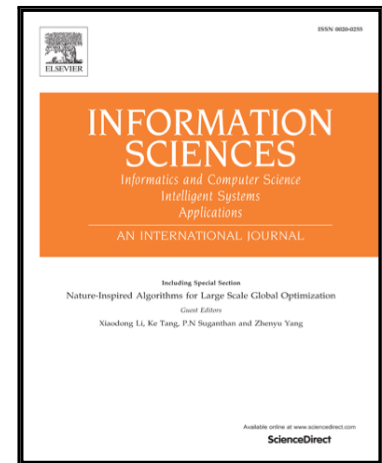


## Accepted Manuscript

### A Multiobjective Adaptive Approach for the Inference of Evolutionary Relationships in Protein-Based Scenarios

Sergio Santander-Jiménez, Miguel A. Vega-Rodríguez,  
Leonel Sousa

PII: S0020-0255(19)30130-6  
DOI: <https://doi.org/10.1016/j.ins.2019.02.020>  
Reference: INS 14290



To appear in: *Information Sciences*

Received date: 28 September 2018  
Revised date: 10 February 2019  
Accepted date: 11 February 2019

Please cite this article as: Sergio Santander-Jiménez, Miguel A. Vega-Rodríguez, Leonel Sousa, A Multiobjective Adaptive Approach for the Inference of Evolutionary Relationships in Protein-Based Scenarios, *Information Sciences* (2019), doi: <https://doi.org/10.1016/j.ins.2019.02.020>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# A Multiobjective Adaptive Approach for the Inference of Evolutionary Relationships in Protein-Based Scenarios

Sergio Santander-Jiménez<sup>a,\*</sup>, Miguel A. Vega-Rodríguez<sup>b</sup>, Leonel Sousa<sup>a</sup>

<sup>a</sup>*Instituto de Engenharia de Sistemas e Computadores - Investigação e Desenvolvimento em Lisboa (INESC-ID), Instituto Superior Técnico, Universidade de Lisboa, Lisboa, 1000-029, Portugal*

<sup>b</sup>*Instituto de Investigación en Tecnologías Informáticas Aplicadas de Extremadura (INTIA), Universidad de Extremadura, Avda. de la Universidad s/n, Cáceres 10003, Spain*

---

## Abstract

Complex optimization problem solving is a constant issue in a wide range of scientific domains. Robust bioinspired procedures with accurate search capabilities are therefore required to address the challenge that such optimization problems represent. This work explores different design alternatives for the metaheuristic Multiobjective Shuffled Frog-Leaping Algorithm, a novel method that combines parallel searches and swarm-based operators to undertake the processing of complex search spaces. Three variants of the metaheuristic are adopted: a dominance-based approach, an indicator-based alternative, and an adaptive proposal that incorporates both multiobjective strategies (dynamically assigning during the execution more resources to the most successful strategy). The performance of the proposed designs is examined when tackling, as a case study, the inference of ancestral relationships from protein data, using different multiobjective metrics and bio-statistical testing procedures. Experimental results show the additional robustness that the adaptive technique provides to the metaheuristic, allowing its search engine to exploit the most fitting multiobjective approach according to the status of the optimization process.

---

\*Corresponding author.

*Email addresses:* `sergio.jimenez@tecnico.ulisboa.pt`  
(Sergio Santander-Jiménez), `mavega@unex.es` (Miguel A. Vega-Rodríguez),  
`leonel.sousa@ist.utl.pt` (Leonel Sousa)

*Keywords:* Bioinspired computing, multiobjective optimization, adaptive algorithms, bioinformatics.

## 1. Introduction

Bioinspired computing represents one of the most promising research directions to deal with NP-hard optimization problems. In the current context, there is special emphasis on solving this kind of problems in accordance with formulations of multiobjective nature, which combine information from different criteria or datasets. A multiobjective optimization problem (MOP) is aimed at finding those solutions  $s = [s_1, s_2, \dots, s_k]$ , characterized by  $k$  decision variables, that optimize  $\eta \geq 2$  objective functions  $\vec{f}(s) = [f_1(s), f_2(s), \dots, f_\eta(s)]$ . The set of all possible solutions to the problem is denoted as the decision space  $S$  and the image of these solutions under the objective functions gives rise to the objective space  $Z = \mathbb{R}^\eta$ . The optimization process in a MOP consists in exploring  $S$  for those solutions that represent the best trade-offs in  $Z$ , commonly designated as Pareto-optimal solutions. In spite of the current interest in solving MOPs, these problems involve important complexity factors [8] from the perspective of problem hardness (e.g. exponentially growing decision spaces) as well as from the time perspective (time-consuming objective functions and computationally demanding operations). Such factors play a key role when real-world MOPs are involved, therefore demanding robust, accurate strategies to address them.

Different algorithmic trends can be found in the design of multiobjective optimization procedures. Among them, the dominance-based [10] and indicator-based [45] approaches represent two of the most commonly adopted strategies for solving MOPs. The first class is based on the concept of Pareto dominance, which stipulates conditions to compare solutions under a multiobjective framework. Given two solutions  $s_1$  and  $s_2$ , the Pareto dominance specifies that  $s_1$  dominates  $s_2$  ( $s_1 \succ s_2$ ) iff 1)  $\forall i \in [1, 2, \dots, \eta]$ ,  $f_i(s_1)$  is not worse than  $f_i(s_2)$  and 2)  $\exists i \in [1, 2, \dots, \eta]$  such that  $f_i(s_1)$  is better than  $f_i(s_2)$ . Dominance-based algorithms employ this relationship to classify the solutions handled during the optimization process, while also adopting measurements of solution density information as a complementary criterion. Regarding indicator-based designs, the quality of the solutions managed in these approaches is examined by integrating multiobjective quality metrics into the

search engine. Such metrics, known as quality indicators, allow the mapping of one or several sets of solutions  $S' \subseteq S$  to a real number  $I(S') \in \mathbb{R}$  that measures multiobjective quality properties of the solutions under evaluation. The application of these two multiobjective strategies in different optimization scenarios presents divergent performance in the case of both benchmark [43, 48] and real-world problems [28, 44]. An open question is therefore how to accurately use the information provided by these different multiobjective mechanisms to improve search capabilities, along with their progressive, dynamic adaptation to the status of the optimization process when difficult MOPs are considered.

This work is focused on the study of multiobjective strategies to improve the Multiobjective Shuffled Frog-Leaping Algorithm (MO-SFLA). MO-SFLA is a novel algorithmic approach proposed in [35] that combines swarm techniques and parallel searches for solutions. A comparative study between dominance-based and indicator-based versions of the proposal is undertaken, with the aim of identifying how the multiobjective behaviour of the algorithm changes according to the design variant under consideration. Furthermore, an additional alternative based on the concept of adaptation is proposed, in order to carry out parallel searches dynamically in accordance with the strategy (dominance or indicator) that generates better solution quality throughout the execution of the algorithm. These MO-SFLA variants will be examined on one of the most challenging problems in the bioinformatics research field: the inference of ancestral relationships from protein sequences [47]. The performance of each variant will be assessed through experimentation over real-world amino acid datasets, carrying out the evaluation of the obtained solutions attending to different multiobjective and bio-statistical quality metrics.

The main novelty and contributions of this paper are as follows. Using as a baseline the dominance-based version from [35], an alternative indicator-based design for MO-SFLA is first introduced, replacing the previous dominance-based mechanisms by fitness assessment procedures based on multiobjective quality indicator information. On the basis of these initial design alternatives, the question on how to take advantage of both strategies simultaneously (dominance-based and indicator-based techniques) is addressed. To this end, a new design variant for MO-SFLA is proposed, where both dominance-based and indicator-based strategies are included and dynamically used under a novel adaptive approach. In this sense, the adaptive MO-SFLA design goes beyond a simple combination of previous techniques and is

not limited to just statically mixing information from different mechanisms as in traditional hybrid schemes. The significance of the adaptive design lies in the dynamic exploitation of the knowledge provided by different multiobjective strategies, putting more effort on the most satisfying strategy in each stage of the optimization process. Particularly, feedback on the success of the parallel searches conducted under each multiobjective strategy is retrieved at each generation in order to adapt the behaviour of MO-SFLA search engine, assigning at each step more resources (in terms of partitions of individuals) to the most successful strategy. In order to justify this approach, an in-depth analysis on the potential divergence in the Pareto fronts from the baseline dominance-based and indicator-based versions is conducted. The evaluation of the adaptive MO-SFLA is then undertaken in accordance with the insights of this previous analysis, verifying the relevance of the dynamic, adaptive approach attending to the quality of the generated Pareto fronts.

Summarizing, the key contributions can be classified in the following way:

- Introduction of an indicator-based multiobjective design for the meta-heuristic MO-SFLA;
- Proposal of a novel adaptive design for MO-SFLA that incorporates different multiobjective strategies (dominance and indicator) to manage parallel searches, dynamically adapting them in accordance with the performance attained during the execution of the algorithm;
- Comparative evaluation of the three design alternatives of MO-SFLA in hard-to-solve real-world scenarios, tackling as a case study the inference of phylogenetic relationships from challenging protein data;
- Analysis of multiobjective behaviour patterns for each design variant, verifying the effect of each strategy over the attained Pareto fronts and examining the relevance of the adaptive design to maximize multiobjective quality;
- Evaluation of the quality of the attained solutions by undertaking multiobjective and bio-statistical analyses with two reference multiobjective algorithms and six state-of-the-art biological tools.

This paper is organized as follows. The next section summarizes relevant works on the topics of this research. The formulation of the tackled problem

is detailed in Section 3, while Section 4 describes the proposed MO-SFLA approaches. Section 5 explains the experimental methodology followed in this study, along with reporting and discussing the obtained results. Finally, Section 6 draws conclusions and future lines of research.

## 2. Related Work

The proposal of adaptive bioinspired methods is becoming an increasingly adopted approach to solve difficult problems in different research fields. This section reviews relevant works on the two main topics that represent the scope of this paper: the use of adaptive approaches (that is, methods that dynamically change their behaviour during the execution) in multiobjective optimization and bioinspired techniques to infer ancestral relationships.

In the multiobjective context, adaptive methods have mostly been used for the inner control and selection of parameter values and/or evolutionary operators, in accordance with the performance feedback gathered throughout the execution of the algorithm. Tan *et al.* proposed an adaptive control method for crossover and mutation coordination, defining an operator that synergized these two operations according to the status of the optimization process [42]. For an application of architectural design, Bittermann and Sariyildiz studied adaptive relaxation of the Pareto dominance to classify solutions during the search process [4]. Charlet *et al.* proposed a hybrid genetic algorithm including up to three different fitness definitions and adaptive control of the selection and mutation rate values [6]. In addition, Hadka and Reed introduced a framework for multiobjective evolutionary computation that incorporates adaptive techniques to select the most accurate recombination operations while also controlling population and tournament sizes [16]. Later on, Zavoianu *et al.* included adaptive allocation of fitness evaluations among populations in the co-evolutionary proposal DECMO2 [49]. Adaptive fitness schemes for multiobjective approaches, based on the weighted sum of objective functions, were studied by Wang *et al.* [46], while Jiang *et al.* addressed the adaptive control of the kappa parameter in the IBEA algorithm by using the simplex method [19]. More recently, Azzouz *et al.* proposed an adaptive algorithm based on three population management strategies to handle dynamic MOPs [2]. In 2018, Lin *et al.* integrated mechanisms for the adaptive selection of differential evolution operators into a multiobjective immune-inspired proposal [21].

Regarding the reconstruction of evolutionary relationships, initial studies already pointed out the additional hardness associated to the consideration of protein sequence data in this problem [25]. Matsuda proposed the use of bioinspired computing to address this issue, designing a genetic algorithm for carrying out phylogenetic analyses of 17 EF-1 $\alpha$  protein sequences under the maximum likelihood criterion [24]. Reijmers *et al.* investigated phylogenetic relationships for sequences of 37 G protein-coupled receptors by means of a genetic algorithm with distance-based individual representation [31]. Hill *et al.* also examined the performance of genetic algorithms for the analysis of protein sequences under the maximum parsimony criterion [18]. More recent research was aimed at proposing robust approaches with improved search capabilities to tackle complex phylogenetic analyses on current protein-based datasets. In this sense, Stamatakis integrated simulated annealing optimization and models of protein evolution into the biological software RAxML [39, 40], which represents one of the most commonly used tools in phylogenetics along with other approaches such as TNT [14], IQ-TREE [29] and MrBayes [33]. Combinations of search strategies were also adopted in GARLI, a genetic algorithm hybridized with a variety of local search operators [3]. Focusing on adaptive techniques, the most well-known approach was due to Skourikhine, who reported a self-adaptive genetic algorithm for the analysis of DNA data [38]. Guo *et al.* proposed the use of adaptive ant colony optimization for the protein case, showing better efficiency than classic neighbour-joining techniques [15].

In the last years, research efforts have put emphasis on tackling multi-objective formulations of the problem to solve the conflicts that arise from different optimality criteria [34] and datasets with divergent phylogenetic evidence [30]. Proposals like PhyloMOEA [5] and omni-aiNet [7] represent two well-known evolutionary approaches to deal with multiobjective phylogenetic analyses, yet only in nucleotide scenarios. An initial, dominance-based design of MO-SFLA was proposed in [35] to address multiobjective reconstructions from protein sequence data, according to the likelihood and parsimony criteria. This paper aims to go a step further in this research direction, proposing first an alternative indicator-based design and detailing a novel MO-SFLA approach where both dominance and indicator-based strategies are dynamically used under adaptive techniques.

### 3. Problem Formulation

Phylogenetic analyses establish evolutionary relationships among natural organisms by processing data about genetic and molecular characteristics [47]. Through the study of divergence and similarities in such characteristics, ancestral relationships can be inferred and illustrated by using a tree-shaped data structure  $T = (V, E)$ , known as phylogenetic tree. The node set  $V$  in  $T$  locates in the leaf nodes those organisms whose molecular features are available at the input of the procedure, while internal nodes are used to represent hypothetical ancestors. The phylogenetic topology is organized in accordance with the branch set  $E$ , which defines ancestor-descendant linkages between related organisms. The biological data to be processed during the inference process is given by a multiple sequence alignment of size  $N \times M$ , where  $N$  refers to the number of sequences and  $M$  is the sequence length. Each position in the sequences is usually denoted as ‘character’, while the value contained in that position is known as ‘state’. The possible states that a character can take are given by an alphabet  $\Lambda$  that corresponds to the amino acid state alphabet in the case of protein-based phylogenetic reconstructions.

The inference of ancestral relationships can be addressed as an optimization problem, which implies the exploration of the phylogenetic tree space to obtain those evolutionary hypotheses that optimize one or several optimality criteria. This work tackles a multiobjective formulation of the problem based on the simultaneous consideration of multiple criteria, an approach that has shown significant benefits in conflicting biological scenarios [34, 5, 7]. More specifically, the problem is formulated as a bi-objective MOP based on the parsimony  $P(T)$  and likelihood  $L(T)$  criteria:

$$\begin{aligned} & \text{optimize } \vec{f}(T) = \{f_1(T), f_2(T)\}, \\ & \text{where } f_1(T) = \text{minimize } P(T), \\ & \quad \quad \quad f_2(T) = \text{maximize } L(T). \end{aligned} \tag{1}$$

The parsimony criterion aims to minimize the number of state changes in the sequences that belong to related nodes. This objective function returns an integer value quantifying the amount of changes or mutations observed in the phylogenetic topology. Formally, the parsimony score  $P(T)$  (also known as parsimony length) for a phylogenetic tree  $T = (V, E)$  that models the



evolution of the organisms characterized in the alignment is expressed as:

$$P(T) = \sum_{i=1}^M \sum_{(u,v) \in E} C(u_i, v_i), \quad (2)$$

where  $(u, v) \in E$  refers to the branch between two nodes  $u, v \in V$ ,  $u_i, v_i \in \Lambda$  are the states at the  $i$ th character of the sequences for  $u$  and  $v$ , and  $C(u_i, v_i)$  quantifies if a mutation event between  $u_i$  and  $v_i$  has taken place ( $C(u_i, v_i) = 1$ ) or not ( $C(u_i, v_i) = 0$ ). In order to carry out parsimony calculations, it is first required the assignment of character states for each node in  $V$ . For a leaf node  $l \in V$ , these states are given by the sequence in the alignment that corresponds to the organism  $l$ . For an internal node  $u \in V$  with children  $v, w$ , a set of potential character states  $A_i(u)$  is computed, for each character  $i = 1$  to  $M$ , such that:

$$A_i(u) = \begin{cases} A_i(v) \cap A_i(w) & \text{if } A_i(v) \cap A_i(w) \neq \emptyset, \\ A_i(v) \cup A_i(w) & \text{if } A_i(v) \cap A_i(w) = \emptyset. \end{cases} \quad (3)$$

Final character states are defined from these sets, starting from the root node  $r$  by choosing a random state among the ones contained in  $A_i(r)$  as the final state value  $r_i$ . For the remaining internal nodes e.g.  $v \in V$  with ancestor  $u$ , the final state  $v_i$  is given by  $u_i$  iff  $u_i$  is included in  $A_i(v)$ . Otherwise,  $v_i$  will be given by a randomly chosen state from  $A_i(v)$ .

Regarding the second objective, the likelihood criterion aims to maximize the probability of observing the evolutionary results expressed in the sequence alignment, given an evolutionary hypothesis described by a phylogenetic tree and a model of sequence evolution. The output of this objective function is a floating-point value that measures how likely the phylogenetic hypothesis under evaluation gave rise to the biological evidence observed in the input data. The likelihood score  $L(T)$  for a phylogenetic tree  $T = (V, E)$  is given by the following expression:

$$L(T) = \prod_{i=1}^M \sum_{x \in \Lambda} \pi_x L_p(r_i = x), \quad (4)$$

where  $\pi_x$  is the stationary probability of state  $x$  in the alphabet  $\Lambda$  and  $L_p(r_i = x)$  the partial likelihood of observing  $x$  at the  $i$ th character of the

root  $r \in V$ .  $L_p(u_i = x)$  for an internal node  $u$  with children  $v, w$  is given by:

$$L_p(u_i = x) = \left( \sum_{y \in \Lambda} P_{xy}(t_{uv}) L_p(v_i = y) \right) \times \left( \sum_{y \in \Lambda} P_{xy}(t_{uw}) L_p(w_i = y) \right), \quad (5)$$

where  $t_{uv}, t_{uw}$  refer to the evolutionary times between  $u$  and its children  $v, w$ , given by the length values of the branches  $(u, v), (u, w) \in E$ , and  $P_{xy}(t)$  the probability of observing a mutation event from the state  $x$  to  $y$  within a time  $t$ . For a leaf node  $l$ , its partial likelihood  $L_p(l_i = x)$  will be 1 if  $l_i = x$  or 0 otherwise. These partial likelihood calculations are carried out for each state  $x \in \Lambda$  prior to the application of Equation 4, being the obtained likelihood score usually reported in logarithmic scale (log-likelihood).

The reconstruction of evolutionary relationships is considered a grand computational challenge due to the exponential growth of the number of possible solutions with the number of sequences  $N$ . More specifically, the size of the phylogeny space is governed by the double factorial  $(2N - 5)!!$  [47], showing for  $N > 50$  a number of alternative solutions higher than the Eddington number. Furthermore, the study of this problem in protein alignments intensifies its time-consuming nature due to the consideration of a state alphabet involving 20 possible amino acids [32]. This issue has an impact in the calculations and data types involved in the evaluation procedures, whose temporal costs are also generally influenced by the sequence length. Nevertheless, phylogenetic analyses over protein sequence data are gaining increasing attention in the current research context. The higher resolution provided by protein data is useful when large evolutionary distances are involved in the analysis, also playing a key role in the tasks of gene family reconstruction, gene discovery, and function prediction [47, 32]. Consequently, the increasing demands of biological processing justify the proposal of novel robust, efficient methods to address the problem.

#### 4. MO-SFLA Approaches

This work studies different MO-SFLA design alternatives to tackle the reconstruction of evolutionary relationships from protein data. The general features of the metaheuristic and the proposed variants are described in this section, also detailing the introduction of adaptive strategies into MO-SFLA.

#### 4.1. General Features

MO-SFLA is a bioinspired metaheuristic proposed in [35] that adapts the Shuffled Frog-Leaping Algorithm (SFLA) [36] to the multiobjective context. The proposal combines properties from Particle Swarm Optimization and Shuffled Complex Evolution, in such a way that swarm-based learning techniques are integrated into multiple parallel searches for solutions. The main idea consists in exploring multiple directions of the solution space simultaneously through the parallel searches, which are implemented by partitioning the population into subsets of individuals, known as memeplexes. The definition of memeplexes is based on shuffling techniques that lead to an equitable distribution of promising individuals among the population partitions. At each generation, the processing of each memeplex involves a certain number of learning steps to generate new candidate solutions, which are obtained by sharing information within the memeplex via swarm strategies. Once the memeplexes have been processed, they are merged to spread the knowledge attained during the parallel searches, thus boosting the population evolution.

The individual representation introduced in MO-SFLA to tackle phylogenetic reconstructions is based on the concept of distance matrix. The search engine of the algorithm will operate over solutions codified by symmetric, matrix-shaped data structures  $\delta$  of size  $N \times N$ , where each entry  $\delta[x, y]$  contains a floating-point value representing the evolutionary distance between two organisms  $x$  and  $y$ . When initializing the population, starter distance matrices are generated from randomly selected topologies that are contained in a repository of 1,000 phylogenetic trees (generated from bootstrap samples of the input sequences). Each entry in the initial matrices is calculated as  $\delta[x, y] = \sum_{u,v \in PT_{x,y}} t_{uv}$ , where  $PT_{x,y}$  are the nodes contained in the path between  $x$  and  $y$  inside the topology, and  $t_{uv}$  is the length of the branch connecting the nodes  $u, v \in PT_{x,y}$ . As an indirect encoding strategy is employed, the matrix-shaped solutions processed throughout the execution of the algorithm must be mapped to the phylogenetic tree space. With this purpose in mind, the tree-building method BIONJ [47] is used to obtain the corresponding phylogenies.

It is worth remarking that, although some of the descriptions given in Section 3 employ rooted trees to make easier the comprehension of the tackled problem, the considered distance-based representation is independent from the rooted / unrooted tree form question. The search engine of the metaheuristic operates over an auxiliary distance matrix space, so the form of

the resulting trees depends on the tree-building method used for the matrix-phylogeny mapping. In the case of BIONJ, the reported trees are unrooted, which is the usual form employed in parsimony and likelihood optimization and, in general terms, when no molecular clock is assumed.

#### 4.2. Dominance-based and Indicator-based Designs

The general scheme for the dominance-based and indicator-based designs of MO-SFLA is shown in Algorithm 1. This metaheuristic examines at each generation the multiobjective quality of the *popSize* individuals contained in the population (line 3 of Algorithm 1). This operation represents the step in which the fundamental differences between the two considered design alternatives are introduced. On the one hand, multiobjective quality in the dominance-based version is quantified by using two well-known procedures from the NSGA-II algorithm [11]: fast non-dominated sort and crowding-based density estimation. Such procedures lead to the computation of Pareto rank and crowding distance values that allow the algorithm to distinguish the multiobjective quality of the solutions handled during the optimization process. On the other hand, the indicator-based version follows the guidelines established in the IBEA algorithm [50], which defines multiobjective fitness values based on the performance of each solution attending to a given multiobjective quality indicator. Particularly, the fitness for an individual  $P_i$  is given by the following expression:

$$\sum_{P_j \in P \setminus \{P_i\}} -e^{-I_{HD}(\{P_j\}, \{P_i\})/c\kappa}, \quad (6)$$

where  $P_j$  refers to any other individual in the population,  $I_{HD}$  the binary hypervolume indicator [50],  $c$  the maximum absolute indicator value, and  $\kappa$  a scaling factor = 0.05.

After assessing multiobjective quality, the population is partitioned into  $m$  memplexes (each one including  $n$  individuals, where  $n = \text{popSize}/m$ ). In this step (line 4), the shuffling and distribution of solutions among memplexes is carried out as follows: the first best individual  $P_1$  is sent to the first memplex  $Mem_1$ , the second best individual  $P_2$  is assigned to  $Mem_2$ ,  $P_m$  goes to  $Mem_m$ ,  $P_{m+1}$  to  $Mem_1$ , and so on. For the case of the dominance-based design, the best individuals are given by those solutions showing lower rank values and higher crowding distance within the same rank, while the indicator-based design gives priority to those individuals with higher indicator-based fitness values.

**Algorithm 1** General MO-SFLA Scheme

---

**Input parameters:**  $popSize$  (population size),  $m$  (number of memplexes),  $n$  ( $popSize/m$ , individuals per memplex),  $n_l$  (learning steps within a memplex),  $maxEval$  (stop criterion, maximum number of evaluations).

**Output:**  $PF$  (set of non-dominated solutions generated throughout the search).

- 1:  $P \leftarrow$  Set Initial Individuals ( $popSize$ ),  $PF \leftarrow \emptyset$
- 2: **while** ! stop criterion is reached ( $maxEval$ ) **do**
- 3: Multiobjective Fitness Assignment ( $P$ ,  $popSize$ ) /\* Rank and Crowding for Dominance, Hypervolume-based Fitness for Indicator \*/
- 4:  $\{Mem_1 \dots Mem_m\} \leftarrow$  Shuffling and Memplex Distribution ( $P$ ,  $popSize$ )
- 5: **for each** memplex  $Mem_i \in \{Mem_1 \dots Mem_m\}$  **do**
- 6:     **for**  $j = 1$  to  $n_l$  **do**
- 7:         **switch** ( $Mem_{i(n-j)}$ .counter)
- 8:             **case 0:**  $P'_{new} \leftarrow$  Learn from Best Local ( $Local$ ,  $Mem_{i(n-j)}$ )
- 9:             **case 1:**  $P'_{new} \leftarrow$  Learn from Best Global ( $Global$ ,  $Mem_{i(n-j)}$ )
- 10:            **case 2:**  $P'_{new} \leftarrow$  Local Search ( $Mem_{i(n-j)}$ )
- 11:            **if**  $P'_{new}$  improves  $Mem_{i(n-j)}$  ||  $Mem_{i(n-j)}$ .counter == 2 **then**
- 12:                  $Mem_{i(n-j)} \leftarrow P'_{new}$
- 13:                  $Mem_{i(n-j)}$ .counter  $\leftarrow 0$
- 14:             **else**
- 15:                  $Mem_{i(n-j)}$ .counter  $\leftarrow Mem_{i(n-j)}$ .counter + 1
- 16:             **end if**
- 17:         **end for**
- 18:     **end for**
- 19: Memplex Merging ( $P$ ,  $\{Mem_1 \dots Mem_m\}$ )
- 20: Pareto Front Update ( $PF$ ,  $P$ )
- 21: **end while**
- 22: **return**  $PF$

---

MO-SFLA conducts next the parallel searches implemented through the independent processing of memplexes (lines 5-18). Within each memplex  $Mem_i$ , different learning strategies are applied to generate new candidate solutions  $P'_{new}$ . The first one involves the generation of  $P'_{new}$  from an individual  $Mem_{ij}$  by using the information provided by the best local solution in the currently processed memplex (line 8):

$$D_{xy} = rand() * (Local.\delta[x, y] - Mem_{ij}.\delta[x, y]), \quad (7)$$

$$P'_{new}.\delta[x, y] = Mem_{ij}.\delta[x, y] + D_{xy}, \quad (8)$$

where  $Local$  refers to the best local individual in  $Mem_i$ ,  $\delta[x, y]$  the entries of the distance matrices associated to the involved individuals, and  $rand()$  a uniformly-distributed random number in the range  $[0,1]$ . In the dominance version, the best local individual is given by the solution in  $Mem_i$  with the best ranking and crowding distance values. For the indicator design, the solution in  $Mem_i$  with the highest indicator-based fitness is selected as

the best local individual. Since distance matrices are symmetric structures (with the main diagonal entries  $\delta[x, x] = 0$ ), Equations 7 and 8 are applied over  $\delta[x, y] \mid y < x$ , using the computed result to update both  $\delta[x, y]$  and  $\delta[y, x]$ . In addition, Equations 7 and 8 operate over positive distance values  $\delta[x, y]$ , thus guaranteeing that the final calculated distance  $P'_{new} \cdot \delta[x, y]$  is also positive. In this way, the generation of the corresponding phylogenetic tree via BIONJ leads to a feasible solution, whose quality is then evaluated by using the considered objective functions.

The second learning strategy involves the application of Equations 7 and 8 using the best global individual in the population (denoted as *Global* in line 9) instead of the best local one. In both versions, the identification of the best global individual follows the same principles as in the case of the best local one, defining the scope of the selection over the whole population instead of over the currently processed memplex. The third and last strategy is based on a local search procedure that applies topological rearrangements (subtree pruning-regrafting and nearest neighbour interchange) and gradient-based branch length optimization [47] over the phylogenetic tree associated to the currently processed individual (line 10).

The choice of the search strategy to be applied is governed by means of trial counters. Each individual  $Mem_{ij}$  is associated to a counter variable that stores the number of attempts in which the corresponding solution has not been improved. A search strategy is then selected according to the value of this counter: learning from the best local (counter = 0), learning from the best global (counter = 1), or local search (counter = 2). When the new candidate solution  $P'_{new}$  improves the one in  $Mem_{ij}$ ,  $Mem_{ij}$ .counter is reinitialized and the new solution accepted (lines 11-13). Otherwise,  $Mem_{ij}$ .counter is increased in order to make a new attempt with a different search strategy in the next generation (line 15). Whenever the local search is carried out, the new solution is accepted and the counter set to 0 in order to promote solution diversity.

The search for new candidate solutions in  $Mem_i$  is repeated for  $n_i$  learning steps, proceeding upon termination with the next memplex  $Mem_{i+1}$ . Once all the memplexes have been processed, they are merged to compose the population for the next generation (line 19). Finally, the Pareto front structure containing the non-dominated solutions found throughout the optimization process is updated (line 20).

### 4.3. Adaptive Design

As an additional design alternative, the adoption of adaptive techniques in MO-SFLA is herein proposed. The algorithmic design of MO-SFLA is naturally oriented towards the integration of adaptive procedures, as the memplexes defined by the metaheuristic can be dynamically managed under different multiobjective assessment approaches. The key idea consists in incorporating information from the two multiobjective strategies under study (dominance and indicator) to conduct the parallel searches, adjusting them adaptively attending to the performance achieved by each strategy during the optimization process. Starting from an equal distribution of memplexes ( $m/2$  managed by a dominance-based procedure and the remaining  $m/2$  by an indicator-based one), the algorithm retrieves feedback data about the success of each approach during the searches. By using this data, the memplexes will be re-distributed in an adaptive way, assigning more (and thereby a higher number of candidate solutions) to the most successful strategy.

The adaptive design proposed for MO-SFLA is presented in Algorithm 2. Once the population has been initialized, it is split into two sub-populations, each one comprising  $popSize/2$  randomly selected individuals (line 2 of Algorithm 2). The first sub-population is assessed and managed according to dominance-based techniques (fast non-dominated sort and crowding estimation), while the second one follows indicator-based principles (with indicator-based fitness calculations) (lines 5 and 6). In this version, the shuffling and distribution of individuals to memplexes is carried out separately for each sub-population (lines 7 and 8, where  $Mem^D$  denotes the  $m^D$  memplexes obtained from the dominance sub-population and  $Mem^I$  refers to the  $m^I$  memplexes associated to the indicator sub-population).

The parallel searches implemented through the processing of memplexes are conducted in the next step (lines 9-12, 13-16). New candidate solutions are generated according to the learning techniques and operators described in Section 4.2. The selection of the best individuals for the learning procedures depends on the multiobjective strategy used to manage each memplex. Memplexes from the dominance-based sub-population apply dominance-based mechanisms for this purpose. On the other hand, memplexes from the indicator-based sub-population employ indicator-based fitness values. The identification of the best global individuals is carried out over the sub-population where each memplex belongs (dominance-based sub-population for dominance-based memplexes, indicator-based sub-population for indicator-based memplexes).

**Algorithm 2** Adaptive MO-SFLA**Input parameters:**  $popSize, m, n, n_l, maxEval, ac$  (adaptive adjustment control).**Output:**  $PF$ .

---

```

1:  $P \leftarrow$  Set Initial Individuals ( $popSize$ ),  $PF \leftarrow 0$ 
2:  $P^D, P^I \leftarrow$  Split Population ( $P, popSize/2$ )
3:  $m^D, m^I \leftarrow m/2$ 
4: while ! stop criterion is reached ( $maxEval$ ) do
5:   Pareto Rank and Crowding Assignment ( $P^D, m^D * n$ )
6:   Indicator-based Fitness Assignment ( $P^I, m^I * n$ )
7:    $\{Mem_1^D \dots Mem_{m^D}^D\} \leftarrow$  Shuffling and Distribution ( $P^D, m^D, n$ )
8:    $\{Mem_1^I \dots Mem_{m^I}^I\} \leftarrow$  Shuffling and Distribution ( $P^I, m^I, n$ )
9:   for each memplex  $Mem_i^D \in \{Mem_1^D \dots Mem_{m^D}^D\}$  do
10:     Memplex Processing and Solution Generation ( $Mem_i^D, n_l$ )
11:      $PerfD_i \leftarrow$  Memplex Performance Assessment ( $Mem_i^D$ )
12:   end for
13:   for each memplex  $Mem_i^I \in \{Mem_1^I \dots Mem_{m^I}^I\}$  do
14:     Memplex Processing and Solution Generation ( $Mem_i^I, n_l$ )
15:      $PerfI_i \leftarrow$  Memplex Performance Assessment ( $Mem_i^I$ )
16:   end for
17:   if Current Generation %  $ac == 0$  then
18:      $PerfD_{mean}, PerfI_{mean} \leftarrow$  Mean Performance Computation ( $PerfD, PerfI, ac$ )
19:      $m^D, m^I \leftarrow$  Adaptive Memplex Assignment ( $PerfD_{mean}, PerfI_{mean}$ )
20:   end if
21:   Memplex Merging ( $P^D, \{Mem_1^D \dots Mem_{m^D}^D\}$ )
22:   Memplex Merging ( $P^I, \{Mem_1^I \dots Mem_{m^I}^I\}$ )
23:   Pareto Front Update ( $PF, P$ )
24: end while
25: return  $PF$ 

```

---

The key difference in this design lies in the fact that the loop for processing memplexes also involves the calculation of performance indicators  $PerfD_i / PerfI_i$ , in order to measure the success of the searches undertaken within each memplex  $Mem_i^D / Mem_i^I$ . Whenever a new candidate solution  $P'_{new}$  improves the one currently under processing  $Mem_{ij}^D / Mem_{ij}^I$ , the normalized distance between  $P'_{new}$  and  $Mem_{ij}^D / Mem_{ij}^I$  is calculated for each objective. In this way, it is possible to quantify how much the new solution has improved the previous one, using this information as a measurement of the success of the searches managed by each multiobjective strategy.

The adaptive adjustment of memplexes (lines 17-20) is carried out by taking into account the performance feedback gathered in the previous step. The number of memplexes  $m^D, m^I$  assigned to each multiobjective strategy is updated by using the following expressions:

$$m^{D'} = Round \left( \frac{\sum_{i=1}^{m^D} PerfD_i}{\sum_{i=1}^{m^D} PerfD_i + \sum_{i=1}^{m^I} PerfI_i} * m \right), \quad (9)$$



$$m^I = \text{Round} \left( \frac{\sum_{i=1}^{m^I} \text{Perf}I_i}{\sum_{i=1}^{m^D} \text{Perf}D_i + \sum_{i=1}^{m^I} \text{Perf}I_i} * m \right). \quad (10)$$

The idea behind Equations 9 and 10 is to assign a higher number of memplexes to the strategy that reported better performance when generating new candidate solutions. The memplexes moved from one strategy to another are those that obtained the lowest  $\text{Perf}I_i$  scores, since lower scores imply a higher degree of stagnation that can be tackled by using the alternative multiobjective strategy. Nevertheless, the proposed implementation guarantees that at least one memplex will be managed by the least successful multiobjective strategy, in order to give it additional opportunities in later stages of the optimization process. The adaptive step is carried out in accordance with a control parameter  $ac$ , which defines the number of generations in which performance feedback is registered prior to the adaptive re-distribution of memplexes. The introduction of this control parameter allows the algorithm to conduct the adaptive adjustment with certain statistical reliability, using as a reference the average performance obtained by each strategy throughout  $ac$  generations.

Finally, the memplexes assigned to each strategy are merged into the corresponding sub-populations (lines 21 and 22), in order to repeat the evolutionary process following the guidance of the adaptive procedure. A graphical representation of the evolution of data structures in the adaptive version of MO-SFLA can be found in Figure 1. It is worth remarking that the adaptive MO-SFLA works with two separate sub-population structures during the entire execution. The expressions 'partitioning' and 'shuffling and distribution' are used for the memplexes definition to highlight that these memplexes belong to a particular sub-population and therefore to a particular multiobjective strategy. The adaptive approach is in charge of updating dynamically the assignment of memplexes per sub-population according to the success of each multiobjective approach during the optimization process.

## 5. Experimental Evaluation

This section undertakes the comparative evaluation of the dominance-based, indicator-based, and adaptive versions of MO-SFLA. The experimental methodology, performance metrics, and statistical tests used in this study are first described. Afterwards, the quality of the outcomes reported by the

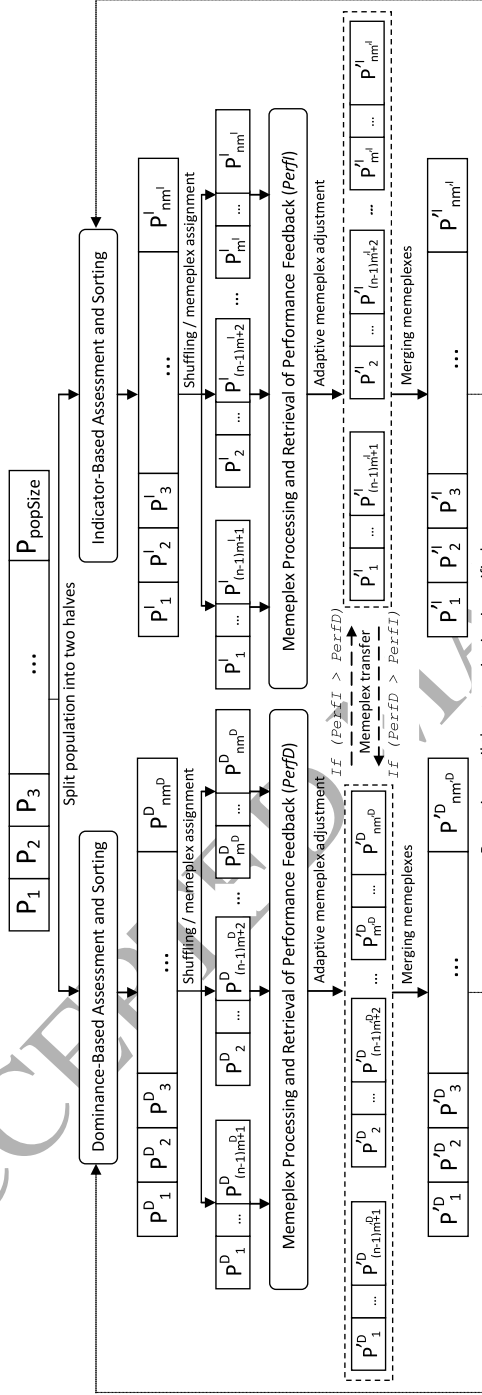


Figure 1: Evolution of population structures in the adaptive version of MO-SFLA. Starting from a balanced distribution of  $popSize/2$  solutions to each multiobjective strategy ('D' for dominance, 'I' for indicator), the corresponding fitness assessment and sorting mechanisms are applied over the assigned individuals. Memplexes are then defined through the shuffling technique ( $m^D$  memplexes at the dominance side and  $m^I$  at the indicator side, each memplex containing  $n$  individuals) and processed to generate new candidate solutions. During this step, performance metrics  $PerfD$  and  $PerfI$  are computed to find out the success of the searches performed within each memplex. This feedback is employed to adaptively adjust the assignment of individuals to each multiobjective strategy, transferring stagnated memplexes from the least successful strategy to the most successful one. Afterwards, the memplexes managed by each strategy ( $m^D$  and  $m^I$ ) are merged separately, repeating the previous steps until the stop condition is satisfied

dominance and indicator-based alternatives will be examined, putting emphasis on the potential divergence observed in the attained Pareto fronts. Then, the evaluation of results from the adaptive approach will be performed, making additional comparisons of solution quality with state-of-the-art methods for phylogenetic reconstruction.

The hardware platform used in the experimentation is supported on AMD Opteron 6174 twelve-core processors running at 2.2 GHz, with 12MB of L3 cache and 32GB of DDR3 RAM. This setup includes Ubuntu 14.04 LTS as operating system and the compiler GCC 5.2.1 was used to compile the software tested in this research (with the -O3 optimization flag enabled). The proposed implementations of MO-SFLA are based on OpenMP+MPI parallel schemes described in [35]. The experimentation has involved five problem instances containing real-world protein data with divergent characteristics attending to the number of sequences and sequence length:

1. M67x11333 [17], 67 sequences of bacterial ancestry euBac proteins (11,333 characters per sequence);
2. M88x3329 [27], 88 sequences of MCM7 and RPB1/RPB2 proteins from Thermophilic fungi (3,329 characters per sequence);
3. M187x814 [20], 187 sequences of ABC-B transporters from Mycorrhiza fungi (814 characters per sequence);
4. M260x1781 [41], 260 sequences of proto-oncogene proteins from *Beta vulgaris* (1,781 characters per sequence);
5. M355x1263 [12], 355 sequences of DHA2, ARN, and GEX proteins from hemiascomycete yeasts (1,263 characters per sequence).

These datasets exhibit different features attending to the two main dimensions of the problem -sequences and lengths-, thus representing evaluation scenarios with decision and objective spaces of different complexity. Consequently, they provide a suitable range of problem sizes for evaluation purposes. As the likelihood objective function specifically requires the use of phylogenetic evolutionary models to compute substitution probabilities along phylogeny branches, each dataset was analyzed with ProtTest [9] to determine the most fitting models. According to the output of this tool, likelihood scores for M67x11333, M88x3329, M187x814, and M355x1263 have been computed by using the LG+ $\Gamma$  model of protein evolution, while JTT+ $\Gamma$  was employed for M260x1781. Details on these models can be found in [1].

### 5.1. Performance Assessment Methodology

This experimental study is oriented towards examining the performance of the three reported designs of MO-SFLA. With this purpose in mind, the Pareto sets generated by each design alternative will be evaluated by using two multiobjective quality metrics. The first one is the set coverage  $SC$  (also known as ‘coverage relation’), which allows pairwise comparisons of the outcomes from the approaches under evaluation. Given two Pareto approximation sets  $X$  and  $Y$ , the  $SC$  metric calculates the fraction of solutions from  $Y$  that are covered (weakly-dominated) by the solutions from  $X$ :

$$SC(X, Y) = \frac{|\{y \in Y, \exists x \in X : x \succeq y\}|}{|Y|}. \quad (11)$$

Under this metric, priority must be given to the alternative that maximizes the achieved  $SC$  (that is, the one that covers a higher fraction of solutions). The second metric herein adopted is the widely-used hypervolume indicator  $I_H$ . Hypervolume is used to measure the  $\eta$ -dimensional volume (area in the case of bi-objective optimization problems) of the objective space  $Z = \mathbb{R}^\eta$  that is covered by at least one solution  $x$  from the Pareto approximation set  $X$  under evaluation. This value corresponds to the volume of the orthogonal polytope  $\Pi^\eta$ :

$$\Pi^\eta = \{p \in \mathbb{R}^\eta : p \preceq x \text{ for some } x \in X\}. \quad (12)$$

Higher hypervolume scores imply better multiobjective quality, thus being preferred those alternatives that maximize this metric. In order to avoid the influence of divergent scales in the objective functions, the ideal and nadir points from Table 1 are used to normalize objective scores in the scale  $[0,1]$  prior to the calculation of hypervolume. Using these normalized scores, hypervolume is calculated with regard to the reference point  $Z_{ref} = (1, 1)$ . The ideal and nadir points herein used were updated with regard to the previous work in [35] to improve the accuracy of the metric. Particularly, they were set by adding / subtracting a 0.05% to the best and worst objective scores observed in the experiments, instead of the former 0.25% used previously.

The configuration of input parameters of MO-SFLA was carried out via parametric studies, which were conducted taking into account the relationships between parameters analyzed in [35]. Particularly, sets of uniformly-distributed candidate values were checked for each input parameter, examining the multiobjective quality obtained by each configuration through the

Table 1: Hypervolume indicator: ideal and nadir points. These points (expressed in terms of  $P(T)$  and  $L(T)$  values) allow the normalization of objective scores to avoid bias in hypervolume due to different objective scales

Dataset	Ideal Point		Nadir Point	
	$P(T)$	$L(T)$	$P(T)$	$L(T)$
M67x11333	171,540	-473,023.58	196,869	-491,215.09
M88x3329	33,456	-149,020.30	33,668	-149,450.07
M187x814	29,832	-133,804.97	30,213	-134,944.68
M260x1781	43,507	-163,813.35	44,519	-165,660.07
M355x1263	54,795	-231,064.52	55,294	-233,866.76

Table 2: MO-SFLA input parameter values. These values represent the most satisfying configuration of the metaheuristic found in the conducted parametric studies

Parameter	Value
Population size ( $popSize$ )	128
Number of memplexes ( $m$ )	32
Individuals per memplex ( $n$ )	4
Number of learning steps ( $n_l$ )	4
Adaptive adjustment ( $ac$ , adaptive version only)	5

previously described multiobjective metrics. Table 2 shows input parameter values for the configuration that led to the best overall multiobjective behaviour. The stop criterion was established to 12,000 evaluations. For testing purposes, the dataset M67x11333 was employed as the reference instance in the parameterization.

The comparative study of the proposed design alternatives of MO-SFLA involved, for each experiment, 31 independent runs per dataset. The generated result samples were analyzed with the following statistical methodology [37] (considering a confidence level of 95%). Kolmogorov-Smirnov normality tests were first used to check for Gaussian distributions in the samples under comparison. If so, Levene tests were performed to analyze potential homoscedasticity. In case of detecting homogeneity in variances, the verification of statistically significant differences among samples was conducted via ANOVA. On the other hand, in case of not detecting Gaussian distributions or homogeneity in variances, such verifications were carried out through Wilcoxon-Mann-Whitney tests. This hybrid methodology was employed instead of a pure non-parametric one to avoid losing information about the characteristics of the examined result samples [26]. In fact, this kind of

methodology represents one of the most widely-used approaches to ensure accurate statistical reliability in evolutionary computation studies [23].

### 5.2. Evaluation of Dominance and Indicator-based Designs

The first step in this comparative analysis consists in verifying the potential divergence in multiobjective performance when different multiobjective strategies (dominance or indicator, separately) are adopted in the algorithmic design of MO-SFLA. To this end, differences in the sets of solutions generated by the dominance-based and indicator-based designs from Section 4.2 are examined. Figure 2 presents the Pareto fronts reported by the two alternatives under evaluation in their median-hypervolume executions.

The analysis of these Pareto fronts shows that the design alternatives under comparison exhibit differences from a multiobjective perspective. In overall terms, it can be observed that the indicator-based design of MO-SFLA shows a stronger focus on attaining high-quality results in the left side of the front, that is, in the region containing the best solutions attending to the parsimony objective. On the other hand, the dominance-based approach succeeds in achieving more effectiveness in the right side of the front, which comprises the best solutions attending to the likelihood objective.

In order to provide quantitative measurements of such multiobjective divergence, the set coverage scores obtained by the considered MO-SFLA alternatives have been calculated for each region of the Pareto front. The parsimony region is defined as the one containing those solutions included in the first half of the front, while the likelihood region comprises the points belonging to the second half of the front. The criterion used to determine the boundary (or ‘middle point’) was the middle parsimony score, as this function depends exclusively on the topology thus allowing the distinction between parsimony-oriented topologies (parsimony region) and likelihood-oriented ones (likelihood region). Table 3 shows the results reported by the set coverage metric when applied over each one of these regions separately, where  $SC(Domin, Indic)$  refers to the fraction of solutions from the indicator-based design that are covered by the ones from the dominance-based one and  $SC(Indic, Domin)$  represents the solutions from the dominance-based version that are improved by the indicator-based counterpart.

Focusing on the parsimony region of the front, it can be observed that the indicator-based design manages to obtain set coverage scores of 59.3% (M187x814) – 95.0% (M355x1263), thus covering average percentages over 77% of the solutions reported by the dominance-based strategy. In this

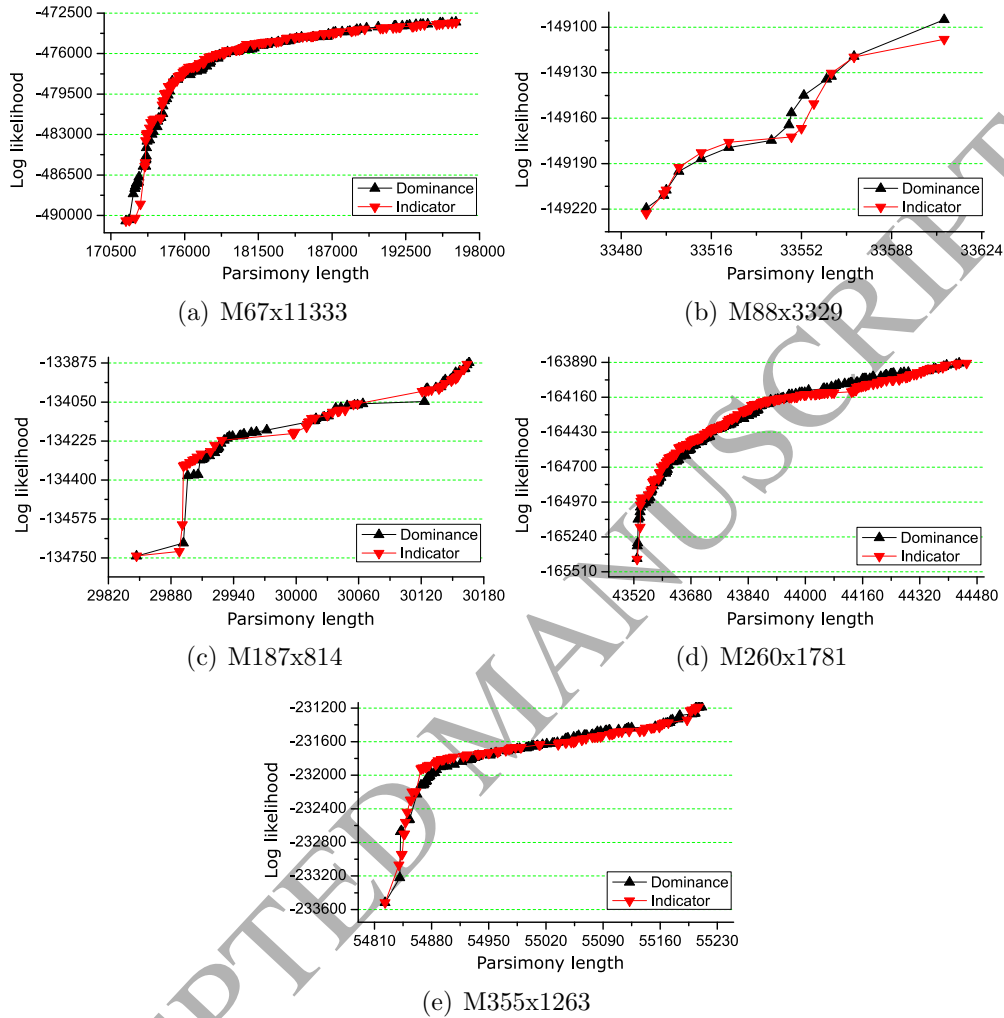


Figure 2: Comparison of Pareto fronts from the dominance-based and indicator-based designs of MO-SFLA. It can be observed the divergent performance of the considered strategies in the two halves of the fronts, being the parsimony region better exploited by the indicator-based design (mean set coverage of 77.4%) and the likelihood region by the dominance-based one (set coverage of 66.8%)

region, the dominance-based design underperforms in comparison to the indicator-based technique, obtaining in the best-case scenario a set coverage value of 17.7% (M187x814). The dominance-based approach plays a more significant role in the case of the likelihood region, according to the set

Table 3: Dominance vs. Indicator: set coverage for each front region.  $SC(X, Y)$  denotes the percentage of solutions from  $Y$  that are covered (weakly-dominated) by the solutions generated by  $X$ . In order to better understand the behaviour of each approach, two regions of the front (parsimony and likelihood regions) are distinguished, reporting set coverages for each of them. The best values in the comparison are highlighted in bold

Dataset	Parsimony region		Likelihood region	
	$SC(Domin, Indic)$	$SC(Indic, Domin)$	$SC(Domin, Indic)$	$SC(Indic, Domin)$
M67x11333	11.25%	<b>77.46%</b>	<b>51.22%</b>	37.50%
M88x3329	16.67%	<b>71.43%</b>	<b>83.33%</b>	14.29%
M187x814	17.65%	<b>59.26%</b>	<b>51.61%</b>	31.03%
M260x1781	10.26%	<b>83.65%</b>	<b>79.66%</b>	14.63%
M355x1263	11.11%	<b>95.00%</b>	<b>68.18%</b>	12.96%

coverage values of 51.2% – 83.3% attained over the indicator-based alternative. An average coverage fraction of 66.8% is reported in the likelihood side in comparison to the 22.1% achieved with the indicator strategy. Since each design alternative reports good performance in different regions of the front, these results support the idea of considering an adaptive dynamic approach to promote the attainment of improved solution quality in MO-SFLA.

### 5.3. Evaluation of the Adaptive Design

Next, insight is provided into the performance of the adaptive MO-SFLA in comparison with the dominance-based and indicator-based alternatives. Taking into account the previous observations, multiobjective behaviour will first be examined through the generated Pareto fronts. Figure 3 compares the Pareto fronts obtained by each alternative in their median-hypervolume executions. For the case of the dataset M67x11333, a more detailed representation of the distribution of solutions in the parsimony and likelihood regions is provided in Figure 4. In addition, Table 4 reports the results of evaluating the considered versions of MO-SFLA under the set coverage metric.

The upper side of Table 4 shows set coverage scores for each region of the front, where  $SC(Adapt, Y)$  refers to the percentage of solutions from the approach  $Y$  that are covered by the ones from the adaptive proposal. According to the obtained results, the adaptive version of MO-SFLA achieves significant solutions in both parsimony and likelihood regions of the Pareto front. In fact, the comparison suggests that the consideration of multiobjective strategies under adaptive techniques represents a promising idea from



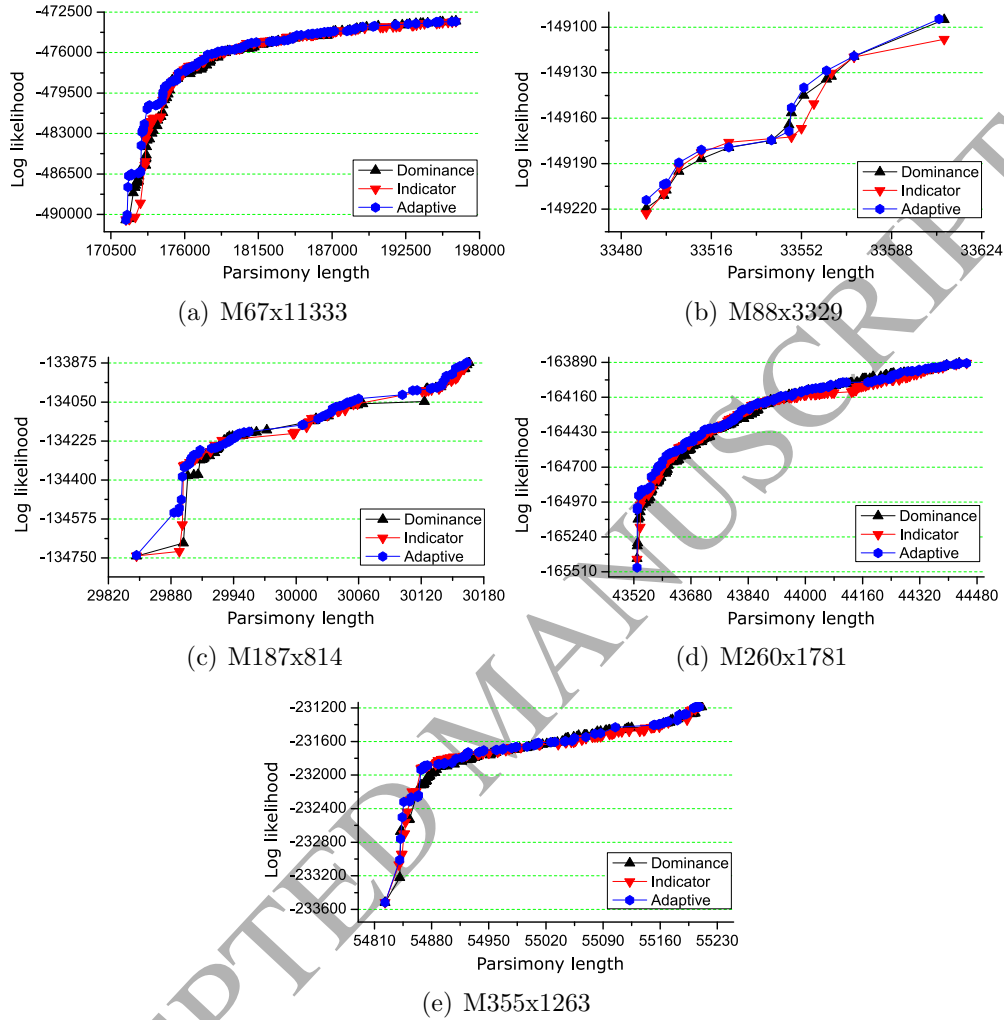
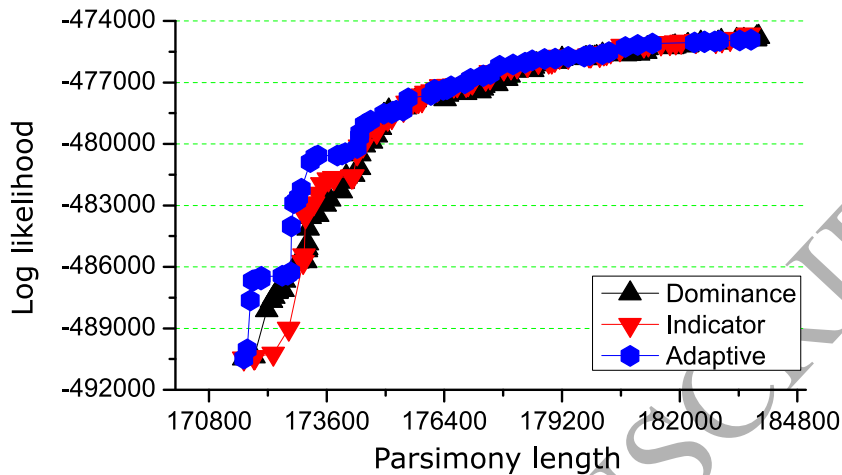
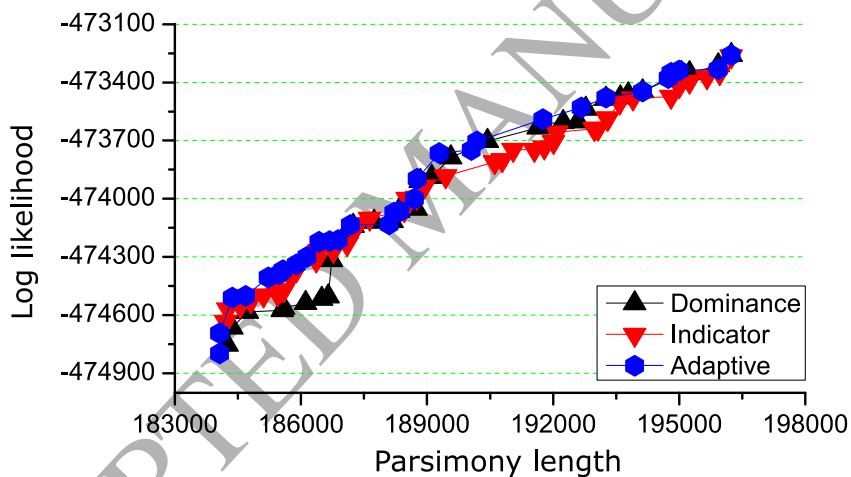


Figure 3: Representation of Pareto fronts from the adaptive MO-SFLA and comparisons with the dominance-based and indicator-based counterparts. The solutions generated by the adaptive MO-SFLA show significant quality throughout the whole front, going a step further with regard the dominance-based and indicator-based versions in both regions of the front. More specifically, set coverages up to 100% (over dominance) and 83.3% (over indicator) are achieved in the parsimony region of the front, while coverages up to 85.7% (over dominance) and 100% (over indicator) are observed in the likelihood one

this perspective, since the adaptive version is not limited to just matching the convergence reported by the best isolated approach at each region. More



(a) Parsimony region



(b) Likelihood region

Figure 4: Visualization of Pareto front regions for M67x11333. It can be distinguished the distribution of solutions from the different versions of MO-SFLA and how the adaptive approach accomplishes relevant results in both parsimony and likelihood regions

specifically, the adaptive design in the parsimony region achieves set coverage scores in the intervals 74.1% – 100.0% (over dominance) and 53.3% – 83.3% (over indicator), while also attaining significant coverage values in the likelihood region: 51.9% – 85.7% (over dominance) and 75.6% – 100.0% (over indicator). The bottom side of Table 4 shows the results of applying the cov-

Table 4: Adaptive version: set coverage evaluation. The upper side of the table reports the coverage values achieved, for each front region, by the adaptive version over the dominance-based and indicator-based counterparts. Additionally, the overall set coverage results observed in the whole front are provided in the bottom side. The best values in the comparison for the whole front are highlighted in bold

Dataset	Parsimony region		Likelihood region	
	$SC(Adapt, Domin)$	$SC(Adapt, Indic)$	$SC(Adapt, Domin)$	$SC(Adapt, Indic)$
M67x11333	87.32%	57.50%	59.38%	75.61%
M88x3329	100.00%	83.33%	85.71%	100.00%
M187x814	74.07%	64.71%	75.86%	80.65%
M260x1781	95.19%	65.81%	56.10%	93.22%
M355x1263	92.50%	53.33%	51.85%	86.36%
Whole-front coverage				
Dataset	$SC(Adapt, Domin)$	$SC(Domin, Adapt)$	$SC(Adapt, Indic)$	$SC(Indic, Adapt)$
M67x11333	<b>78.64%</b>	9.47%	<b>63.64%</b>	21.05%
M88x3329	<b>92.86%</b>	7.69%	<b>91.67%</b>	7.69%
M187x814	<b>75.00%</b>	11.48%	<b>75.00%</b>	14.75%
M260x1781	<b>84.14%</b>	10.64%	<b>76.14%</b>	17.73%
M355x1263	<b>69.15%</b>	15.28%	<b>69.66%</b>	20.83%

erage relation when both regions of the front are considered. In this scenario, coverage percentages of 69.2% – 92.9% and 63.6% – 91.7% are observed with regard to the dominance-based and indicator-based versions, respectively.

The assessment of the three proposed designs of MO-SFLA under hypervolume is given by Table 5. This table reports the median hypervolume scores and quartile deviations observed for each problem instance. In all these evaluation scenarios, the adaptive proposal shows more satisfactory behaviour than the dominance-based and indicator-based alternatives, achieving improved hypervolume scores in the range 60.4% (M88x3329) – 84.4% (M67x1133). Furthermore, the hypervolume samples from the adaptive version verify less variability in overall terms, in accordance with the observed quartile deviations. In order to find out if the improvement obtained by the adaptive version of MO-SFLA is statistically significant, statistical tests were conducted over the resulting hypervolume samples. The output of the statistical analysis is shown in Table 6, which details the P-values reported when comparing the adaptive version with the dominance-based and indicator-

Table 5: Comparison of hypervolume performance between the adaptive and the dominance-based / indicator-based versions of MO-SFLA. Hypervolumes are reported in the format  $I_H \pm qd$ , where  $I_H$  is the observed median hypervolume score and  $qd$  the quartile deviation. The best values in the comparison are highlighted in bold

Dataset	Adaptive	Dominance	Indicator
	$I_H$ Score	$I_H$ Score	$I_H$ Score
M67x11333	<b>84.38%±0.15</b>	82.42%±0.27	82.65%±0.22
M88x3329	<b>60.40%±0.14</b>	59.97%±0.17	59.43%±0.20
M187x814	<b>66.89%±1.09</b>	65.09%±1.34	65.18%±1.17
M260x1781	<b>77.93%±0.37</b>	76.68%±0.31	76.64%±0.32
M355x1263	<b>75.30%±0.55</b>	74.41%±0.68	74.78%±0.77

Table 6: Statistical assessment of hypervolume samples from the adaptive version of MO-SFLA. Considering a confidence level of 95%, P-values and the corresponding test output ( $\times$  = non-significant differences,  $\checkmark$  = significant differences) are reported

Dataset	Vs. Dominance		Vs. Indicator	
	P-value	Stat. Sign.?	P-value	Stat. Sign.?
M67x11333	1.33E-11	$\checkmark$ (<0.05)	1.33E-11	$\checkmark$ (<0.05)
M88x3329	3.15E-06	$\checkmark$ (<0.05)	7.29E-10	$\checkmark$ (<0.05)
M187x814	4.33E-08	$\checkmark$ (<0.05)	5.38E-07	$\checkmark$ (<0.05)
M260x1781	5.10E-10	$\checkmark$ (<0.05)	3.25E-10	$\checkmark$ (<0.05)
M355x1263	1.18E-04	$\checkmark$ (<0.05)	0.001	$\checkmark$ (<0.05)

based counterparts. The obtained P-values (below 0.05) show the relevance of the adaptive design, as this method leads to a statistically significant improvement in multiobjective quality over the remaining approaches in all the protein datasets under study.

Therefore, both the set coverage and hypervolume confirm the significant results attained by the adaptive version of MO-SFLA in terms of multiobjective quality. In this sense, it is important to emphasize that the inclusion of adaptive techniques in MO-SFLA does not have a noticeable impact in the execution time of the metaheuristic. To illustrate this point, Table 7 provides the median execution times reported by each version of MO-SFLA. In average terms, the adaptive method only implies additional time percentages of 1.96% and 1.51% with regard to the other design alternatives considered.

Figure 5 graphically shows how the adaptive adjustment of memeplexes evolves throughout the execution of MO-SFLA in each dataset. It can be observed that the proposal is able to adapt itself to different search scenarios

Table 7: Comparison of execution times (in seconds, using 32 cores).  $T_{exec}$  refers to the median execution time reported by each MO-SFLA design, while  $\Delta(Domin)$ ,  $\Delta(Indic)$  represent the difference between the adaptive version and the dominance-based and indicator-based counterparts, respectively

Dataset	$T_{exec}(Domin)$	$T_{exec}(Indic)$	$T_{exec}(Adapt)$	$\Delta(Domin)$	$\Delta(Indic)$
M67x11333	9135.04	9187.75	9239.93	1.15%	0.57%
M88x3329	1549.27	1562.86	1595.82	3.00%	2.11%
M187x814	1201.49	1202.12	1234.52	2.75%	2.70%
M260x1781	1960.13	1965.07	1990.57	1.55%	1.30%
M355x1263	3572.26	3589.38	3620.24	1.34%	0.86%

represented by the considered problem instances. The assignment of memplexes also suggests the relevance of the two adopted multiobjective strategies (dominance and indicator), as they are able to manage significant numbers of memplexes (at least 10 in the worst case scenario) without none of them being totally neglected throughout the entire optimization process.

#### 5.4. Comparisons with Other Algorithms and Tools

Once examined the adaptive MO-SFLA over other design alternatives, this section undertakes next comparisons of solution quality with other multiobjective methods and state-of-the-art tools for the tackled problem.

Focusing on multiobjective performance, two standard multiobjective evolutionary algorithms have been considered in the comparisons: NSGA-II [11] and IBEA [50]. These algorithms represent two of the most commonly employed reference metaheuristics for dominance-based and indicator-based multiobjective optimization, respectively. NSGA-II and IBEA have been adapted to phylogenetic reconstruction by using a matrix-shaped individual representation. As evolutionary operators, they include binary tournament selection (according to rank and crowding values in NSGA-II and  $I_{HD}$ -based fitness scores in IBEA), uniform crossover, and gamma-distributed mutation of evolutionary distances.

Comparisons of set coverage and hypervolume scores from the adaptive MO-SFLA, NSGA-II, and IBEA are introduced in Tables 8 and 9 (median results from 31 independent runs). In addition, Table 10 includes the statistical evaluation of hypervolume samples from the adaptive MO-SFLA, reporting if statistically significant differences were observed over NSGA-II and IBEA. Attending to the set coverage metric, the solutions reported by the proposed adaptive approach are able to dominate significant percentages

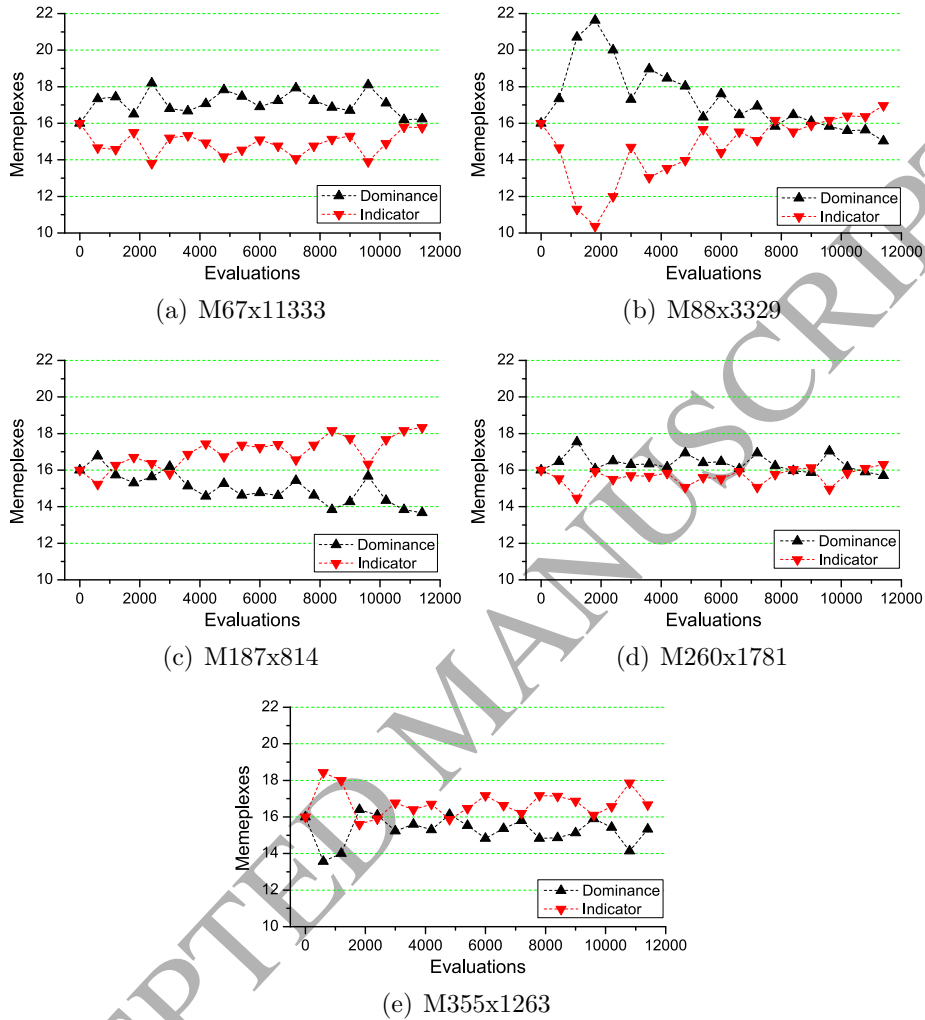


Figure 5: Evolution of the adjustment of memplexes in the adaptive version of MO-SFLA (mean values observed in the experimentation). It can be observed the roles assigned to each multiobjective strategy throughout the optimization process in the considered evaluation scenarios. All to all, both strategies are able to actively collaborate during the whole execution of the metaheuristic, without none of them being totally stagnated

of the fronts generated by the two other multiobjective algorithms under comparison. More specifically, set coverage values in the intervals 84.6% – 100% and 69.2% – 98.2% are respectively achieved over NSGA-II and IBEA. The attained hypervolume scores also highlight the remarkable multiobjec-

Table 8: Adaptive MO-SFLA: set coverage comparisons with NSGA-II and IBEA.  $SC(MO-SFLA, NSGA-II)$  and  $SC(MO-SFLA, IBEA)$  denote the percentage of solutions from NSGA-II and IBEA that are covered (weakly-dominated) by the solutions from MO-SFLA, while  $SC(NSGA-II, MO-SFLA)$  and  $SC(IBEA, MO-SFLA)$  refer to the solutions from MO-SFLA that are covered by NSGA-II and IBEA, respectively. The best values in the comparisons are highlighted in bold

Dataset	$SC(MO-SFLA, NSGA-II)$	$SC(NSGA-II, MO-SFLA)$	$SC(MO-SFLA, IBEA)$	$SC(IBEA, MO-SFLA)$
M67x11333	<b>97.17%</b>	2.11%	<b>98.21%</b>	1.05%
M88x3329	<b>100.00%</b>	0.00%	<b>85.71%</b>	15.39%
M187x814	<b>100.00%</b>	0.00%	<b>72.73%</b>	4.92%
M260x1781	<b>87.34%</b>	11.35%	<b>89.01%</b>	2.13%
M355x1263	<b>84.62%</b>	8.33%	<b>69.23%</b>	8.33%

Table 9: Adaptive MO-SFLA: hypervolume comparisons with NSGA-II and IBEA. Hypervolumes are reported in the format  $I_H \pm qd$ , where  $I_H$  is the median hypervolume score observed for each algorithm and  $qd$  the quartile deviation. The best values in the comparisons are highlighted in bold

Dataset	MO-SFLA	NSGA-II	IBEA
	$I_H$ Score	$I_H$ Score	$I_H$ Score
M67x11333	<b>84.38%±0.15</b>	78.79%±0.33	75.85%±0.39
M88x3329	<b>60.40%±0.14</b>	55.98%±0.54	57.42%±0.56
M187x814	<b>66.89%±1.09</b>	58.09%±1.82	59.38%±2.21
M260x1781	<b>77.93%±0.37</b>	65.15%±1.32	71.32%±0.78
M355x1263	<b>75.30%±0.55</b>	71.77%±0.76	72.45%±0.74

Table 10: Adaptive MO-SFLA: statistical comparison of hypervolume samples with NSGA-II and IBEA. Under a confidence level of 95%, P-values and the corresponding test output ( $\times$  = non-significant differences,  $\checkmark$  = significant differences) are reported

Dataset	Vs. NSGA-II		Vs. IBEA	
	P-value	Stat. Sign.?	P-value	Stat. Sign.?
M67x11333	1.29E-11	$\checkmark$ (<0.05)	1.33E-11	$\checkmark$ (<0.05)
M88x3329	1.33E-11	$\checkmark$ (<0.05)	1.33E-11	$\checkmark$ (<0.05)
M187x814	1.33E-11	$\checkmark$ (<0.05)	1.33E-11	$\checkmark$ (<0.05)
M260x1781	1.33E-11	$\checkmark$ (<0.05)	1.33E-11	$\checkmark$ (<0.05)
M355x1263	1.31E-11	$\checkmark$ (<0.05)	4.59E-06	$\checkmark$ (<0.05)

tive performance reported by the adaptive MO-SFLA (60.4% – 84.4%), going a step forward with regard to NSGA-II (55.9% – 78.8%) and IBEA (57.4%

– 75.9%) in all the problem instances under study.

In fact, the statistical testing of these results in Table 10 confirms that the adaptive MO-SFLA leads to a statistically significant improvement over these two reference methods for dominance-based and indicator-based multiobjective optimization. In this sense, the significance of the proposed adaptive MO-SFLA lies in the boosted search capabilities provided by the different components of the algorithm, namely: 1) parallel searches to process different directions of the search space simultaneously; 2) memplexes merging and shuffling techniques for information sharing; 3) definition of multiple search operators based on swarm techniques; and more importantly 4) adaptive mechanisms to exploit the most effective multiobjective strategy at each stage of the execution of the algorithm. The integration of all these components gives as a result improved multiobjective behaviour for this hard-to-tackle problem in comparison to the two standard multiobjective methods herein considered.

With regard to the evaluation of phylogenetic quality, the proposed adaptive approach is compared with up to six methods for protein-based phylogenetic reconstruction, each one with different characteristics:

1. Two methods for parsimony-based analysis: TNT [14] and ProtPars [13]. TNT is a reference tool that integrates multiple heuristics and operators, such as sectorial searches and tree fusing. On the other hand, ProtPars is the standard method provided in the PHYLIP phylogenetic package to carry out parsimony analysis from protein data.
2. Four methods for likelihood-based analysis: RAxML [40], IQ-TREE [29], GARLI [3], and MrBayes [33]. RAxML is a high-performance tool that combines different low-level optimizations and search techniques, such as rapid hill climbing and lazy subtree rearrangements. IQ-TREE provides a hill climbing-based stochastic approach to undertake efficient tree reconstructions. On the other hand, GARLI is a hybrid evolutionary algorithm that integrates local searches into a genetic algorithm scheme. Finally, MrBayes uses Markov Chain Monte Carlo techniques and stepping-stone algorithms to define a Bayesian framework for inferring phylogenies.

The quality of the extreme points in the median-hypervolume execution of MO-SFLA has been compared with the median solutions generated by each method (from 31 independent runs, using parametric configurations that matched the execution time of MO-SFLA). In order to provide statistical



robustness to these comparisons, two bio-statistical testing procedures have been applied [47]. For the parsimony comparisons, the Kishino-Hasegawa (KH) test is employed. This test examines divergence in the substitutions or mutations detected in the compared phylogenetic topologies to report a T-value that denotes if statistically significant differences were found among them. For the likelihood case, the CONSEL approximately unbiased test (AU) is applied to classify solutions attending to their statistical chance of representing the most fitting likelihood hypothesis for the input data (where a higher AU-value denotes higher statistical chance).

Tables 11 and 12 report parsimony / likelihood scores for each method, along with the output of the corresponding bio-statistical tests. According to these results, accurate parsimony solutions are generated by the adaptive proposal as it is able to match the scores found by the state-of-the-art tool TNT. In addition, both MO-SFLA and TNT succeeded in achieving a statistically significant improvement over ProtPars in all the datasets, as pointed out by the T-values reported by the KH test. As for likelihood, the adaptive proposal reaches the scores reported by RAxML and IQ-TREE in M67x11333 and M88x3329, while improving all the methods under comparison for the remaining datasets. In this sense, the CONSEL test sheds light on the significant quality of the likelihood results generated by the adaptive MO-SFLA, attaining the highest AU-values in all the problem instances under study. Furthermore, the differences observed in the likelihood scores are more noticeable in problem instances with increased complexity. In M260x1781 and M355x1263, the adaptive proposal is able to report significant likelihood values not only with regard to the reference methods but also when comparing with the previous dominance-based version of MO-SFLA [35] (which originally attained values of -163,895.81 and -231,186.40).

In order to provide further insight into the benefits of the adaptive MO-SFLA over the previous version, the biological quality of representative solutions from the median Pareto fronts for the two mentioned harder instances, M260x1781 and M355x1263, has been compared. The L2 metric was employed in this context to select the solutions that minimize the Euclidean distance to the ideal point in Table 1. The phylogenies from the adaptive MO-SFLA showed scores of  $P(T)=43,717$ ,  $L(T)=-164,414.86$  in M260x1781 and  $P(T)=54,870$ ,  $L(T)=-231,924.79$  in M355x1263. On the other hand, the original MO-SFLA achieved solutions with scores of  $P(T)=43,739$ ,  $L(T)=-164,425.51$  in M260x1781 and  $P(T)=54,886$ ,  $L(T)=-231,932.95$  in M355x1263. Hence, the solutions reported by the adaptive proposal dominate the ones

Table 11: Adaptive MO-SFLA: comparisons of solution quality (parsimony) with reference biological tools (TNT and ProtPars). Parsimony scores ( $P(T)$ ) and the results of the Kishino-Hasegawa test (T-value) are provided for each dataset. The best values in the comparison are highlighted in bold

		M67x11333		M88x3329	
Method	$P(T)$	T-value	$P(T)$	T-value	
MO-SFLA	<b>171,623</b>	<b>Best</b>	<b>33,490</b>	<b>Best</b>	
TNT	<b>171,623</b>	<b>0.00</b>	<b>33,490</b>	<b>0.00</b>	
ProtPars	173,768	15.07	33,944	9.28	
		M187x814		M260x1781	
Method	$P(T)$	T-value	$P(T)$	T-value	
MO-SFLA	<b>29,847</b>	<b>Best</b>	<b>43,529</b>	<b>Best</b>	
TNT	<b>29,847</b>	<b>0.00</b>	<b>43,529</b>	<b>0.00</b>	
ProtPars	29,955	3.01	44,479	14.81	
		M355x1263			
Method	$P(T)$	T-value			
MO-SFLA	<b>54,823</b>	<b>Best</b>			
TNT	<b>54,823</b>	<b>0.00</b>			
ProtPars	55,328	9.27			

from the original version, since an improvement is observed in both parsimony and likelihood objectives. This comparison suggests that the proposed adaptive approach is able to go a step further in the biological quality of the solution, boosting the accuracy of MO-SFLA when complex problem instances are considered.

As in the assessment of multiobjective performance, the improvement observed over the biological methods under comparison shows the relevant optimization capabilities of the adaptive MO-SFLA. This algorithmic design addresses the processing of highly complex search spaces in parallel through the memplex concept. These memplexes evolve independently for a certain number of learning steps, sharing afterwards the attained knowledge by means of the merging and shuffling techniques. Moreover, the accuracy of the method is improved by adaptively using different strategies to assess solutions and search operators, in such a way that the algorithm is able to put more effort on the processing of promising solutions without dismissing the exploration of alternative candidates.

In conclusion, the obtained results suggest the relevance of the proposed adaptive design of MO-SFLA. Considering different multiobjective strategies

Table 12: Adaptive MO-SFLA: comparisons of solution quality (likelihood) with reference biological tools (RAxML, IQ-TREE, GARLI, and MrBayes). Likelihood scores ( $L(T)$ ) and the results of the CONSEL approximately unbiased test (AU-value) are provided for each dataset. The best values in the comparison are highlighted in bold

		M67x11333		M88x3329	
Method	$L(T)$	AU-value	$L(T)$	AU-value	
MO-SFLA	<b>-473,260.21</b>	<b>0.68</b>	<b>-149,094.84</b>	<b>0.68</b>	
RAxML	<b>-473,260.21</b>	0.48	<b>-149,094.84</b>	0.64	
IQ-TREE	<b>-473,260.21</b>	0.63	<b>-149,094.84</b>	0.35	
GARLI	-473,264.64	0.26	-149,110.96	0.24	
MrBayes	-473,404.42	0.00	-149,809.72	0.00	
		M187x814		M260x1781	
Method	$L(T)$	AU-value	$L(T)$	AU-value	
MO-SFLA	<b>-133,871.90</b>	<b>0.61</b>	<b>-163,895.63</b>	<b>0.62</b>	
RAxML	-133,877.80	0.33	-163,911.41	0.28	
IQ-TREE	-133,871.96	0.60	-163,899.10	0.52	
GARLI	-133,875.58	0.44	-163,975.20	0.00	
MrBayes	-134,008.88	0.01	-165,929.91	0.00	
		M355x1263			
Method	$L(T)$	AU-value			
MO-SFLA	<b>-231,185.96</b>	<b>0.79</b>			
RAxML	-231,199.54	0.36			
IQ-TREE	-231,299.10	0.12			
GARLI	-231,660.77	0.00			
MrBayes	-233,060.07	0.00			

managed with adaptive procedures, this metaheuristic is able to tackle MOPs by dynamically adjusting memplex-based searches according to the most successful strategy in each stage of the optimization process. The adoption of such techniques gives as a result a robust metaheuristic engine, with improved search capabilities that lead to relevant results when reconstructing ancestral relationships in challenging protein scenarios.

## 6. Conclusions

This work focused on the study of different design alternatives for the multiobjective metaheuristic MO-SFLA, a novel optimization method that combines swarm-based search operators and parallel searches to address hard-to-tackle MOPs. Using as a case study the inference of evolutionary rela-

tionships from protein data, two alternative implementations for MO-SFLA based on different multiobjective mechanisms were examined – dominance-based and indicator-based variants. Moreover, an additional adaptive design that dynamically makes use of both multiobjective strategies was introduced to improve the optimization capabilities of the metaheuristic. The adaptive version of MO-SFLA is based on the idea of gathering feedback from the parallel searches managed by each multiobjective strategy, in order to adjust them in accordance with the technique that exhibited more satisfying performance at each stage of the optimization process.

The proposed design alternatives have been assessed with a thorough experimental study that involved the evaluation of multiobjective performance in five problem instances comprising real-world protein data. The comparison of the Pareto fronts generated by the dominance-based and indicator-based approaches has shown divergences in the performance of each design alternative at different regions of the front. On the other hand, the proposed adaptive approach has successfully allowed MO-SFLA to take advantage of the capabilities of both multiobjective strategies, achieving significant results throughout the entire Pareto front. The hypervolume metric has also confirmed the relevance of the adaptive MO-SFLA, which reports statistically significant improvements over the dominance-based and indicator-based counterparts along with less variability in the generated results samples. Such improvements are also observed over two representative algorithms in multiobjective optimization, NSGA-II and IBEA. Finally, the evaluation of solution quality under bio-statistical tests (comparing with six state-of-the-art biological tools) has shown that the adaptive MO-SFLA successfully addresses challenging optimization scenarios, leading to significant results in increasingly difficult problem instances.

Future research involves the exploration of other techniques to provide additional robustness to MO-SFLA. Within the adaptive framework, the integration of additional multiobjective mechanisms specifically aimed at improving certain multiobjective properties will be pursued, as well as other search operators to boost the generation of new solutions during the memplex processing. Multiobjective behaviour in complex optimization scenarios will be further analyzed, using problem instances comprising both protein and genome data. Finally, the implementation and analysis of heterogeneous parallel schemes to run the metaheuristic in CPU+co-processor hardware setups will be undertaken. Such parallel designs are of particular interest when tackling many-objective optimization problems, due to the additional

challenges that could imply the consideration of a high number of objectives, e.g. over the indicator-based side of the application. In this sense, it is necessary an in-depth evaluation of multiple factors (number of objectives, quality indicator complexity, problem-specific operations, etc.) to identify potential implications in execution time. On the basis of this evaluation, parallel implementations [22] of the indicator-based strategy can be included to boost the efficiency of the method in hard-to-tackle many-objective scenarios.

### Acknowledgments

This work was partially funded by the AEI (State Research Agency, Spain) and the ERDF (European Regional Development Fund, EU), under the contract TIN2016-76259-P (PROTEIN project), as well as Portuguese national funds through FCT (Fundação para a Ciência e a Tecnologia, Portugal) projects UID/CEC/50021/2019 and LISBOA-01-0145-FEDER-031901 (PT-DC/CCI-COM/31901/2017, HiPErBio). Sergio Santander-Jiménez is supported by the Post-Doctoral Fellowship from FCT under Grant SFRH/BPD/-119220/2016. Finally, the authors would like to thank Dr. José Jasnau Caeiro for his comprehensive comments on the contents of the paper.

**Declarations of interest:** none

### References

- [1] M. Arenas, Trends in substitution models of molecular evolution, *Frontiers in Genetics* 6:319 (2015) 1–9.
- [2] R. Azzouz, S. Bechikh, L. B. Said, A dynamic multi-objective evolutionary algorithm using a change severity-based adaptive population management strategy, *Soft Comput.* 21 (4) (2017) 885–906.
- [3] A. L. Bazinet, D. J. Zwickl, M. P. Cummings, A Gateway for Phylogenetic Analysis Powered by Grid Computing Featuring GARLI 2.0, *Systematic Biology* 63 (5) (2014) 812–818.
- [4] M. S. Bittermann, I. S. Sariyildiz, An Adaptive Multi-Objective Evolutionary Algorithm with Human-Like Reasoning for Enhanced Decision-Making in Building Design, in: *Proc. of the 2011 IEEE Symposium on Computational Intelligence in Multicriteria Decision-Making*, IEEE, 2011, pp. 1–8.

- [5] W. Cancino, A. C. B. Delbem, A Multi-Criterion Evolutionary Approach Applied to Phylogenetic Reconstruction, in: *New Achievements in Evol. Comp.*, InTech, 2010, pp. 135–156.
- [6] D. Charlet, F. Spies, C. Bloch, W. Abdou, Adaptive Multi-objective Genetic Algorithm using Multi-Pareto-Ranking, in: *Genetic and Evolutionary Computation Conference*, 2012, pp. 449–456.
- [7] G. P. Coelho, A. E. A. Silva, F. J. V. Zuben, An Immune-Inspired Multi-Objective Approach to the Reconstruction of Phylogenetic Trees, *Neural Comput. Appl.* 19 (8) (2010) 1103–1132.
- [8] C. Coello, Multi-objective Evolutionary Algorithms in Real-World Applications: Some Recent Results and Current Challenges, in: *Advances in Evolutionary and Deterministic Methods for Design, Optimization and Control in Engineering and Sciences*, Springer, 2015, pp. 3–18.
- [9] D. Darriba, G. L. Taboada, R. Doallo, D. Posada, ProtTest 3: fast selection of best-fit models of protein evolution, *Bioinformatics* 27 (8) (2011) 1164–1165.
- [10] K. Deb, Multi-Objective Evolutionary Algorithms, in: *Springer Handbook of Computational Intelligence*, Springer, 2015, pp. 995–1015.
- [11] K. Deb, A. Pratap, S. Agarwal, T. Meyarivan, A Fast and Elitist Multi-Objective Genetic Algorithm: NSGA-II, *IEEE Trans. Evol. Comput.* 6 (2) (2002) 182–197.
- [12] P. J. Dias, I. Sá-Correia, The drug:H<sup>+</sup> antiporters of family 2 (DHA2), siderophore transporters (ARN) and glutathione:h<sup>+</sup>antiporters (GEX) have a common evolutionary origin in hemiascomycete yeasts, *BMC Genomics* 14 (901) (2013) 1–22.
- [13] J. Felsenstein, PHYLIP (phylogeny inference package), <http://evolution.genetics.washington.edu/phylip.html> (2014).
- [14] P. A. Goloboff, S. A. Catalano, TNT version 1.5, including a full implementation of phylogenetic morphometrics, *Cladistics* 32 (3) (2016) 221–238.

- [15] J. Guo, L. Chen, L. Qin, C. Wang, A Self-adaptive Ant Colony Algorithm for Phylogenetic Tree Construction, in: Proc. of ICHIT 2006, IEEE, 2006, pp. 1–6.
- [16] D. Hadka, P. Reed, Borg: An Auto-Adaptive Many-Objective Evolutionary Computing Framework, *Evolutionary Computation* 21 (2) (2013) 231–259.
- [17] D. He, O. Fiz-Palacios, C. Fu, J. Fehling, C. Tsai, S. L. Baldauf, An Alternative Root for the Eukaryote Tree of Life, *Current Biology* 24 (4) (2014) 465–470.
- [18] T. Hill, A. Lundgren, R. Fredriksson, H. B. Schiöth, Genetic algorithm for large-scale maximum parsimony phylogenetic analysis of proteins, *Biochim. Biophys. Acta.* 1725 (1) (2005) 19–29.
- [19] S. Jiang, et al., Adaptive Indicator-based Evolutionary Algorithm for Multiobjective Optimization Problems, in: Proc. of the IEEE Congress on Evolutionary Computation, CEC 2016, IEEE, 2016, pp. 492–499.
- [20] A. Kovalchuk, A. Kohler, F. Martin, F. O. Asiegbu, Diversity and evolution of ABC proteins in mycorrhiza-forming fungi, *BMC Evolutionary Biology* 15 (249) (2015) 1–19.
- [21] Q. Lin, et al., An adaptive immune-inspired multi-objective algorithm with multiple differential evolution strategies, *Information Sciences* 430–431 (2018) 46–64.
- [22] E. M. López, L. M. Antonio, C. Coello, A GPU-Based Algorithm for a Faster Hypervolume Contribution Computation, in: EMO 2015: Evolutionary Multi-Criterion Optimization, Vol. 9019 of LNCS, Springer Verlag, 2015, pp. 80–94.
- [23] G. Luque, E. Alba, *Parallel Genetic Algorithms: Theory and Real World Applications*, Springer Verlag, Berlin / Heidelberg, 2011.
- [24] H. Matsuda, Protein phylogenetic inference using maximum likelihood with a genetic algorithm, in: Proc. of the Pacific Symposium on Bio-computing 96, World Scientific, 1996, pp. 512–523.

- [25] H. Matsuda, H. Yamashita, Y. Kaneda, Molecular Phylogenetic Analysis using both DNA and Amino Acid Sequence Data and Its Parallelization, *Genome Informatics* 5 (1994) 120–129.
- [26] Q. Minella, R. Ruiz, M. Ciavotta, A Review and Evaluation of Multi-objective Algorithms for the Flowshop Scheduling Problem, *INFORMS Journal on Computing* 20 (3) (2008) 451–471.
- [27] I. Morgenstern, et al., A molecular phylogeny of thermophilic fungi, *Fungal Biology* 116 (4) (2012) 489–502.
- [28] C. A. Nicolaou, N. Brown, Multi-objective optimization methods in drug design, *Drug Discovery Today: Technologies* 10 (3) (2013) e427–e435.
- [29] L. Nguyen, H. Schmidt, A. Haeseler, B. Minh, IQ-TREE: A fast and effective stochastic algorithm for estimating maximum likelihood phylogenies, *Mol. Biol. Evol.* 32 (1) (2015) 268–274.
- [30] L. Poladian, L. Jermiin, Multi-Objective Evolutionary Algorithms and Phylogenetic Inference with Multiple Data Sets, *Soft Computing* 10 (4) (2006) 359–368.
- [31] T. H. Reijmers, R. Wehrens, F. D. Daeyaert, P. J. Lewi, L. M. C. Buydens, Using genetic algorithms for the construction of phylogenetic trees: application to G-protein coupled receptor sequences, *Biosystems* 49 (1) (1999) 31–43.
- [32] A. Rokas, Phylogenetic Analysis of Protein Sequence Data Using the Randomized Axelerated Maximum Likelihood (RAxML) Program, *Current Protocols in Molecular Biology* 96 (2011) 19.11.1–19.11.14.
- [33] F. Ronquist, et al., MrBayes 3.2: Efficient Bayesian Phylogenetic Inference and Model Choice Across a Large Model Space, *Systematic Biology* 61 (3) (2012) 539–542.
- [34] S. Santander-Jiménez, M. A. Vega-Rodríguez, Applying a Multiobjective Metaheuristic Inspired by Honey Bees to Phylogenetic Inference, *BioSystems* 114 (1) (2013) 39–55.
- [35] S. Santander-Jiménez, M. A. Vega-Rodríguez, L. Sousa, Multiobjective Frog-Leaping Optimization for the Study of Ancestral Relationships in



- Protein Data, *IEEE Transactions on Evolutionary Computation* 22 (6) (2018) 879–893.
- [36] A. Sarkheyli, A. M. Zain, S. Sharif, The role of basic, modified and hybrid shuffled frog leaping algorithm on optimization problems: a review, *Soft Computing* 19 (7) (2015) 2011–2038.
- [37] D. J. Sheskin, *Handbook of Parametric and Nonparametric Statistical Procedures*. 5th edition, Chapman & Hall/CRC, NY, USA, 2011.
- [38] A. Skourikhine, Phylogenetic Tree Reconstruction Using Self-Adaptive Genetic Algorithm, in: *Proc. of the 1st IEEE International Symposium on Bioinformatics and Biomedical Engineering*, 2000, pp. 129–134.
- [39] A. Stamatakis, An Efficient Program for Phylogenetic Inference Using Simulated Annealing, in: *Proc. of the 19th IEEE International Parallel and Distributed Processing Symposium*, IEEE, 2005, pp. 1–8.
- [40] A. Stamatakis, RAXML Version 8: A Tool for Phylogenetic Analysis and Post-Analysis of Large Phylogenies, *Bioinformatics* 30 (9) (2014) 1312–1313.
- [41] R. Stracke, et al., Genome-wide identification and characterisation of R2R3-MYB genes in sugar beet (*Beta vulgaris*), *BMC Plant Biology* 14 (249) (2014) 1–17.
- [42] K. C. Tan, S. C. Chiam, A. A. Mamun, C. K. Goh, Balancing exploration and exploitation with adaptive variation for evolutionary multi-objective optimization, *European Journal of Operational Research* 197 (2) (2009) 701–713.
- [43] R. Tanabe, H. Ishibuchi, A. Oyama, Benchmarking Multi- and Many-Objective Evolutionary Algorithms Under Two Optimization Scenarios, *IEEE Access* 5 (2017) 19597–19619.
- [44] P. Tangpattanakul, N. Jozefowicz, P. Lopez, Multi-Objective Optimization for Selecting and Scheduling Observations by Agile Earth Observing Satellites, in: *Parallel Problem Solving From Nature XII*, Vol. 7492 of LNCS, Springer Verlag, 2012, pp. 112–121.

- [45] L. Thiele, Indicator-Based Selection, in: Springer Handbook of Computational Intelligence, Springer, 2015, pp. 983–994.
- [46] M. Wang, Y. Wang, X. Wang, A Space Division Multiobjective Evolutionary Algorithm Based on Adaptive Multiple Fitness Functions, *Int. J. Patt. Recogn. Artif. Intell.* 30 (3) (2016) 1–25.
- [47] T. Warnow, Computational Phylogenetics: An Introduction to Designing Methods for Phylogeny Estimation, Cambridge Univ. Press, Cambridge, 2017.
- [48] G. G. Yen, Z. He, Performance Metric Ensemble for Multiobjective Evolutionary Algorithms, *IEEE Trans. Evol. Comput.* 18 (1) (2014) 131–144.
- [49] A. C. Zavoianu, E. Lughofer, G. Bramerdorfer, W. Amrhein, E. P. Klement, DECMO2: a robust hybrid and adaptive multi-objective evolutionary algorithm, *Soft Computing* 19 (12) (2015) 3551–3569.
- [50] E. Zitzler, S. Künzli, Indicator-Based Selection in Multiobjective Search, in: Parallel Problem Solving From Nature VIII, Vol. 3242 of LNCS, Springer Verlag, 2004, pp. 832–842.

## Vitae



**Sergio Santander-Jiménez** received the Ph.D. degree in Computer Engineering from the University of Extremadura, Spain, in 2016. He is currently a Post-doctoral Fellow and a Senior Researcher at the R&D Instituto de Engenharia de Sistemas e Computadores (INESC-ID), Instituto Superior Técnico (IST), Universidade de Lisboa (UL), Portugal. He has co-organized several international workshops on high-performance computing, computational intelligence, computational biology and bioinformatics, reviewing articles on these topics for different international JCR-indexed journals. His main research interests include evolutionary and bioinspired computing, multi-objective optimization, parallel and distributed computing, and their applications to real-world biological problems.



**Miguel A. Vega-Rodríguez** received the Ph.D. degree in Computer Engineering from the University of Extremadura, Spain, in 2003. He is currently an Associate Professor (accredited as Full Professor) of computer architecture in the Department of Computer and Communications Technologies, University of Extremadura. He has authored or co-authored more than 630 publications including journal papers (more than 120 JCR-indexed journal papers), book chapters, and peer-reviewed conference proceedings, for which he got several awards - such as Best Paper Awards in ISDA'11, IBER-GRID'11, ICEC'09, and IEA-AIE'08. His main research interests include parallel and distributed computing, evolutionary computation, bioinformatics, and reconfigurable and embedded computing.



**Leonel Sousa** received the Ph.D. degree in Electrical and Computer Engineering from the Instituto Superior Técnico (IST), Universidade de Lisboa (UL), Lisbon, Portugal, in 1996. He is currently a Full Professor with UL and a Senior Researcher with the R&D Instituto de Engenharia de Sistemas e Computadores (INESC-ID). He has authored or coauthored more than 200 papers in journals and international conferences, and has edited four special issues of in-

ternational journals. Prof. Sousa is a Fellow of the IET and a Distinguished Scientist of the ACM. His research interests include VLSI architectures, computer architectures and arithmetic, parallel computing, and signal processing systems.

ACCEPTED MANUSCRIPT