 (±0.91)	60.48	55.48 (±2.96)	1232.1 (±56.0)	4836.4 (±82.0)	116.7 (±1.8)
	DMA-Final	<b>3.22</b>	<b>4.25 (±0.70)</b>	<b>65.32</b>	<b>60.40 (±3.87)</b>	1033.1 (±38.8)	5418.7 (±82.9)	151.7 (±3.3)
	Rosetta	*0.57	*1.98 (±0.97)	*94.35	*78.79 (±10.34)	-30.73 (±5.8)	-	-
2F4K	DMA-Original	4.55	5.48 (±0.78)	63.64	<b>52.58 (±6.27)</b>	1323.1 (±67.3)	6272.4 (±117.8)	108.5 (±1.6)
	DMA-Final	<b>3.17</b>	<b>4.90 (±1.33)</b>	<b>69.70</b>	52.12 (±8.87)	1085.6 (±33.7)	6800.5 (±90.6)	142.5 (±2.9)
	Rosetta	5.02	5.79 (±0.43)	-	-	-39.42 (±5.61)	-	-
2MR9	DMA-Original	4.79	7.70 (±1.42)	54.55	46.76 (±5.52)	1738.2 (±109.1)	8075.3 (±243.7)	80.0 (±1.7)
	DMA-Final	<b>2.21</b>	<b>5.80 (±1.83)</b>	<b>67.05</b>	<b>51.59 (±7.13)</b>	1342.1 (±59.1)	8690.1 (±217.8)	104.2 (±2.2)
	Rosetta	*1.43	*2.27 (±0.65)	*83.52	*73.13 (±6.23)	-72.25 (±3.59)	-	-
2P81	DMA-Original	4.55	6.26 (±1.46)	<b>64.20</b>	52.22 (±7.15)	2110.3 (±118.1)	9196.3 (±371.9)	79.6 (±1.3)
	DMA-Final	<b>2.95</b>	<b>4.03 (±0.62)</b>	<b>64.20</b>	<b>55.68 (±5.10)</b>	1562.3 (±62.8)	10243.7 (±241.3)	106.4 (±2.5)
	Rosetta	5.61	6.85 (±0.78)	36.93	34.09 (±1.05)	-63.86 (±5.26)	-	-
3P7K	DMA-Original	2.42	3.46 (±0.86)	72.78	65.67 (±3.52)	2105.9 (±137.7)	9139.0 (±132.2)	78.0 (±1.2)
	DMA-Final	<b>1.65</b>	<b>2.89 (±1.54)</b>	<b>81.11</b>	<b>72.28 (±9.81)</b>	1848.2 (±21.3)	9585.2 (±230.4)	104.7 (±3.2)
	Rosetta	15.04	15.94 (±0.49)	54.44	48.39 (±4.68)	-65.42 (±5.58)	-	-
1AB1	DMA-Original	4.91	8.72 (±1.67)	52.17	43.15 (±4.61)	1770.7 (±63.6)	6841.7 (±142.9)	76.7 (±1.1)
	DMA-Final	<b>2.96</b>	<b>6.14 (±2.22)</b>	<b>64.67</b>	<b>52.66 (±8.66)</b>	1349.2 (±50.4)	7805.4 (±158.9)	100.9 (±1.5)
	Rosetta	3.45	*5.37 (±1.09)	*67.93	*58.1 (±5.24)	-44.21 (±5.23)	-	-
3V1A	DMA-Original	6.22	10.07 (±2.83)	50.00	46.67 (±1.88)	2129.6 (±120.7)	8635.9 (±144.3)	73.3 (±1.2)
	DMA-Final	<b>1.65</b>	<b>3.23 (±0.90)</b>	<b>78.12</b>	<b>61.77 (±8.47)</b>	1604.3 (±40.9)	9500.4 (±174.6)	97.4 (±1.9)
	Rosetta	*0.7	*2.18 (±1.82)	55.21	52.09 (±4.35)	-84.29 (±8.36)	-	-
2P6J	DMA-Original	4.66	8.41 (±2.30)	48.56	43.90 (±4.09)	2551.8 (±204.7)	10879.1 (±324.0)	66.7 (±1.5)
	DMA-Final	<b>3.59</b>	<b>5.40 (±2.02)</b>	<b>56.73</b>	<b>51.83 (±4.35)</b>	1910.1 (±84.2)	11940.8 (±266.4)	90.2 (±2.0)
	Rosetta	*2.25	*3.37 (±1.06)	*74.52	*62.93 (±5.57)	-82.89 (±7.51)	-	-
1ENH	DMA-Original	3.54	9.22 (±3.14)	59.72	45.33 (±6.45)	2549.4 (±142.2)	11434.3 (±364.4)	63.6 (±1.4)
	DMA-Final	<b>2.81</b>	<b>5.07 (±2.57)</b>	<b>65.74</b>	<b>55.19 (±8.87)</b>	1939.8 (±63.0)	12554.4 (±333.4)	87.4 (±2.0)
	Rosetta	*1.7	*2.88 (±0.94)	46.3	43.61 (±1.22)	-86.32 (±7.47)	-	-
2JUC	DMA-Original	7.43	12.26 (±3.67)	38.18	32.91 (±3.11)	2726.6 (±223.8)	10993.5 (±262.1)	64.4 (±1.3)
	DMA-Final	<b>5.84</b>	<b>9.27 (±1.55)</b>	<b>45.45</b>	<b>37.86 (±4.33)</b>	1956.0 (±89.7)	11962.2 (±309.6)	84.4 (±1.8)
	Rosetta	6.1	10.23 (±1.56)	42.8	35.97 (±3.4)	-47.09 (±7.63)	-	-
1ROP	DMA-Original	2.73	13.32 (±5.53)	68.75	49.51 (±7.62)	2447.2 (±179.1)	10423.3 (±197.5)	62.0 (±1.2)
	DMA-Final	<b>2.68</b>	<b>3.94 (±1.00)</b>	<b>70.54</b>	<b>60.00 (±7.43)</b>	1904.4 (±26.8)	11369.7 (±235.0)	81.9 (±2.2)
	Rosetta	*1.98	4.98 (±2.6)	*81.25	*65.49 (±12.42)	-88.13 (±11.25)	-	-
2P5K	DMA-Original	6.42	<b>8.74 (±1.92)</b>	44.44	37.22 (±4.10)	2847.2 (±164.7)	12212.0 (±294.5)	55.1 (±0.9)
	DMA-Final	<b>4.35</b>	8.79 (±2.68)	<b>54.37</b>	<b>43.13 (±5.32)</b>	2190.3 (±101.9)	13456.7 (±260.6)	74.6 (±1.8)
	Rosetta	*1.53	*2.21 (±0.92)	53.97	*51.67 (±1.74)	-105.6 (±4.02)	-	-
1AIL	DMA-Original	12.32	16.43 (±3.64)	42.86	37.57 (±3.28)	3127.8 (±137.4)	13347.8 (±249.9)	49.8 (±0.6)
	DMA-Final	<b>8.26</b>	<b>9.46 (±0.94)</b>	<b>53.93</b>	<b>41.50 (±4.82)</b>	2319.0 (±101.2)	14720.4 (±248.9)	66.3 (±1.3)
	Rosetta	*6.85	*9.34 (±0.97)	48.93	38.96 (±4.92)	-97.26 (±6.84)	-	-
2PMR	DMA-Original	11.60	17.66 (±3.32)	48.03	39.34 (±5.93)	3575.7 (±171.3)	14426.0 (±318.3)	45.5 (±1.0)
	DMA-Final	<b>2.02</b>	<b>8.45 (±3.05)</b>	<b>72.37</b>	<b>47.11 (±10.98)</b>	2644.2 (±82.0)	16238.0 (±463.3)	62.1 (±1.3)
	Rosetta	2.46	*3.85 (±0.78)	45.39	40.2 (±3.33)	-129.86 (±6.4)	-	-

Table 6: Results of the experiments that compare the original, the latest implementation of the DMA-3DPSP, and Rosetta. DMA-Final combines the approach with the best results of the previous experiments. For each one of the 20 target peptides, we show the results of the best solution and the average of the tests for both RMSD values and GDT percentage vs experimental structures, as well as the average results of free energy, runtime and generations performed. The values in bold represent the best results between the first two implementations. The (\*) denotes the case studies where Rosetta outperformed all of the others.

Sequence	Structure	Steric clashes	Unsatisfied hydrogen bonds		SASA	Void volume
			% in shell	% in core		
<b>Gaia Target</b>		<b>0.02</b>	<b>9.56</b>	<b>1.45</b>	<b>221.64</b>	<b>0.0097</b>
1K43	DMA-Original	0.13	11.30	0.00	208.25	0.00
	DMA-Final	0.10	18.10	0.00	203.04	0.00
1L2Y	DMA-Original	0.14	17.60	0.00	190.27	0.00
	DMA-Final	0.07	9.80	0.00	174.48	0.00
2MTW	DMA-Original	0.14	10.40	0.00	234.52	0.00
	DMA-Final	0.08	17.90	0.00	208.65	0.00
1WQC	DMA-Original	0.14	12.60	0.00	220.85	0.00
	DMA-Final	0.13	13.90	0.00	192.67	0.00
1ACW	DMA-Original	0.19	19.30	0.00	208.09	1.30
	DMA-Final	0.11	19.30	0.00	194.27	0.89
1DFN	DMA-Original	0.10	15.90	0.00	255.56	0.00
	DMA-Final	0.06	17.00	0.00	212.61	1.75
1Q2K	DMA-Original	0.11	14.90	0.00	220.23	0.00
	DMA-Final	0.07	20.60	0.00	192.60	1.58
2F4K	DMA-Original	0.12	16.40	0.00	238.06	0.00
	DMA-Final	0.08	17.50	0.00	202.65	2.48
2MR9	DMA-Original	0.14	11.80	0.00	227.07	0.00
	DMA-Final	0.09	7.60	1.60	178.10	0.51
2P81	DMA-Original	0.18	13.80	0.00	271.61	0.00
	DMA-Final	0.10	13.10	0.00	217.16	0.00
3P7K	DMA-Original	0.18	6.00	0.00	259.35	0.00
	DMA-Final	0.10	9.80	0.00	263.35	0.00
1AB1	DMA-Original	0.20	14.80	0.00	211.79	1.37
	DMA-Final	0.09	14.80	0.00	173.52	1.12
3V1A	DMA-Original	0.10	6.20	0.00	253.64	0.00
	DMA-Final	0.09	8.20	0.00	208.14	0.00
2P6J	DMA-Original	0.17	13.90	0.00	261.18	0.00
	DMA-Final	0.09	9.30	0.00	221.77	0.00
1ENH	DMA-Original	0.14	11.30	0.00	250.69	0.00
	DMA-Final	0.08	6.80	0.00	217.13	1.26
2JUC	DMA-Original	0.16	16.10	0.00	273.61	0.00
	DMA-Final	0.11	15.60	0.50	220.80	4.78
1ROP	DMA-Original	0.10	9.20	0.00	229.96	0.00
	DMA-Final	0.09	5.70	0.00	223.51	0.00
2P5K	DMA-Original	0.19	9.40	0.00	254.14	0.18
	DMA-Final	0.11	12.00	0.00	224.20	0.81
1AIL	DMA-Original	0.19	13.50	0.00	265.12	0.00
	DMA-Final	0.08	9.10	0.00	211.24	1.09
2PMR	DMA-Original	0.16	7.30	0.00	257.61	0.33
	DMA-Final	0.11	11.20	0.00	219.04	5.38

Table 7: Results of the structure quality test of *Gaia Web Tool* for the best solutions found by the two DMA-3DPSP implementations tested in the final experiments. For each peptide and DMA implementation, we show the values of steric clashes, unsatisfied hydrogen bonds, solvent accessible surface area and void volume from the Gaia report. The values that exceed the optimum margins in Gaia report are presented with red background.

CASP ID	PDB ID	Size	Folding pattern
T0769	2MQ8	112	2 $\alpha$ -helix, 4 $\beta$ -sheet
T0773	2N2U	77	2 $\alpha$ -helix, 1 $3_{10}$ helix, 4 $\beta$ -sheet
T0829	4RGI	70	2 $\alpha$ -helix, 5 $\beta$ -sheet
T0837	2MQD	128	4 $\alpha$ -helix, 6 $\beta$ -sheet
T0855	5TF3	119	7 $\alpha$ -helix, 1 $3_{10}$ helix

Table 8: List of protein sequences from CASP 11 targets used in computational tests.

Sequence	RMSD ( $\text{\AA}$ )				GDT_TS (%)			
	Best	Average	Zhang	Robetta	Best	Average	Zhang	Robetta
T0769 (2MQ8)	13.503	19.286 ( $\pm 4.365$ )	1.739	1.963	29.69	25.91 ( $\pm 2.33$ )	59.80	53.35
T0773 (2N2U)	8.961	12.593 ( $\pm 2.161$ )	1.180	2.268	40.26	33.44 ( $\pm 3.90$ )	90.67	78.73
T0829 (4RGI)	11.560	12.727 ( $\pm 1.131$ )	6.271	6.461	29.78	24.45 ( $\pm 2.83$ )	64.92	62.31
T0837 (2MQD)	15.014	17.958 ( $\pm 3.384$ )	7.847	8.506	22.90	19.18 ( $\pm 2.02$ )	44.56	39.78
T0855 (5TF3)	13.723	16.749 ( $\pm 2.238$ )	9.271	11.902	30.66	23.93 ( $\pm 2.71$ )	38.84	25.00

Table 9: Comparison of the last version of the DMA-3DPSP method with the CASP 11 protein targets. For each peptide, we show the best solution found and the average of 10 runs regarding RMSD and GDT. We also present the results of RMSD and GDT obtained by Zhang-Server and Robetta Server in the CASP 11 experiments.