# **High-dimensional Black-box Optimization** via Divide and Approximate Conquer

Peng Yang<sup>1</sup>, Ke Tang<sup>1\*</sup>, and Xin Yao<sup>2</sup>

<sup>1</sup>UBRI, School of Computer Science and Technology, University of Science and Technology of China, Hefei, China, 230027 <sup>2</sup>CERCIA, School of Computer Science,

University of Birmingham, Birmingham B15 2TT, U.K.

Emails: trevor@mail.ustc.edu.cn; ketang@ustc.edu.cn; x.yao@cs.bham.ac.uk

## **Abstract**

Divide and Conquer (DC) is conceptually well suited to high-dimensional optimization by decomposing a problem into multiple small-scale subproblems. However, appealing performance can be seldom observed when the sub-problems are interdependent. This paper suggests that the major difficulty of tackling interdependent sub-problems lies in the precise evaluation of a partial solution (to a sub-problem), which can be overwhelmingly costly and thus makes sub-problems nontrivial to conquer. Thus, we propose an approximation approach, named Divide and Approximate Conquer (DAC), which reduces the cost of partial solution evaluation from exponential time to polynomial time. Meanwhile, the convergence to the global optimum (of the original problem) is still guaranteed. The effectiveness of DAC is demonstrated empirically on two sets of non-separable high-dimensional problems.

#### 1 Introduction

Developing Artificial Intelligence (AI) applications often encounters a key task of solving challenging optimization problems. Formally, it can be stated as:

$$\mathbf{x}^* = \arg\max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}),$$

 $\mathbf{x}^* = \arg\max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}),$  where  $f: \mathcal{X} \to \mathbb{R}$  denotes a function on a bounded solution space  $\mathcal{X} \subseteq \mathbb{R}^D$  and  $\mathbf{x}^*$  denotes the global optimum in  $\mathcal{X}$ . We consider f as a black-box function that the problem information is completely unknown beforehand, where only the function value of  $\mathbf{x} \in \mathcal{X}$  can be directly provided if explicitly queried. Therefore, only derivative-free approaches can be brought to bear, which, however, often suffer from the curse of dimensionality if D is considerably large.

An intuitive idea to handle a high-dimensional optimization problem is to project its solution space onto lower dimensions, where traditional approaches perform well [Kabán et al., 2015]. However, it is nontrivial to identify an appropriate projection. Typical approaches in this category, e.g., Random Embedding techniques [Wang et al., 2013; Qian and Yu, 2016], consider the high-dimensional problem having low effective dimensionality, for which a random projection would suffice to find the global optimal solution of a high-dimensional problem in a low-dimensional space. Although these algorithms also showed appealing performance in case the low effective dimensionality assumption is mildly violated, their performance may not be satisfactory on the wider range of irreducible problems.

Divide-and-Conquer (DC) is another general idea for tackling large-scale problems. In case of high-dimensional blackbox optimization, DC can be implemented by dividing the original problem into multiple sub-problems (say of dimensionality  $d_i$ , where i = 1, ..., M, and  $d_i < D$ ) [Yang et al., 2008a]. For the i-th sub-problem, **x** is optimized along  $d_i$  dimensions, while its values on the other  $D-d_i$  dimensions are fixed. That is, each sub-problem concerns a low-dimensional subspace of the original solution space. Applying an existing search method to the *i*-th sub-problem leads to a  $d_i$ dimensional partial solution to the original high-dimensional problem. The solution to original problem can be obtained by combining the partial solutions achieved on all sub-problems.

Given an appropriate sub-problems optimizer, the abovedescribed DC strategy works well on the so-called separable problems, for which the global optimal optimum can be found by optimizing one dimension at a time regardless of the values taken on the other dimensions [Chen et al., 2010]. If this condition does not hold, the performance of DC heavily relies on the decomposition method [Omidvar et al., 2014], which aims to divide a black-box high-dimensional problem in such a way that the global optimum can still be obtained by solving the sub-problems in a fully independent manner. In the past few years, a large variety of decomposition methods have been proposed [Mahdavi et al., 2015]. Despite the performance enhancement brought by them, none of these methods are guaranteed to achieve the desired sub-problems. Meanwhile, a practical problem of interest may be fully nonseparable such that the ideal decomposition mentioned above does not even exist. Therefore, how to deal with (conquer) the interdependent sub-problems remains a challenge to DC in the context of high-dimensional optimization.

In this paper, we suggest that the major challenge of tackling interdependent sub-problems lies in the difficulty of evaluating the quality of a partial solution (to a sub-problem) during the search course. To be specific, as the quality of a partial solution (to the original problem) depends on the values taken on the dimensions involved in other sub-problems, precisely evaluating a partial solution requires overwhelming computation cost, which increases exponentially with the number of interacting variables in other sub-problems. We propose an approximation approach for partial solution evaluation, which yields a novel framework, named Divide and Approximate Conquer (DAC). With DAC, the computational cost increases polynomially with the number of solutions generated in each iteration while the convergence to global optimum (of the original problem) is still guaranteed.

The major difficulty of tackling interdependent subproblems is analyzed in Section 2. The proposed DAC is detailed in Section 3. Section 4 reports empirical studies of DAC on five synthetic high-dimensional optimization problems and the hyper-parameters fine-tuning task for multiclass SVM. Section 5 concludes this work.

## 2 Major challenge of dealing with interdependent sub-problems

The DC strategy consists of three steps:

- 1) (**Divide**) Decompose a problem into M  $d_i$ -dimensional sub-problems, where i = 1, ..., M and  $d_i < D$ ;
- 2) (**Conquer**) Search the best partial solution in each subproblem by applying an existing search approach;
- 3) Merge the best partial solutions obtained on all subproblems as the final output.

We restrict our discussions here to the Conquer phase. Usually, a derivative-free search process in each sub-problem is guided by the solutions with better function values, despite a few exceptions [Kirkpatrick et al., 1983; Tang et al., 2016]. Before a  $d_i$ -dimensional partial solution is evaluated, it must be complemented to D-dimensional by fixing the values for the variables in other sub-problems. Specially, we call a vector of such fixed  $D-d_i$  values as a complement to a partial solution. A partial solution will receive different function values with different complements. Fortunately, the indeterminate function values of partial solutions will not influence the search process unless the rank of partial solutions changes, on which the search direction is actually determined. The rank of partial solutions may change if sub-problems are interdependent, which is defined as follows:

**Definition 1.** (Interacting Variables) [Chen et al., 2010] Given a D-dimensional problem, its i-th and j-th variables are said to be interacting, if the rank of two partial solutions  $x_i$  and  $x'_i$  in the i-th dimension may change with different complements, e.g.,  $x_j$  and  $x'_j$ , in the j-th variable:

$$\begin{split} \exists \mathbf{x}, x_i', x_j' : \\ f(x_1, ..., x_i, ..., x_j, ..., x_D) < f(x_1, ..., x_i', ..., x_j, ..., x_D) \wedge \\ f(x_1, ..., x_i, ..., x_j', ..., x_D) > f(x_1, ..., x_i', ..., x_j', ..., x_D). \end{split}$$

#### **Definition 2.** (Interdependent Sub-problems)

Given two arbitrary sub-problems, they are said to be interdependent if at least one variable in a sub-problem is interacting with at least one variable in the other sub-problem.

Intuitively, two interacting variables of the 2-D Schwefel function are depicted in Figure 1, where the rank of two partial solutions  $x_1$  and  $x_1'$  in the first dimension varies by fixing different values in the second dimension, i.e.,  $x_2$  and  $x_2'$ .

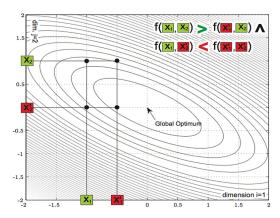


Figure 1: Two interacting variables of the 2-D Schwefel function.

If a problem is separable, where no interdependency exists between sub-problems, partial solutions can be complemented by arbitrary identical values in  $\mathcal{X}$ , without perturbing their rank. However, for many non-separable problems, the sub-problems are interdependent, where the rank of partial solutions significantly relies on their complements. As a result, the search direction has a close relation to the choice of complements. Hence, the choice of complements to partial solutions should be carefully addressed. Otherwise, the search process in a sub-problem will run the risk of being misled, eventually resulting in an ineffective search.

Unfortunately, we find that how to accurately complement partial solutions is a difficult optimization problem.

**Lemma 1.** (The Difficulty of Accurately Complementing) Given a set of partial solutions in the  $d_i$ -dimensional i-th sub-problem, let  $\widehat{D}_i$  be the set of all D variables except the ones in the i-th sub-problem,  $|\widehat{D}_i|$  be the cardinality of  $\widehat{D}_i$ ,  $P_{\widehat{D}_i(j)}$  be the probability of fixing a correct value for the j-th variable in  $\widehat{D}_i$ , and P be the arithmetic mean of all  $P_{\widehat{D}_i(j)}$ ,  $j=1,...,D-d_i$ , then the probability of accurately complementing those partial solutions so that they can be correctly ranked, denoted as  $P_i$ , is:

$$P_{i} = \prod_{j=1}^{|\widehat{D}_{i}|} P_{\widehat{D}_{i}(j)} \le \left(\frac{\sum_{j=1}^{|\widehat{D}_{i}|} P_{\widehat{D}_{i}(j)}}{|\widehat{D}_{i}|}\right)^{|\widehat{D}_{i}|} = P^{D-d_{i}} \quad (1)$$

*Proof.* By Definition 1, we learn that fixing the correct values for variables are independent events. By Definition 2, we know that variables in the same sub-problem are non-interacting. Thus, we can directly have  $P_i = \prod_{j=1}^{|\widehat{D_i}|} P_{\widehat{D_i}(j)}$ . Finally, according to the AM-GM Inequality, we have Eq.(1).

Notice that,  $P_{\widehat{D_i}(j)}=1$  if the  $\widehat{D_i}(j)$ -th variable is not interacting with any variable in the i-th sub-problem. Lemma 1 shows that the computational cost of accurately complementing partial solutions increases exponentially with the number of variables that are interacting with the current sub-problem. Hence, the interdependent sub-problems cannot be accurately conquered within a reasonable time budget.

## 3 Divide and Approximate Conquer

## 3.1 The Accurate Complement

According to Lemma 1, only the brute-force method is applicable to accurately complement partial solutions on interdependent sub-problems. However, the required computational costs are beyond being acceptable. On the other hand, it is an effective way to derive the mathematical formulation of a problem by observing its corresponding brute-force method, as such a method usually has to scan the whole solution space and thus reflects the problem characteristics naturally.

The core idea of the brute-force method is mathematically described as follow ( *the maximization case is considered*):

$$\forall \mathbf{x}_i \in S_i : f(\mathbf{x}^*) = f(\mathbf{x}_i^*, \mathbf{x}_r^*) = \max_{\mathbf{x}_r \in S_r} f(\mathbf{x}_i^*, \mathbf{x}_r) \ge \max_{\mathbf{x}_r \in S_r} f(\mathbf{x}_i, \mathbf{x}_r), \quad (2)$$

where  $\mathbf{x}_i$  denotes a partial solution in the *i*-th sub-problem, i=1,...,M, and  $\mathbf{x}_{\mathbf{r}}$  denotes its candidate complement, where  $\mathbf{r}=[1,...,i-1,i+1,...,M]$  denotes all the sub-problems except the *i*-th one.  $S_i$  and  $S_{\mathbf{r}}$  denote the corresponding subspace of the solution space subjecting to  $|S_i| \cdot |S_{\mathbf{r}}| = |S|$ , where  $S \subseteq \mathcal{X}$  and  $\mathbf{x}^* = (\mathbf{x}_i^*, \mathbf{x}_{\mathbf{r}}^*)$  is the optimal solution in S.

Eq.(2) states a fact that the correct rank of partial solutions can be obtained by comparing their best function values among all combinations with all possible complements. On this basis, we mathematically define the problem of accurately complementing partial solutions as follow:

## **Definition 3.** (The Accurate Complement)

Given arbitrary  $\mathbf{x}_i \in S_i$ , a complement  $\mathbf{x}_{\mathbf{r}}^{\dagger} \in S_{\mathbf{r}}$  is said to be the accurate complement  $\iff \mathbf{x}_{\mathbf{r}}^{\dagger} = \underset{\mathbf{x}_{\mathbf{r}} \in S_{\mathbf{r}}}{\operatorname{max}} f(\mathbf{x}_i, \mathbf{x}_{\mathbf{r}})$ .

Notice that, every partial solution has its own accurate complement. To identify the accurate complement, the combinations of  $\mathbf{x}_i$  and all possible complements in  $S_{\mathbf{r}}$  should be evaluated by f, among which the complement with the largest function value is chosen.

## 3.2 DAC: an approximate approach to DC

In fact, Definition 3 allows us to approximate the accurate complements of partial solutions with only a limited set of candidate complements  $S'_{\mathbf{r}} \subseteq S_{\mathbf{r}}$ . That is,

$$\mathbf{x}_{\mathbf{r}}^{\dagger} = \underset{\mathbf{x}_{\mathbf{r}} \in S_{\mathbf{r}}}{\operatorname{max}} f(\mathbf{x}_{i}, \mathbf{x}_{\mathbf{r}}) \succeq \widetilde{\mathbf{x}}_{\mathbf{r}}^{\dagger} = \underset{\mathbf{x}_{\mathbf{r}}' \in S_{\mathbf{r}}' \subseteq S_{\mathbf{r}}}{\operatorname{max}} f(\mathbf{x}_{i}, \mathbf{x}_{\mathbf{r}}').$$
(3)

where  $x_r^\dagger \succeq \widetilde{x}_r^\dagger$  means that  $x_r^\dagger$  is more accurate than  $\widetilde{x}_r^\dagger$ , since,

$$\max_{\mathbf{x_r} \in S_r} f(\mathbf{x}_i, \mathbf{x_r}) \ge \max_{\mathbf{x_r'} \in S_r' \subseteq S_r} f(\mathbf{x}_i, \mathbf{x_r'}). \tag{4}$$

Based on Eq.(4), it is reasonable to assume that good *approximate complements* will perturb the rank of partial solutions slightly. Eq.(3) thus gives rise to the proposed Divide and Approximate Conquer (DAC), as shown in Algorithm 1.

DAC shares the same framework as the basic DC. The only difference between them is that: while complementing a partial solution, DAC always selects the complement associated with the largest function value among a given set of candidate

## Algorithm 1 DAC $(f, T_{max}, N)$

```
1: Randomly initialize N solutions \mathbf{x}_{1:N}.
    2: Divide f into M sub-problems.
             For t = 1 to T_{max}
   4:
                      For i=1 to M
                             or i = 1 to M
\mathbf{x}'_{1:N;i} = \text{SearchOperator}(\mathbf{x}_{1:N;i}).
For j = 1 to N
\widetilde{\mathbf{x}}^{\dagger}_{j;\mathbf{r}} = \underset{\mathbf{x}_{\mathbf{r}} \in \mathbf{x}_{1:N;\mathbf{r}}}{\operatorname{arg}} \max_{\mathbf{x}} f(\mathbf{x}_{j;i}, \mathbf{x}_{\mathbf{r}}).
\widetilde{\mathbf{x}}'^{\dagger}_{j;\mathbf{r}} = \underset{\mathbf{x}_{\mathbf{r}} \in \mathbf{x}_{1:N;\mathbf{r}}}{\operatorname{arg}} \max_{\mathbf{x}} f(\mathbf{x}'_{j;i}, \mathbf{x}_{\mathbf{r}}).
EndFor
   5:
    6:
    7:
   8:
   9:
                              \mathbf{x}_{1:N;i} \leftarrow \{\mathbf{x}_{1:N;i}; \mathbf{x}_{1:N:i}' \mid \widetilde{\mathbf{x}}_{1:N:r}^{\dagger}; \widetilde{\mathbf{x}'}_{1:N:r}^{\dagger} \} .
 10:
                              \mathbf{x}_{1:N;\mathbf{r}} \leftarrow \{\widetilde{\mathbf{x}}_{1:N;\mathbf{r}}^{\dagger}; \widetilde{\mathbf{x}'}_{1:N;\mathbf{r}}^{\dagger}\}.
 11:
                      EndFor
 12:
13: EndFor
14: Output the best solution found so far.
```

complements. Specifically, DAC works by first randomly initializing N solutions  $\mathbf{x}_{1:N}$  (step 1). The problem f is decomposed into M sub-problems with a certain decomposition strategy (step 2). After that, without loss of generality, let us consider the i-th sub-problem. N new partial solutions  $\mathbf{x}'_{1:N;i}$  are generated by applying some search operator to the current ones, i.e.,  $\mathbf{x}_{1:N;i}$  (step 5). To identify the approximate complement to the j-th partial solution  $\mathbf{x}_{j;i}$ , j=1,...,N, all the combinations of  $\mathbf{x}_{j;i}$  and the vectors of partial solutions in other sub-problems  $\mathbf{x}_{1:N;\mathbf{r}}$  will be evaluated by f. The vector associated with the largest function value is chosen as the approximate complement  $\widetilde{\mathbf{x}}_{j;\mathbf{r}}^{\dagger}$  to  $\mathbf{x}_{j;i}$  (step 7). The same strat-

egy is used to obtain the approximate complement  $\widetilde{\mathbf{x}'}_{j;\mathbf{r}}^{\dagger}$  to the j-th new partial solution  $\mathbf{x}'_{j;i}$  (step 8). After that, according to a certain selection criterion, N partial solutions will be remained for the next iteration (step 10). Notice that, the selection of a partial solution is conditioned by its corresponding approximate complement. At last, for the j-th selected partial solution  $\mathbf{x}_{j;i}$ , its corresponding partial solutions in the rest sub-problems  $\mathbf{x}_{j;\mathbf{r}}$  will be replaced with its approximate complement for further optimizing (step 11).

As a result, DAC consumes  $2MN^2$  Function Evaluations (FEs) in each iteration, which is a significant reduction to the brute force method. Meanwhile, albeit the computational time is cut down, the convergence of DAC is still guaranteed.

## **Lemma 2.** (The Convergence of DAC)

Given a search algorithm that can converge to the global optimum of each sub-problem (regarded as independent problems), DAC can approximately converge to the global optimum  $\mathbf{x}^*$  of the original problem.

*Proof.* With out loss of generality, let us consider the function values of the j-th solution at the t-th and t+1-th iteration, i.e.,  $f(\mathbf{x}_j^t)$  and  $f(\mathbf{x}_j^{t+1})$ . For the i-th sub-problem, we have:

$$f(\mathbf{x}_{j;i}^{t}, \mathbf{x}_{j;\mathbf{r}}^{t}) \leq \max_{\mathbf{x}_{\mathbf{r}}^{t} \in \mathbf{x}_{1:N;\mathbf{r}}^{t}} f(\mathbf{x}_{j;i}^{t}, \mathbf{x}_{\mathbf{r}}^{t}) \leq \max_{\mathbf{x}_{\mathbf{r}}^{t} \in \mathbf{x}_{1:N;\mathbf{r}}^{t}} f(\mathbf{x}_{j;i}^{t+1}, \mathbf{x}_{\mathbf{r}}^{t}).$$
(5)

where the first '\(\leq'\) indicates the procedure of identifying the approximate complements for the current partial solutions (step 7 in Algorithm 1), and the second '\(\leq'\) represents the procedures of generating new partial solutions (step 5 in Algorithm 1), identifying their approximate complements (step 8 in Algorithm 1), and selecting better ones from candidate partial solutions, conditioned by their approximate complements (step 10 in Algorithm 1).

Then by repeating Eq.(5) for i = 1, ..., M, we have that:

$$f(\mathbf{x}_{i;1}^t, ..., \mathbf{x}_{i;M}^t) \le f(\mathbf{x}_{i;1}^{t+1}, ..., \mathbf{x}_{i:M}^{t+1}),$$
 (6)

which means the function value of the j-th solution  $f(\mathbf{x}_j^t)$  monotonically increases with the iteration index t.

Note that, the equality of Eq.(6) holds in two cases:

- 1) The approximate complement of a partial solution happens to be its corresponding partial solutions in the rest sub-problems, i.e.,  $\mathbf{x}_{j:\mathbf{r}} = \widetilde{\mathbf{x}}_{i:\mathbf{r}}^{\dagger}$ ;
- 2) The search algorithm fails to produce new better solutions, i.e.,  $\max_{\mathbf{x_r} \in S_{\mathbf{r}}} f(\mathbf{x}_{j;i}, \mathbf{x_r}) \geq \max_{\mathbf{x_r} \in S_{\mathbf{r}}} f(\mathbf{x}'_{j;i}, \mathbf{x_r})$ .

The first case actually explains the term "approximate" in DAC, as it happens at a probability of at least  $\frac{1}{N}$ . Hence, if the sub-problems optimizer of DAC can optimally solve each sub-problem separately, the global optimum value  $f(\mathbf{x}^*)$  can be approximately approached by DAC.

#### 3.3 DAC-HC: an instantiation of DAC

An instantiation of DAC is presented to illustrate the detail steps of a DAC algorithm and for further empirical studies.

To instantiate DAC, both the decomposition strategy and the sub-problems optimizer should be specified. In order to highlight the advantages of the DAC framework further in empirical studies, the improvement of performance introduced by these two specified components should be kept to minimal. On this basis, we first decompose a problem f via random grouping [Yang  $et\ al.$ , 2008a]. That is, in the beginning of each iteration, M equal-sized sub-problems are randomly generated. For the sub-problems optimizer, a Parallel Hill Climbing (PHC) method is employed, which thus yields the DAC-Hill Climbing (DAC-HC). The DAC-HC conducts N RLS processes on each sub-problem. Specifically, at each iteration of the i-th sub-problem, the j-th RLS produces one new partial solution  $\mathbf{x}'_{j;i}$  by applying the Gaussian mutation operator to the current partial solution  $\mathbf{x}_{j;i}$ , using Eq.(7):

$$\mathbf{x}'_{j;i} = \mathbf{x}_{j;i} + \mathbf{I} \cdot \mathcal{N}(0, \sigma_{j;i}). \tag{7}$$

where  $\mathcal{N}(0, \sigma_{j;i})$  denotes a Gaussian random variable with zero mean and standard deviation  $\sigma_{j;i}$ , and **I** is the identity matrix of size  $d_i$ . Generally, the value of  $\sigma_{j;i}$  represents the search step-size that can be adaptively varied during the search and may also be distinct over RLSs or even dimensions. To keep it simple, all RLSs in DAC-HC are initially set to the same search step-size, i.e., 1.00. After that, each search step-size is adapted at every iteration in terms of the 1/5 successful rule [Kern *et al.*, 2004], using Eq.(8):

$$\sigma_{j;i} = \sigma_{j;i} \times \exp^{\frac{1}{\sqrt{D+1}}} (\mathbb{I}_{f(\mathbf{x}_{j;i}', \widetilde{\mathbf{x}}_{j;r}^{\dagger}) \ge f(\mathbf{x}_{j;i}, \widetilde{\mathbf{x}}_{j;r}^{\dagger})} - \frac{1}{5}), \quad (8)$$

```
Algorithm 2 DAC-HC(f, T_{max}, N, M)
```

```
1: Randomly initialize N solutions \mathbf{x}_{1:N}.
  2: For t = 1 to T_{max}
  3:
                Randomly divide f into M equal-sized sub-problems.
  4:
                For i=1 to M
  5:
                       For j = 1 to N
                            \mathbf{x}'_{j;i} = \mathbf{x}_{j;i} + \mathbf{I} \cdot \mathcal{N}(0, \sigma_{j;i}).
\mathbf{\tilde{x}}^{\dagger}_{j;\mathbf{r}} = \underset{\mathbf{x}_{\mathbf{r}} \in \mathbf{x}_{1:N;\mathbf{r}}}{\arg \max} f(\mathbf{x}_{j;i}, \mathbf{x}_{\mathbf{r}}).
  6:
  7:
                             \sigma_{j;i} = \sigma_{j;i} \times \exp^{\frac{1}{\sqrt{D+1}}} (\mathbb{I}_{f(\mathbf{x}_{j:i}',\widetilde{\mathbf{x}}_{j:r}^\dagger) \geq f(\mathbf{x}_{j:i},\widetilde{\mathbf{x}}_{j:r}^\dagger)} - \frac{1}{5}).
  8:
                            \mathbf{x}_{j;i} \leftarrow \{\mathbf{x}_{j;i}, \mathbf{x}_{j;i}' \mid \widetilde{\mathbf{x}}_{j;\mathbf{r}}^{\dagger}\}.
  9:
                             \mathbf{x}_{j;\mathbf{r}} \leftarrow \widetilde{\mathbf{x}}_{i;\mathbf{r}}^{\dagger}.
10:
                       EndFor
11:
12:
                EndFor
13: EndFor
14: Output the best solution found so far.
```

where  $\mathbb{I}_a$  is a indicator function that returns 1 if a is true and 0 otherwise.

During the selection procedure, each new partial solution in PHC only competes with its corresponding old one for survival. Based on this one-on-one relation, we further reduce the FEs consumption of DAC-HC to a half of DAC, i.e.,  $MN^2$ . This is conducted by letting two competing partial solutions share the same approximate complement. The reason behind this is that, by adopting RLSs with small search stepsizes, pairwise partial solutions can be close to each other in the solution space, in which case their respective approximate complements may also be similar. Lastly, the pseudo-code of DAC-HC is given in Algorithm 2 for illustration.

## 4 Empirical Studies

DAC is proposed for solving non-separable high-dimensional optimization problems. That is where the empirical studies should concentrate on to verify the effectiveness of DAC-HC. For this purpose, two sets of non-separable high-dimensional optimization problems are employed.

#### 4.1 Varied numbers of interacting variables tests

The first set of problems is based on the fully non-separable functions, i.e., Schwefel's 1.2 and Rosenbrock [Tang et al., 2009], which are formulated as:  $f_{\rm sch}(\mathbf{x}) = \sum_{i=1}^D (\sum_{j=1}^i x_j)^2$  and  $f_{\rm ros}(\mathbf{x}) = \sum_{i=1}^{D-1} [100(x_i^2 - x_{i+1})^2 + (x_i - 1)^2]$ . In these two functions, all variables are interacting. Meanwhile, it has also been observed that, in many real-world problems, only parts of variables are interacting [Friesen and Domingos, 2015]. Hence, it is a necessity to test DAC-HC with varied numbers of interacting variables. For this purpose, we further consider three problems that combine Schwefel's 1.2 function and the fully separable sphere function, i.e.,  $f_{\rm sph}(\mathbf{x}) = \sum_{i=1}^D x_i^2$ , in different formations. The dimensionality is set to 1000 for all 5 problems. All variables are randomly perturbed to avoid any potential bias. All problems are expected to be minimized to the global optimal value 0.00.

Table 1: The formulations of 5 tested functions.  $\mathbf{z} = \mathbf{x} - \mathbf{o}$  is a shifted solution and  $\mathbf{o}$  is the global optimum.  $\mathbf{P}$  is a random permutation of [1, ..., D] and D = 1000, m = 50.

$$f_{1}(\mathbf{x}) = f_{\text{sch}}(\mathbf{z}(P_{1}:P_{m})) * 10^{6} + f_{\text{sph}}(\mathbf{z}(P_{m+1}:P_{D}))$$

$$f_{2}(\mathbf{x}) = \sum_{k=1}^{\frac{D}{2m}} f_{\text{sch}}(\mathbf{z}(P_{(k-1)*m+1}:P_{k*m}))$$

$$+ f_{\text{sph}}(\mathbf{z}(P_{\frac{D}{2}+1}:P_{D}))$$

$$f_{3}(\mathbf{x}) = \sum_{k=1}^{\frac{D}{m}} f_{\text{sch}}(\mathbf{z}(P_{(k-1)*m+1}:P_{k*m}))$$

$$f_{4}(\mathbf{x}) = f_{\text{sch}}(\mathbf{z})$$

$$f_{5}(\mathbf{x}) = \sum_{k=1}^{\frac{D}{m}} f_{\text{ros}}(\mathbf{z}(P_{(k-1)*m+1}:P_{k*m}))$$

Specifically,  $f_1(\mathbf{x})$  consists of a group of 50 interacting variables and 950 independent variables.  $f_2(\mathbf{x})$  has 10 groups of 50 interacting variables and 500 independent variables.  $f_3(\mathbf{x})$  and  $f_5(\mathbf{x})$  compose of 20 groups of 50 interacting variables.  $f_4(\mathbf{x})$  involves a group of 1000 interacting variables. The detailed formulations are listed in Table 1.

On these five problems, two groups of comparisons are conducted for different purposes.

## Advantages of DAC-HC over existing approaches

In the first group of comparison, DAC-HC is compared with CMA-ES [Hansen and Ostermeier, 2001], RESOO [Qian and Yu, 2016], and DECC-I [Omidvar et al., 2014], which are representatives of three basic ideas for high-dimensional optimization: the straightforward method, dimensionality reduction, and DC, respectively. Specifically, CMA-ES is widely endorsed as a powerful global optimizer that has been applied in many aspects. Here the basic version of CMA-ES is utilized. RESOO is a recently proposed approach built on the Random Embedding for reducing dimensionality. It should be noted that, DECC-I is not an algorithm for black-box optimization. It is an ideal approach that perfectly decomposes the problems using the priori knowledge of functions. Hence, all sub-problems of DECC-I are independent, while it is not the case for DAC-HC. Besides, the sub-problems of DECC-I are optimized by a variant of Differential Evolution [Yang et al., 2008b], which has been empirically observed more advanced than the employed RLSs [Tang et al., 2016]. On this basis, if DAC-HC outperforms DECC-I, it is reasonable to infer that the proposed DAC facilitates DC on non-separable high-dimensional optimization problems.

All algorithms are repeated 25 runs on each problem to diminish the noise introduced by their randomized search essence. The time budget for each run is set to 3e6 FEs. CMA-ES is parameterless that no parameter needs to be specified. For RESOO, after some coarse-tuning, the probability  $\eta$  is set to 1/3, the restart times is set to 10, and the reduced dimensionality is set to 100. For DECC-I, the only parameter, i.e., population size N, is set to 100 as Omidvar et~al.~ [2014] suggested. For DAC-HC, two parameters should be specified, i.e., the number of RLSs N and the number of sub-problems M. Recall that, the parameter N>1 generally influences the approximate ability of DAC. To test the extreme case of DAC-HC, it is set to 2. To gain a relatively fair comparison with RESOO, we set M=10 so that each sub-problems optimizer faces a 100-dimensional problem as RESOO does.

Table 2: The mean and standard derivation of final outputs.

Function		CMA-ES	RESOO	DECC-I	DAC-HC
$f_1$	Mean	1.35e+09	2.52e+11	2.97e+01	0.00e+00
	Std	3.29e+08	3.62e+10	8.59e+01	0.00e+00
$f_2$	Mean	0.00e+00	1.28e+07	1.48e+03	1.55e+00
	Std	0.00e+00	9.41e+05	4.28e+02	1.25e+00
$f_3$	Mean	0.00e+00	3.62e+07	3.91e+04	7.78e+02
	Std	0.00e+00	2.89e+06	2.75e+03	7.12e+02
$f_4$	Mean	2.87e + 06	7.80e+07	1.74e+06	3.93e+05
	Std	6.61e + 05	7.10e+06	9.54e+04	2.52e+04
$f_5$	Mean	3.36e+03	1.41e+12	1.17e+03	1.13e+03
	Std	1.81e+03	8.02e+09	9.66e+01	2.32e+02

The mean and standard derivation of the final outputs in 25 runs are shown in Table 2. A gray cell indicates an algorithm achieves the best mean value on a problem, while a light gray cell indicates a second place. DAC-HC outperforms all the compared algorithms on  $f_1$ ,  $f_4$  and  $f_5$ . On  $f_2$  and  $f_3$ , though slightly inferior to CMA-ES, DAC-HC performs significantly better than DECC-I and RESOO. Since DAC-HC dominates DECC-I on all five problems, the effectiveness of DAC for promoting DC on non-separable high-dimensional optimization problems is confirmed. Besides, Lemma 1 is verified by observing that the solution qualities of DAC-HC deteriorates as the number of interacting variables increases. RESOO performs poorly because the tested problems are irreducible.

## **Empirical support to the convergence of DAC**

In the second group of comparison, the PHC with 2 RLSs is compared. The only difference between PHC and DAC-HC is that PHC complements a partial solution  $\mathbf{x}_{j;i}$  merely with the corresponding partial solutions in the rest sub-problems  $\mathbf{x}_{j;\mathbf{r}}$ . On this basis, PHC does not satisfy Eq.(5), and its convergence is not guaranteed. The experimental protocol is set the same to the first group of comparison.

The convergence rates of both algorithms are shown in Figure 2, where the x-axis denotes the FEs and y-axis denotes the logarithm of function values. It can be seen that, DAC-HC always converges faster than PHC. Specially, log-linear convergence of DAC-HC is observed on the first two sphere function based problems. It should be noted that, the employed RLS has also been theoretically proved to converge log-linearly on the sphere function [Jebalia *et al.*, 2008]. This coincidence actually supports the convergence of DAC stated in Lemma 2. Comparatively, the convergence rates of PHC are heavily retarded due to the unfit complements to partial solutions.

## 4.2 Hyper-parameter tuning for multi-class SVMs

Given a set of labelled data  $\{\mathbf{x}_i, y_i\}_{i=1}^l$ , where  $\mathbf{x}_i \in R^n$ , the classification task is to train a classifier in terms of  $\{\mathbf{x}_i, y_i\}_{i=1}^l$  to predict the labels of incoming data. Support Vector Machines (SVMs) [Vapnik, 1998] is often considered as a family of powerful tools for classification. Here the SVM with linear kernel is considered. Let  $\mathbf{y} \in \{-1, +1\}$ , SVM requires to fine-tune three parameters  $\mathbf{w}, b, \lambda$  by solving the following optimization problem:  $\min_{\mathbf{w}, b, \lambda} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \lambda \sum_{i=1}^l \xi_i$ , subject to  $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \wedge \xi_i \geq 0$ . Notice that,  $\lambda > 0$  is a hyper-parameter supplied by the user, which

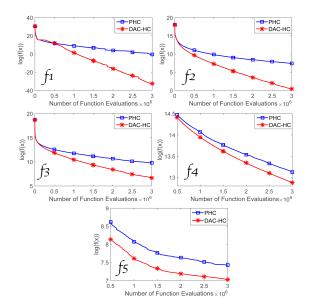


Figure 2: The convergence rates of DAC-HC and PHC. The x-axis denotes the FEs (time) and y-axis denotes  $\log(f(\mathbf{x}))$ .

penalizes the error vector  $\boldsymbol{\xi}$ . When dealing with a multiclass classification problem, a typical idea is to divide it into multiple binary classification problems by adopting the one-on-one strategy. Let K be the number of class, then we have  $\binom{K}{2} = \frac{K(K-1)}{2}$  binary classifiers to train, resulting in  $\frac{K(K-1)}{2}$  hyper-parameters to tune.

Of course, a simple strategy of specifying the same value of  $\lambda$  for all binary classifiers works well [Chang and Lin, 2011]. It is also an intuition that varied  $\lambda$  can facilitate a multi-class SVM better. On this basis, a potentially high-dimensional optimization problem needs to be solved for more advanced performance. Recall that, a final output of a multi-class SVM is based on the majority voting. Due to the overfitting risk on each binary classifier, the votes may introduce interdependencies in between. To sum up, the problem of hyperparameter tuning for multi-class SVMs is non-separable and high-dimensional in essence.

We thus apply the proposed DAC-HC to deal with it. RE-SOO and CMA-ES are again included as the compared algorithms. Since the ideal decomposition is no longer applicable for this problem, DECC-I will not be compared with DAC-HC. The grid search is also tested but on the assumption that all binary classifiers share the same value of  $\lambda$ . That is , the grid search actually solves a 1-dimensional problem. Three data sets, i.e., usps [Hull, 1994], news20 [Lang, 1995], and letter [Hsu and Lin, 2002], are used for comparison, which contain 10, 20, 26 classes in each and yield three problems with 45, 190, 325 dimensions, respectively. All the features in each dataset are scaled to [-1, 1] or [0, 1]. The solution space for each binary classifier is bounded as  $[10^{-3}, 10^{2}]$ . All algorithms are repeated for 20 runs on each problem, except the deterministic grid search. For each run, the hyper-parameters are tuned on the training set with the 5-fold cross-validation. The higher the accuracy is, the better the hyper-parameters

Table 3: Testing accuracies of tuned multi-class SVM.

DataSet	GridSearch	RESOO	CMA-ES	DAC-HC
	93.92%	94.38%	93.33%	94.60%
usps	-	$\pm 0.37\%$	$\pm 0.24\%$	$\pm 0.04\%$
news20	85.16%	85.43%	84.34%	85.40%
news20	-	$\pm 0.21\%$	$\pm 0.17\%$	$\pm 0.11\%$
letter	85.12%	85.36%	84.03%	85.72%
ieiier	-	$\pm 0.12\%$	$\pm 1.44\%$	$\pm 0.24\%$

are supposed to be. The best hyper-parameters obtained in a run will be tested on the testing set, and the testing accuracy is regarded as the performance of an algorithm in a run. The time budget for the training phase in each run is set to 100 FEs. For RESOO, the probability  $\eta$  is set to 1/3, the restart times is set to 2, and the reduced dimensionality is set to 1/5 to the original dimensionality, i.e., 9, 38, 65, respectively. For the grid search, 100 candidate solutions are uniformly selected over the solution space. For DAC-HC, the parameters N and M are set to 2 and 3, respectively.

Table 3 lists the mean and standard derivation of the testing accuracies of the multi-class SVMs tuned by each algorithm. It can be seen that, both RESOO and DAC-HC outperform the grid search, although the improvement is marginal. Hence, it can be inferred that tuning multiple hyper-parameters, rather than one shared hyper-parameter, is beneficial to the multiclass SVMs. DAC-HC outperforms all the compared algorithms on the usps and letter datasets. Although it shows slightly lower accuracy than RESOO on the news20 dataset, a more stable behavior is observed as its standard derivation is smaller. CMA-ES is inferior to the grid search on all three problems. This phenomenon suggests that, for tuning multiple hyper-parameters for multi-class SVMs, an appropriate optimization approach should be employed. Otherwise, it would be better to tune just one shared hyper-parameter. It is also worthwhile to notice that, CMA-ES does not adopt any special treatment for high-dimensional optimization.

## **5** Conclusions and Future Directions

This work investigated the Divide and Conquer idea on high-dimensional black-box optimization problems. We found that the interdependent sub-problems after decomposition actually cannot be accurately conquered. Instead, we proposed the Divide and Approximate Conquer (DAC) to solve each sub-problem approximately. The convergence of DAC was proved and empirically supported. For empirical studies, a simple instantiation of DAC, i.e., DAC-HC, was also proposed. The advantages of DAC-HC over existing representative approaches were verified on two sets of non-separable high-dimensional problems.

For future work, we are interested in:

- Promoting the ability of DAC by adopting more advanced sub-problems optimizers.
- Theoretically analyzing the convergence rate of DAC.

#### References

[Chang and Lin, 2011] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines.

- ACM Transactions on Intelligent Systems and Technology, 2:27:1–27:27, 2011.
- [Chen et al., 2010] Wenxiang Chen, Thomas Weise, Zhenyu Yang, and Ke Tang. Large-scale global optimization using cooperative coevolution with variable interaction learning. In *Parallel Problem Solving from Nature*, *PPSN XI*, pages 300–309. Springer, 2010.
- [Friesen and Domingos, 2015] Abram L Friesen and Pedro Domingos. Recursive decomposition for nonconvex optimization. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence*, 2015.
- [Hansen and Ostermeier, 2001] Nikolaus Hansen and Andreas Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary computation*, 9(2):159–195, 2001.
- [Hsu and Lin, 2002] Chih-Wei Hsu and Chih-Jen Lin. A comparison of methods for multiclass support vector machines. *IEEE Transactions on Neural Networks*, 13(2):415–425, 2002.
- [Hull, 1994] Jonathan J Hull. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(5):550–554, 1994.
- [Jebalia *et al.*, 2008] Mohamed Jebalia, Anne Auger, and Pierre Liardet. Log-linear convergence and optimal bounds for the (1+ 1)-ES. In *Artificial Evolution*, pages 207–218. Springer, 2008.
- [Kabán et al., 2015] Ata Kabán, Jakramate Bootkrajang, and Robert John Durrant. Toward large-scale continuous eda: A random matrix theory perspective. Evolutionary computation, 2015.
- [Kern *et al.*, 2004] Stefan Kern, Sibylle D Müller, Nikolaus Hansen, Dirk Büche, Jiri Ocenasek, and Petros Koumoutsakos. Learning probability distributions in continuous evolutionary algorithms—a comparative review. *Natural Computing*, 3(1):77–112, 2004.
- [Kirkpatrick *et al.*, 1983] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, 1983.
- [Lang, 1995] Ken Lang. Newsweeder: Learning to filter netnews. In *Proceedings of the 12th international conference on machine learning*, pages 331–339, 1995.
- [Mahdavi *et al.*, 2015] Sedigheh Mahdavi, Mohammad Ebrahim Shiri, and Shahryar Rahnamayan. Metaheuristics in large-scale global continues optimization: A survey. *Information Sciences*, 295:407–428, 2015.
- [Omidvar *et al.*, 2014] Mohammad Nabi Omidvar, Xi-aodong Li, Yi Mei, and Xin Yao. Cooperative co-evolution with differential grouping for large scale optimization. *IEEE Transactions on Evolutionary Computation*, 18(3):378–393, 2014.
- [Qian and Yu, 2016] Hong Qian and Yang Yu. Scaling simultaneous optimistic optimization for high-dimensional non-convex functions with low effective dimensions. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence (AAAI 2016)*, Phoenix, AZ, 2016.

- [Tang et al., 2009] Ke Tang, Xiaodong Li, P. N. Suganthan, Zhenyu Yang, and Thomas Weise. Benchmark functions for the cec2010 special session and competition on large-scale global optimization. Technical report, Nature Inspired Computation and Applications Laboratory, 2009.
- [Tang et al., 2016] Ke Tang, Peng Yang, and Xin Yao. Negatively correlated search. *IEEE Journal on Selected Areas in Communications*, 34(3):542–550, March 2016.
- [Vapnik, 1998] Vladimir Vapnik. Statistical learning theory, volume 1. Wiley New York, 1998.
- [Wang et al., 2013] Ziyu Wang, Masrour Zoghi, Frank Hutter, David Matheson, and Nando De Freitas. Bayesian optimization in high dimensions via random embeddings. In Proceedings of the Twenty-Third international joint conference on Artificial Intelligence (IJCAI 2013), pages 1778–1784, Beijing, China, 2013. AAAI Press.
- [Yang et al., 2008a] Zhenyu Yang, Ke Tang, and Xin Yao. Large scale evolutionary optimization using cooperative coevolution. *Information Sciences*, 178(15):2985–2999, 2008.
- [Yang et al., 2008b] Zhenyu Yang, Ke Tang, and Xin Yao. Self-adaptive differential evolution with neighborhood search. In *IEEE Congress on Evolutionary Computation*, 2008 (*IEEE World Congress on Computational Intelligence*), pages 1110–1116. IEEE, 2008.