

Nikolaus Hansen and Stefan Kern

Evaluating the CMA Evolution Strategy on Multimodal  
Test Functions

Parallel Problem Solving from Nature, PPSN 2004

<http://www.springer.de/comp/lncs/index.html>  
© Springer-Verlag

# Evaluating the CMA Evolution Strategy on Multimodal Test Functions

Nikolaus Hansen and Stefan Kern

Computational Science and Engineering Laboratory (CSE Lab),  
Swiss Federal Institute of Technology (ETH) Zurich, Switzerland  
{nikolaus.hansen, skern}@inf.ethz.ch,  
<http://www.icos.ethz.ch/cse/>

**Abstract.** In this paper the performance of the CMA evolution strategy with rank- $\mu$ -update and weighted recombination is empirically investigated on eight multimodal test functions. In particular the effect of the population size  $\lambda$  on the performance is investigated. Increasing the population size remarkably improves the performance on six of the eight test functions. The optimal population size takes a wide range of values, but, with one exception, scales sub-linearly with the problem dimension. The global optimum can be located in all but one function. The performance for locating the global optimum scales between linear and cubic with the problem dimension. In a comparison to state-of-the-art global search strategies the CMA evolution strategy achieves superior performance on multimodal, non-separable test functions without intricate parameter tuning.

## 1 Introduction

The derandomized Evolution Strategy (ES) with Covariance Matrix Adaptation (CMA) [1] adapts the complete covariance matrix of the normal mutation (search) distribution. The CMA-ES exhibits several invariances. Hereunder are (a) invariance against order preserving (i.e. strictly monotonic) transformations of the objective function value; (b) invariance against angle preserving transformations of the search space (rotation, reflection, and translation) if the initial search point is transformed accordingly; (c) scale invariance if the initial scaling is chosen accordingly. Invariances are highly desirable: they imply uniform behavior on classes of functions and therefore generalizability of empirical results.

Originally designed for small population sizes, the CMA-ES was interpreted as a robust local search strategy [2]. It efficiently minimizes unimodal test functions [1] and in particular is superior on ill-conditioned and non-separable problems. It was successfully applied to a considerable number of real world problems.<sup>1</sup> In [3, 4] the CMA-ES was expanded by the so-called rank- $\mu$ -update. The

---

<sup>1</sup> See [www.icos.ethz.ch/software/evolutionary\\_computation/cmaapplications.pdf](http://www.icos.ethz.ch/software/evolutionary_computation/cmaapplications.pdf) for a list of references.

rank- $\mu$ -update exploits the information contained in large populations more effectively without affecting the performance<sup>2</sup> for small population sizes. It can reduce the time complexity of the strategy (i.e. the number of *generations* to reach a certain function value) from quadratic to linear [4]. A recent study [5] showed a surprisingly good performance of this CMA-ES on the multimodal Rastrigin function. Large populations and rank- $\mu$ -update were the prerequisites for this observation. Therefore, we empirically investigate the effect of the population size  $\lambda$  on the global search performance of the CMA-ES.

The remainder is organized as follows: In Sect. 2 we describe the CMA-ES using weighted recombination and rank- $\mu$ -update. In Sect. 3 test functions and methodology for the performance study are outlined. Section 4 examines the performance depending on the population size and compares the CMA-ES with other global search strategies. Sect. 5 gives a summary and conclusion.

## 2 The CMA-ES with Rank- $\mu$ -Update and Weighted Recombination

We thoroughly define the CMA-ES combining weighted recombination [1] and rank- $\mu$ -update of the covariance matrix [3, 4]. In this  $(\mu_{\text{W}}, \lambda)$ -CMA-ES the  $\lambda$  individuals (candidate solutions) of generation  $g + 1$  are generated according to

$$\mathbf{x}_k^{(g+1)} \sim \mathcal{N}\left(\langle \mathbf{x} \rangle_{\text{w}}^{(g)}, \sigma^{(g)2} \mathbf{C}^{(g)}\right), \quad k = 1, \dots, \lambda, \quad (1)$$

where  $\mathcal{N}(\mathbf{m}, \mathbf{C})$  denotes a normally distributed random vector with mean  $\mathbf{m}$  and covariance matrix  $\mathbf{C}$ .<sup>3</sup>

The recombination point  $\langle \mathbf{x} \rangle_{\text{w}}^{(g)} = \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}^{(g)}$  is the weighted mean of the selected individuals,  $w_i > 0$  for all  $i = 1 \dots \mu$  and  $\sum_{i=1}^{\mu} w_i = 1$ . The index  $i : \lambda$  denotes the  $i$ -th best individual. Setting all  $w_i$  to  $1/\mu$  is equivalent to intermediate (multi-)recombination. The adaptation of the mutation parameters consists of two parts: (i) adaptation of the covariance matrix  $\mathbf{C}^{(g)}$ , and (ii) adaptation of the global step size  $\sigma^{(g)}$ . The covariance matrix  $\mathbf{C}^{(g)}$  is adapted by the evolution path  $\mathbf{p}_c^{(g+1)}$  and by the  $\mu$  *weighted* difference vectors between the recent parents and  $\langle \mathbf{x} \rangle_{\text{w}}^{(g)}$ :

$$\mathbf{p}_c^{(g+1)} = (1 - c_c) \cdot \mathbf{p}_c^{(g)} + H_{\sigma}^{(g+1)} \sqrt{c_c(2 - c_c)} \cdot \frac{\sqrt{\mu_{\text{eff}}}}{\sigma^{(g)}} \left( \langle \mathbf{x} \rangle_{\text{w}}^{(g+1)} - \langle \mathbf{x} \rangle_{\text{w}}^{(g)} \right) \quad (2)$$

$$\begin{aligned} \mathbf{C}^{(g+1)} &= (1 - c_{\text{cov}}) \cdot \mathbf{C}^{(g)} + c_{\text{cov}} \frac{1}{\mu_{\text{cov}}} \mathbf{p}_c^{(g+1)} \left( \mathbf{p}_c^{(g+1)} \right)^{\text{T}} \quad (3) \\ &+ c_{\text{cov}} \cdot \left( 1 - \frac{1}{\mu_{\text{cov}}} \right) \sum_{i=1}^{\mu} \frac{w_i}{\sigma^{(g)2}} \left( \mathbf{x}_{i:\lambda}^{(g+1)} - \langle \mathbf{x} \rangle_{\text{w}}^{(g)} \right) \left( \mathbf{x}_{i:\lambda}^{(g+1)} - \langle \mathbf{x} \rangle_{\text{w}}^{(g)} \right)^{\text{T}}, \end{aligned}$$

<sup>2</sup> We define performance as the number of function evaluations needed to reach a certain function value.

<sup>3</sup> Note that  $\mathcal{N}\left(\langle \mathbf{x} \rangle_{\text{w}}^{(g)}, \sigma^{(g)2} \mathbf{C}^{(g)}\right) \sim \langle \mathbf{x} \rangle_{\text{w}}^{(g)} + \sigma^{(g)} \mathbf{B}^{(g)} \mathbf{D}^{(g)} \mathcal{N}\left(\mathbf{0}, \mathbf{I}\right)$ , see below.

where  $H_\sigma^{(g+1)} = 1$  if  $\frac{\|\mathbf{p}_\sigma^{(g+1)}\|}{\sqrt{1-(1-c_\sigma)^{2(g+1)}}} < (1.5 + \frac{1}{n-0.5})\mathbb{E}(\|\mathcal{N}(0, \mathbf{I})\|)$ , and 0 otherwise.  $\mu_{\text{eff}} = 1/\sum_{i=1}^\mu w_i^2$  denotes the ‘‘variance effective selection mass’’ and  $\mu_{\text{eff}} = \mu$  if  $w_i = 1/\mu$ . The weights  $w_i$  are used for the summation term in (3), a matrix with rank  $\min(\mu, n)$ . Parameter  $c_{\text{cov}} \approx \min(1, 2\mu_{\text{eff}}/n^2)$  determines the learning rate for the covariance matrix  $\mathbf{C}$ . The adaptation of the global step size  $\sigma^{(g+1)}$  is based on a ‘‘conjugate’’ evolution path  $\mathbf{p}_\sigma^{(g+1)}$ :

$$\begin{aligned} \mathbf{p}_\sigma^{(g+1)} &= (1 - c_\sigma) \cdot \mathbf{p}_\sigma^{(g)} \\ &\quad + \sqrt{c_\sigma(2 - c_\sigma)} \cdot \mathbf{B}^{(g)} \mathbf{D}^{(g)-1} \mathbf{B}^{(g)\text{T}} \cdot \frac{\sqrt{\mu_{\text{eff}}}}{\sigma^{(g)}} \left( \langle \mathbf{x} \rangle_{\mathbf{w}}^{(g+1)} - \langle \mathbf{x} \rangle_{\mathbf{w}}^{(g)} \right). \end{aligned} \quad (4)$$

The orthogonal matrix  $\mathbf{B}^{(g)}$  and the diagonal matrix  $\mathbf{D}^{(g)}$  are obtained through a principal component analysis of  $\mathbf{C}^{(g)}$ ;  $\mathbf{C}^{(g)} = \mathbf{B}^{(g)} \mathbf{D}^{(g)2} \mathbf{B}^{(g)\text{T}}$  (cf. [1]). The global step size  $\sigma^{(g+1)}$  obeys

$$\sigma^{(g+1)} = \sigma^{(g)} \cdot \exp \left( \frac{c_\sigma}{d_\sigma} \left( \frac{\|\mathbf{p}_\sigma^{(g+1)}\|}{\mathbb{E}(\|\mathcal{N}(0, \mathbf{I})\|)} - 1 \right) \right), \quad (5)$$

where  $\mathbb{E}(\|\mathcal{N}(0, \mathbf{I})\|) = \sqrt{2}\Gamma(\frac{n+1}{2})/\Gamma(\frac{n}{2}) \approx \sqrt{n}(1 - \frac{1}{4n} + \frac{1}{21n^2})$  is the expected length of  $\mathbf{p}_\sigma$  under random selection.

Initial values are  $\mathbf{p}_\sigma^{(0)} = \mathbf{p}_c^{(0)} = \mathbf{0}$  and  $\mathbf{C}^{(0)} = \mathbf{I}$ , while  $\mathbf{x}^{(0)}$  and  $\sigma^{(0)}$  are problem dependent. Default strategy parameter values are

$$\lambda = 4 + \lceil 3 \cdot \ln(n) \rceil, \quad \mu = \lfloor \lambda/2 \rfloor, \quad w_{i=1\dots\mu} = \frac{\ln(\mu+1) - \ln(i)}{\sum_{j=1}^\mu \ln(\mu+1) - \ln(j)}, \quad (6)$$

$$c_\sigma = \frac{\mu_{\text{eff}} + 2}{n + \mu_{\text{eff}} + 3}, \quad d_\sigma = 1 + 2 \max \left( 0, \sqrt{\frac{\mu_{\text{eff}} - 1}{n + 1}} - 1 \right) + c_\sigma, \quad c_c = \frac{4}{n + 4}, \quad (7)$$

$$\mu_{\text{cov}} = \mu_{\text{eff}}, \quad c_{\text{cov}} = \frac{1}{\mu_{\text{cov}}} \frac{2}{(n + \sqrt{2})^2} + \left( 1 - \frac{1}{\mu_{\text{cov}}} \right) \min \left( 1, \frac{2\mu_{\text{eff}} - 1}{(n + 2)^2 + \mu_{\text{eff}}} \right). \quad (8)$$

While  $1/c_\sigma$ , and  $1/c_c$  can be interpreted as memory time constants,  $d_\sigma$  is a damping parameter. Parameters from (7) and (8) are not meant to be in the users choice. A profound discussion of the strategy parameters is given in [1].

We consider weighted recombination to be more natural than intermediate recombination, because the ranking of all  $\lambda/2$  best individuals is fully regarded.<sup>4</sup> Nevertheless, to our experience weighted recombination, where  $\mu \approx \lambda/2$ , only slightly outperforms intermediate recombination, where  $\mu \approx \lambda/4$ .

### 3 Test Functions and Experimental Procedure

#### 3.1 Test Functions

The unconstrained multimodal test problems are summarized in Table 1. All con-

<sup>4</sup> Even the mating success in nature is not well described by two possible outcomes.

**Table 1.** Test functions to be minimized and initialization regions

Name	Function	Init
Ackley	$f_{\text{Ackley}}(\mathbf{x}) = 20 - 20 \cdot \exp\left(-0.2 \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}\right) + e - \exp\left(\frac{1}{n} \sum_{i=1}^n \cos(2\pi x_i)\right)$	$[1, 30]^n$
Bohachevsky	$f_{\text{Bohachevsky}}(\mathbf{x}) = \sum_{i=1}^{n-1} (x_i^2 + 2x_{i+1}^2) - 0.3 \cos(3\pi x_i) - 0.4 \cos(4\pi x_{i+1}) + 0.7$	$[1, 15]^n$
Griewank	$f_{\text{Griewank}}(\mathbf{x}) = \frac{1}{4000} \sum_{i=1}^n x_i^2 - \prod_{i=1}^n \cos\left(\frac{x_i}{\sqrt{i}}\right) + 1$	$[10, 600]^n$
Rastrigin	$f_{\text{Rastrigin}}(\mathbf{x}) = 10n + \sum_{i=1}^n (x_i^2 - 10 \cos(2\pi x_i))$	$[1, 5]^n$
Scaled Rastrigin	$f_{\text{RastScaled}}(\mathbf{x}) = 10n + \sum_{i=1}^n \left( (10^{\frac{i-1}{n-1}} x_i)^2 - 10 \cos(2\pi 10^{\frac{i-1}{n-1}} x_i) \right)$	$[1, 5]^n$
Skew Rastrigin	$f_{\text{RastSkew}}(\mathbf{x}) = 10n + \sum_{i=1}^n (y_i^2 - 10 \cos(2\pi y_i))$ , with $y_i = \begin{cases} 10 \cdot x_i & \text{if } x_i > 0, \\ x_i & \text{otherwise} \end{cases}$	$[1, 5]^n$
Schaffer	$f_{\text{Schaffer}}(\mathbf{x}) = \sum_{i=1}^{n-1} (x_i^2 + x_{i+1}^2)^{0.25} \cdot [\sin^2(50 \cdot (x_i^2 + x_{i+1}^2)^{0.1}) + 1.0]$	$[10, 100]^n$
Schwefel	$f_{\text{Schwefel}}(\mathbf{x}) = 418.9828872724339 \cdot n - \sum_{i=1}^n x_i \cdot \sin(\sqrt{ x_i })$	$[-500, 300]^n$

sidered functions have a high number of local optima, are scalable in the problem dimension, and have a minimal function value of 0. The known global minimum is located at  $\mathbf{x} = \mathbf{0}$ , except for the Schwefel function, where the global minimum within  $[-500, 500]^n$  equals 420.96874636 in each coordinate. Additional bounds are implemented for  $f_{\text{Schwefel}}$  (in  $[-500, 500]^n$ ) and  $f_{\text{Ackley}}$  (in  $[-30, 30]^n$ ) by adding a quadratic penalty term. E.g.,  $f_{\text{Schwefel}}(\mathbf{x}) + 10^4 \cdot \sum_{i=1}^n \theta(|x_i| - 500) x_i^2$  is minimized, where  $\theta(\cdot)$  is the Heaviside function. The skew Rastrigin function was proposed by [6] to be deceptive for the CMA-ES.

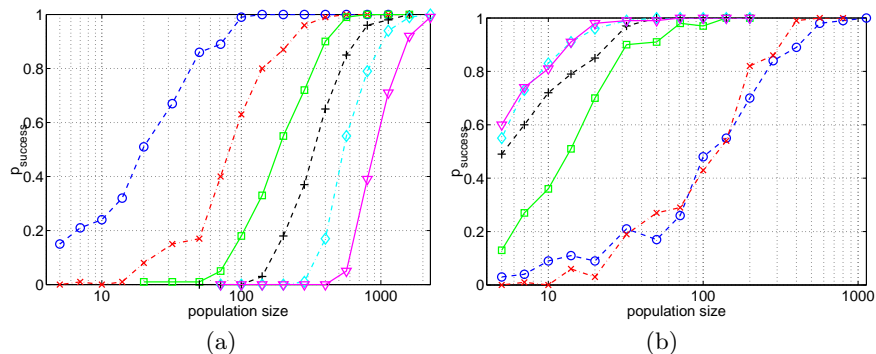
Besides  $f_{\text{RastSkew}}$  and  $f_{\text{Schwefel}}$ , the functions are point symmetrical around the global optimum. To avoid an easy exploitation of the symmetry, we suggest non-symmetrical initialization intervals, see Table 1. The Rastrigin functions and  $f_{\text{Schwefel}}$  are additively separable, while  $f_{\text{Ackley}}$  and  $f_{\text{Bohachevsky}}$  are separable, in that the global optimum can be located by optimizing each variable independently. Recall, that both is not exploited by the CMA-ES, because *all results of the CMA-ES are invariant under orthogonal transformations (rotations) of the coordinate system, given accordingly transformed initial intervals.*

### 3.2 Experimental procedure

The performance of the CMA-ES is tested for dimensions  $n = [2, 5, 10, 20, 40, 80]$ .<sup>5</sup> All runs are performed with the default strategy parameter setting given in Sect. 2, except for the population size  $\lambda$ .<sup>6</sup> Starting from  $\lambda = 5$ , the population

<sup>5</sup> For simulations we used the MATLAB code `cmaes.m`, Version 2.24, available from [http://www.icos.ethz.ch/software/evolutionary\\_computation/cma](http://www.icos.ethz.ch/software/evolutionary_computation/cma).

<sup>6</sup> Note that  $\mu$  is chosen dependently on  $\lambda$  and further parameters depend on  $\mu$ .



**Fig. 1.** Success rate to reach  $f_{\text{stop}} = 10^{-10}$  versus population size for (a) Rastrigin function (b) Griewank function for dimensions  $n = 2$  ('-○-'),  $n = 5$  ('-×-'),  $n = 10$  ('-□-'),  $n = 20$  ('-+-'),  $n = 40$  ('-◇-'), and  $n = 80$  ('-▽-').

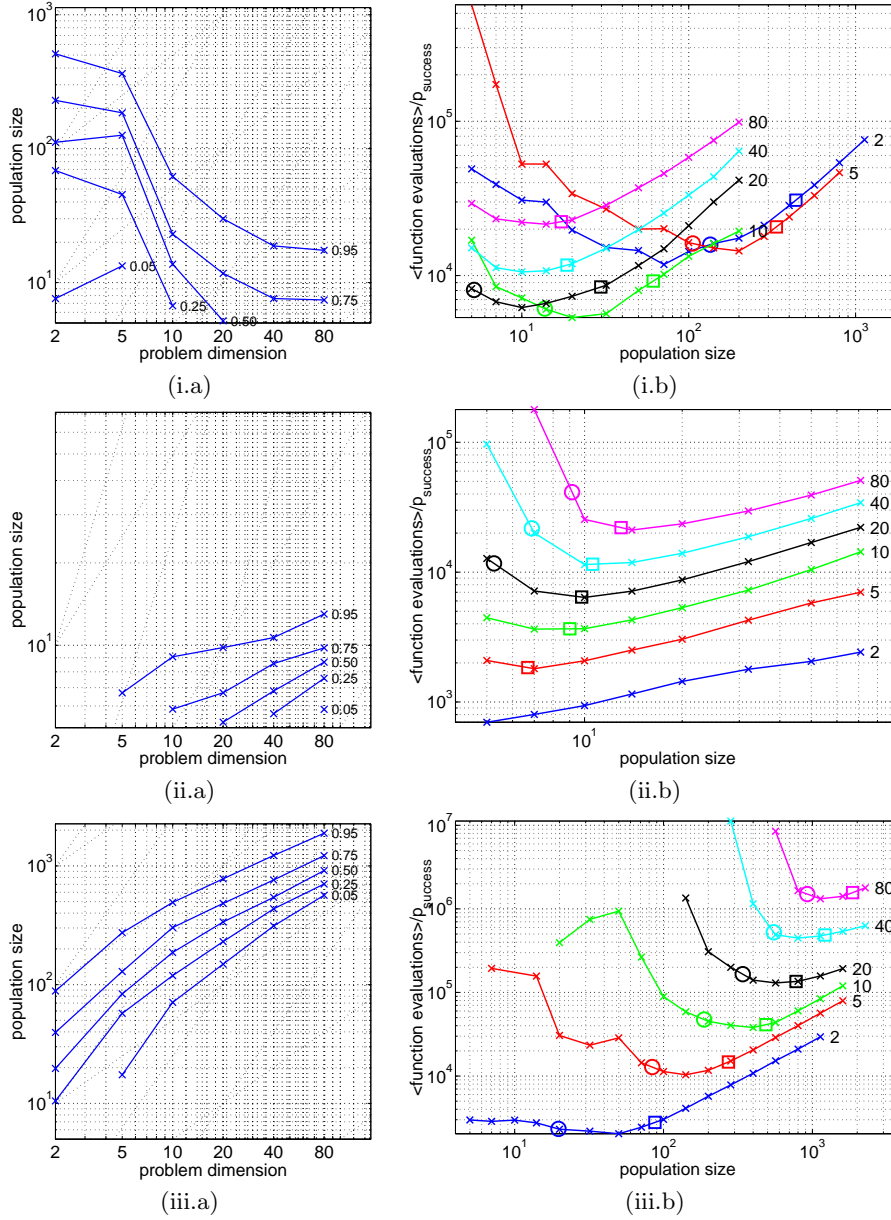
size is increased repeatedly in the sequence 5, 7, 10, 14, 20, 32, 50,  $\lceil 50\sqrt{2} \rceil$ , 100,  $\lceil 100\sqrt{2} \rceil$ , 200,  $\dots$ , and for each setting 100 runs are conducted. The starting point  $\mathbf{x}^{(0)}$  is sampled uniformly within the initialization intervals given in Table 1. The initial step size  $\sigma^{(0)}$  is set to half of the initialization interval. Too small initial step sizes have a considerable impact on the performance on multimodal functions. Each run is stopped and regarded as successful, when the function value is smaller than  $f_{\text{stop}} = 10^{-10}$ . Additionally, the run is stopped after  $10^7$  function evaluations, or when the condition number of the covariance matrix  $\mathbf{C}$  exceeds  $10^{14}$ , or by the option TolX, set to  $10^{-15}$  (for  $f_{\text{Schaffer}} 10^{-30}$ ).

## 4 Simulation Results

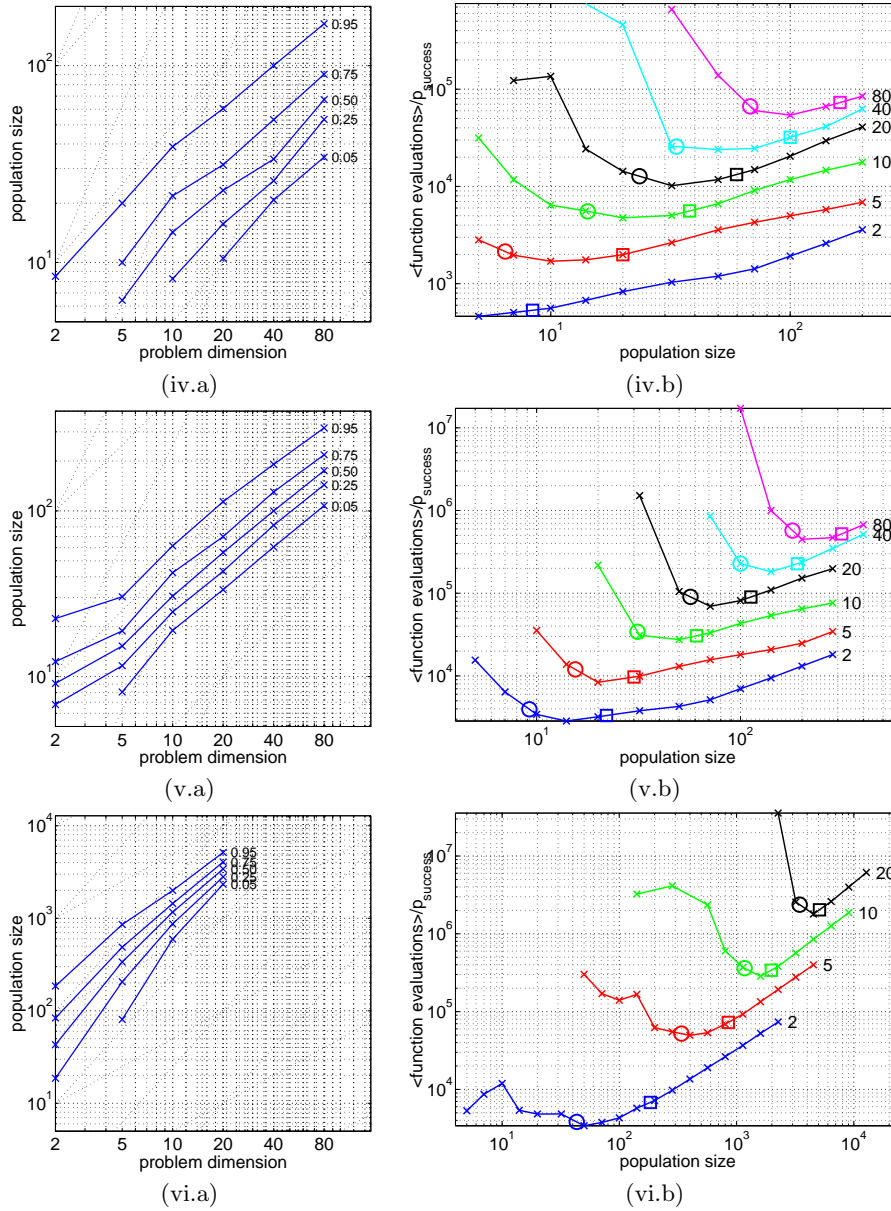
The success rate to reach  $f_{\text{stop}}$  depends strongly on the population size, see Fig. 1, where exemplary results are shown for (a)  $f_{\text{Rastrigin}}$  and (b)  $f_{\text{Griewank}}$ . The Rastrigin function (Fig. 1a) represents the typical picture. The graphs have a sigmoidal shape and larger dimensions require larger population sizes to approach 100% success rate. We observe two exceptions from this behavior. First, on  $f_{\text{RastSkew}}$  the success rates are low with any population size, and, except for very low dimensions,  $f_{\text{RastSkew}}$  is not solvable for the CMA-ES. The second exception is shown in Fig. 1b: on  $f_{\text{Griewank}}$  smaller dimensions require larger population sizes, but success rates of 100% can be achieved in all dimensions.

Figures 2a and 3a show the scaling of the (minimal) population size w.r.t. the problem dimension to achieve success rates of (at least) 5, 25, 50, 75, 95%. Graphs for  $f_{\text{RastScaled}}$  are almost identical with  $f_{\text{Rastrigin}}$  and therefore omitted. Except for  $f_{\text{RastSkew}}$  (not shown), larger success rates require larger population sizes. The figures are sorted from (i.a) to (vi.a) by increasing slopes that indicate the scaling. The steepest slope ( $f_{\text{Schwefel}}$ ) is slightly above linear. For all other functions the slope is sub-linear.

Figures 2b and 3b show performance versus population size. Performance is measured as mean number of function evaluations for successful runs, divided by

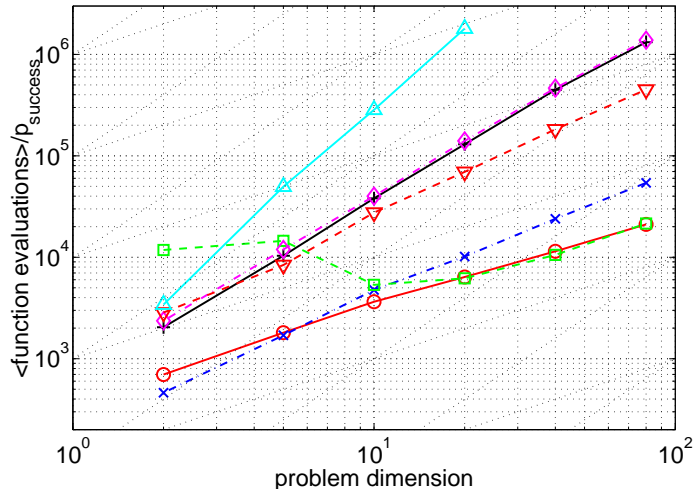


**Fig. 2.** (i) Griewank function, (ii) Ackley function, (iii) Rastrigin function. (a) population size to reach success rates of 0.05, 0.25, 0.5, 0.75, and 0.95 versus problem dimension  $n$ . The sloping grid indicates linear and quadratic dependency. (b) average number of function evaluations to reach  $f_{\text{stop}} = 10^{-10}$  divided by the success rate, versus population size for problem dimensions  $n = 2, 5, 10, 20, 40, 80$ . The symbols  $\circ$  and  $\square$  indicate success rates of 50% and 95%, respectively. Missing points on the left side of a graph indicate that no run (out of 100 runs) was successful.



**Fig. 3.** (iv) Bohachevsky function, (v) Schaffer function, (vi) Schwefel function. (a) population size to reach success rates of 0.05, 0.25, 0.5, 0.75, and 0.95 versus problem dimension  $n$ . The sloping grid indicates linear and quadratic dependency. (b) average number of function evaluations to reach  $f_{\text{stop}} = 10^{-10}$  divided by the success rate, versus population size for problem dimensions  $n = 2, 5, 10, 20, 40, 80$ . The symbols  $\circ$  and  $\square$  indicate success rates of 50% and 95%, respectively. Missing points on the left side of a graph indicate that no run (out of 100 runs) was successful.





**Fig. 4.** Mean number of function evaluations to reach  $f_{\text{stop}}$  versus problem dimension  $n$  for CMA-ES on Ackley ('—○—'), Bohachevsky ('---×---'), Griewank ('--□--'), Rastrigin ('—+—'), Scaled Rastrigin ('- - -◇ - - -'), Schaffer ('- - ▽ - -'), and Schwefel ('—△—') function.

the success rate. This performance measure assumes the same expected number of function evaluations for successful and for unsuccessful runs. The best performance is usually achieved for success rates between 50% and 95%. The impact of a smaller than optimal population size can be high. With increasing population size the performance decreases at most linearly.

Figure 4 shows the scaleup of the performance with optimal population size. The scaleup with  $n$ , put in order, appears to be at most linear for  $f_{\text{Ackley}}$  and  $f_{\text{Griewank}}$ , between linear and quadratic for  $f_{\text{Bohachevsky}}$ ,  $f_{\text{Schaffer}}$ ,  $f_{\text{Rastrigin}}$ , and  $f_{\text{RastScaled}}$ , and slightly below cubic for  $f_{\text{Schwefel}}$ .

Table 2 compares the performance of the CMA-ES with optimal population size to the performance of Differential Evolution (DE) [7], the Robust Evolution Strategy (RES) [8], and a Local Optima Smoothing (LOS) [9] restart BFGS algorithm. For each function the results with the best parameter tuning were taken from the respective publication and additional experiments were performed with DE and BFGS.<sup>7</sup> Only on the additively separable functions  $f_{\text{Rastrigin}}$  and  $f_{\text{Schwefel}}$ , DE outperforms the CMA-ES by a factor of five to 50. Otherwise, the CMA-ES outperforms DE by a factor of at least three, while DE even fails to find the global optimum on the non-separable  $f_{\text{Rastrigin}}(\mathbf{Ax})$ , where an orthogonal transformation  $\mathbf{A}$  of the search space is applied, and performs much worse

<sup>7</sup> In [9] the results for LOS are stated as numbers of necessary local searches to hit the global optimum. The average number of function evaluations for a local search using BFGS is simulated with MATLAB's `fminunc` with `MaxFunEvals` =  $500n$  and `TolX` = `TolFun` =  $10^{-3}$  (0.9 for  $f_{\text{Rastrigin}}$ ). A *random* restart strategy performs much worse than LOS.

**Table 2.** Average number of function evaluations to reach  $f_{\text{stop}}$  of CMA-ES versus DE, RES, and LOS on Griewank, Ackley, Rastrigin, and Schwefel function. If  $f_{\text{stop}}$  is not reached the *best function value/number of function evaluations* are stated. The initialization regions do not apply to CMA-ES for which the more difficult intervals from Table 1 are used. Results are taken from [7–9], except for the cases marked with \* obtained using the DE code available from <http://www.icsi.berkeley.edu/~storn/code.html>. The matrix  $\mathbf{A} = [\mathbf{o}_1, \dots, \mathbf{o}_n]^T$  implements an angle-preserving (i.e. orthogonal) linear transformation of  $\mathbf{x}$  (cf. [1]), chosen anew for each run.

Function	$f_{\text{stop}}$	init	$n$	CMA-ES	DE	RES	LOS
$f_{\text{Ackley}}(\mathbf{x})$	1e-3	[−30, 30] <sup>n</sup>	20	<b>2667</b>	.	.	6.0e4
			30	<b>3701</b>	12481	1.1e5	9.3e4
			100	<b>11900</b>	36801	.	.
$f_{\text{Griewank}}(\mathbf{x})$	1e-3	[−600, 600] <sup>n</sup>	20	<b>3111</b>	8691	.	.
			30	<b>4455</b>	11410 *	$8.5e-3/2e5$	.
			100	<b>12796</b>	31796	.	.
$f_{\text{Rastrigin}}(\mathbf{x})$	0.9	[−5.12, 5.12] <sup>n</sup> DE: [−600, 600] <sup>n</sup>	20	68586	<b>12971</b>	.	9.2e4
			30	147416	<b>20150</b> *	1.0e5	2.3e5
			100	1010989	<b>73620</b>	.	.
$f_{\text{Rastrigin}}(\mathbf{Ax})$	0.9	[−5.12, 5.12] <sup>n</sup>	30	<b>152000</b>	$171/1.25e6$ *	.	.
			100	<b>1011556</b>	$944/1.25e6$ *	.	.
$f_{\text{Schwefel}}(\mathbf{x})$	1e-3	[−500, 500] <sup>n</sup>	5	43810	<b>2567</b> *	.	7.4e4
			10	240899	<b>5522</b> *	.	5.6e5

on  $f_{\text{Schwefel}}(\mathbf{Ax})$  (not shown).<sup>8</sup> This supports our hypothesis that DE strongly exploits the separability of the function. The RES too exploits separability by sampling Cauchy distributions which strongly favor steps in coordinate directions. Even so, on the separable  $f_{\text{Rastrigin}}$  RES outperforms CMA-ES only by a factor of 1.5, while on  $f_{\text{Ackley}}$  and  $f_{\text{Griewank}}$  it performs worse by a factor of 30 or more. The LOS performs between a factor 1.5 (on  $f_{\text{Rastrigin}}$ ) and a factor of 25 (on  $f_{\text{Ackley}}$ ) worse than the CMA-ES.<sup>9</sup>

## 5 Summary and Conclusion

The CMA-ES with rank- $\mu$ -update is investigated on a suit of eight highly multimodal test functions for problem dimensions between 2 and 80. Tuning (that is increasing) the population size considerably improves the performance on six of the functions, compared to the performance with default population size.

<sup>8</sup> On  $f_{\text{Rastrigin}}$ , the parent number equals 20 in DE. Increasing the parent number improves the performance on  $f_{\text{Rastrigin}}(\mathbf{Ax})$ . However, even with 500 parents for  $n = 10$ , the minimum function value reached does not drop below 2 after  $10^7$  function evaluations. Choosing the recombination parameter  $\text{CR} = 1$ , DE becomes invariant against orthogonal transformations, but performs even worse on  $f_{\text{Rastrigin}}$ .

<sup>9</sup> Two parameters of LOS,  $r$  and  $K$ , see [9], are chosen to be optimal for each entry. In particular  $r$  has a considerable impact on the performance.

On seven of the eight functions, the CMA-ES can precisely locate the global optimum. If the local optima can be interpreted as perturbations of an underlying unimodal function, the CMA-ES with a large population size can “detect” the global topology. Then, the global optimum is located within  $300n$  and  $500n^2$  function evaluations. A strong asymmetry of the underlying function jeopardizes a successful detection and can lead to a failure (as on  $f_{\text{RastSkew}}$ ). The optimal population size usually scales sub-linearly with the problem dimension  $n$ , but significantly depends on the test function considered.

The results were compared with other global search strategies, stated to achieve superior results in earlier investigations. Surprisingly, the CMA-ES outperforms these global searchers, typically by a factor of three, with the following exception. Only if the function is *additively separable*, Differential Evolution strongly outperforms the CMA-ES. If the search space is rotated, the performance of the CMA-ES is unchanged, however Differential Evolution massively degrades in performance or even fails to locate the global optimum with a reasonable probability. For the CMA-ES the population size was tuned, while for the compared algorithms up to three parameters had to be tuned to the given objective function. In our opinion, tuning the population size in the CMA-ES is comparatively unproblematic. The results suggest that a CMA-ES restart strategy with a successively increased population size (by a factor of three, initialized with the default population size) constitutes a highly competitive, quasi parameter free global optimization algorithm for non-separable objective functions.

## References

1. Hansen, N., Ostermeier, A.: Completely derandomized self-adaptation in evolution strategies. *Evol. Comput.* **9** (2001) 159–195
2. Hansen, N., Ostermeier, A.: Adapting arbitrary normal mutation distributions in evolution strategies: The covariance matrix adaptation. In: Proceedings of the 1996 IEEE Conference on Evolutionary Computation (ICEC '96). (1996) 312–317
3. Müller, S.D., Hansen, N., Koumoutsakos, P.: Increasing the serial and the parallel performance of the CMA-evolution strategy with large populations. In: Parallel Problem Solving from Nature (PPSN). (2002)
4. Hansen, N., Müller, S.D., Koumoutsakos, P.: Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES). *Evol. Comput.* **11** (2003) 1–18
5. Kern, S., Müller, S.D., Hansen, N., Büche, D., Ocenasek, J., Koumoutsakos, P.: Learning probability distributions in continuous evolutionary algorithms – a comparative review. *Natural Computing* **3** (2004) 77–112
6. Büche, D.: Personal communication (2003)
7. Storn, R., Price, K.: Differential evolution: A simple and efficient heuristic for global optimization over continuous spaces. *J. Glob. Opt.* **11** (1997) 341–359
8. Ohkura, K., Matsumura, Y., Ueda, K.: Robust evolution strategies. *Lect. Notes Comput. Sc.* **1585** (1999) 10–17
9. Addis, B., Locatelli, M., Schoen, F.: Local optima smoothing for global optimization. Technical report dsi 5–2003, Dipartimento di Sistemi e Informatica, Università degli Studi di Firenze (2003)