

Mirrored Variants of the (1,4)-CMA-ES Compared on the Noiseless BBOB-2010 Testbed

[Black-Box Optimization Benchmarking Workshop]

Anne Auger, Dimo Brockhoff, and Nikolaus Hansen
Projet TAO, INRIA Saclay—Ile-de-France
LRI, Bât 490, Univ. Paris-Sud
91405 Orsay Cedex, France
firstname.lastname@inria.fr

ABSTRACT

Derandomization by means of mirrored samples has been recently introduced to enhance the performances of $(1, \lambda)$ -Evolution-Strategies (ESs) with the aim of designing fast and robust stochastic local search algorithms. This paper compares on the BBOB-2010 noiseless benchmark testbed two variants of the (1,4)-CMA-ES where the mirrored samples are used. Independent restarts are conducted up to a total budget of $10^4 D$ function evaluations, where D is the dimension of the search space.

The results show that the improved variants are significantly faster than the baseline (1,4)-CMA-ES on 4 functions in 20D (respectively 7 when using sequential selection in addition) by a factor of up to 3 (on the attractive sector function). In no case, the (1,4)-CMA-ES is significantly faster on any tested target function value in 5D and 20D. Moreover, the algorithm employing both mirroring and sequential selection is significantly better than the algorithm without sequentialism on five functions in 20D with expected running times that are about 20% smaller.

Categories and Subject Descriptors

G.1.6 [Numerical Analysis]: Optimization—*global optimization, unconstrained optimization*; F.2.1 [Analysis of Algorithms and Problem Complexity]: Numerical Algorithms and Problems

General Terms

Algorithms

1. INTRODUCTION

Evolution Strategies (ESs) are robust stochastic search algorithms for numerical optimization where the function to be minimized, f , maps the continuous search space \mathbb{R}^D

into \mathbb{R} where D is the dimension of the search space. Recently, a new derandomization technique replacing the independent sampling of new solutions by mirrored sampling has been introduced to enhance the performances of ESs [1]. In this paper, we assess quantitatively the improvements that can be brought when using mirrored samples instead of independent ones. To do so, we compare on the BBOB-2010 noiseless testbed the (1,4)-Covariance-Matrix-Adaptation Evolution-Strategy (CMA-ES) with two variants: first the $(1,4_m)$ -CMA-ES where mirrored samples are used, and second the $(1,4_m^s)$ -CMA-ES that in addition to the mirrored samples uses the concept of sequential selection [1]. Both variants are described in Sec. 2.

2. THE ALGORITHMS TESTED

The three algorithms tested are variants of the well-known CMA-ES [10, 9, 8] where at each iteration n , λ new solutions are generated by sampling *independently* λ random vectors $(\mathcal{N}_i(\mathbf{0}, \mathbf{C}_n))_{1 \leq i \leq \lambda}$ following a multivariate normal distribution with mean vector $\mathbf{0}$ and covariance matrix \mathbf{C}_n . The vectors are added to the current solution \mathbf{X}_n to create the λ new solutions or offspring $\mathbf{X}_n^i = \mathbf{X}_n + \sigma_n \mathcal{N}_i(\mathbf{0}, \mathbf{C}_n)$ where σ_n is a strictly positive parameter called step-size. In the simple (1,4)-CMA-ES, the number of offspring λ equals 4 and \mathbf{X}_{n+1} is set to the best solution among $\mathbf{X}_n^1, \dots, \mathbf{X}_n^4$, i.e., $\mathbf{X}_{n+1} = \operatorname{argmin}\{f(\mathbf{X}_n^1), \dots, f(\mathbf{X}_n^4)\}$.

In the mirrored variant, denoted $(1,4_m)$ -CMA-ES, the second and fourth offspring are replaced by the offspring symmetric to the first and third with respect to \mathbf{X}_n , namely $\mathbf{X}_n^2 = \mathbf{X}_n - \sigma_n \mathcal{N}_1(\mathbf{0}, \mathbf{C}_n)$ and $\mathbf{X}_n^4 = \mathbf{X}_n - \sigma_n \mathcal{N}_3(\mathbf{0}, \mathbf{C}_n)$, where $\sigma_n \mathcal{N}_1(\mathbf{0}, \mathbf{C}_n)$ and $\sigma_n \mathcal{N}_3(\mathbf{0}, \mathbf{C}_n)$ are the random vectors added to \mathbf{X}_n to get \mathbf{X}_n^1 and \mathbf{X}_n^3 . The first/second and third/fourth offspring are, thus, negatively correlated. The update of \mathbf{X}_{n+1} is then identical to the (1,4)-CMA-ES, namely $\mathbf{X}_{n+1} = \operatorname{argmin}\{f(\mathbf{X}_n^1), \dots, f(\mathbf{X}_n^4)\}$.

In the $(1,4_m^s)$ -CMA-ES, sequential selection is implemented. The four offspring solutions are generated with mirrored sampling. Evaluations are carried out in a sequential manner: after evaluating solution \mathbf{X}_n^i , it is compared to \mathbf{X}_n and if $f(\mathbf{X}_n^i) \leq f(\mathbf{X}_n)$, the sequence of evaluations is stopped and $\mathbf{X}_{n+1} = \mathbf{X}_n^i$. In case all four offspring are worse than \mathbf{X}_n , $\mathbf{X}_{n+1} = \operatorname{argmin}\{f(\mathbf{X}_n^1), \dots, f(\mathbf{X}_n^4)\}$ according to the comma selection. In sequential selection, the number of offspring evaluated is a random variable by itself ranging from 1 to $\lambda = 4$ —allowing to reduce the number of offspring adaptively as long as improvements are easy to achieve [1].

Covariance matrix and step-size are updated using the selected steps [9, 1].

Independent Restarts.

Similar to [3], we independently restarted all algorithms as long as function evaluations were left, where maximally $10^4 \cdot D$ function evaluations have been used.

Parameter setting.

We used the default parameter and termination settings (cf. [1, 5, 8]) found in the source code on the WWW¹ with two exceptions. We rectified the learning rate of the rank-one update of the covariance matrix for small values of λ , setting $c_1 = \min(2, \lambda/3)/((D+1.3)^2 + \mu_{\text{eff}})$. The original value was not designed to work for $\lambda < 5$. We modified the damping parameter for the step-size to $d_\sigma = 0.3 + 2\mu_{\text{eff}}/\lambda + c_\sigma$. The setting was found by performing experiments on the sphere function, f_1 : d_σ was set as large as possible while still showing close to optimal performance, but, at least as large such that decreasing it by a factor of two did not lead to unacceptable performance. For $\mu_{\text{eff}}/\lambda = 0.35$ and $\mu_{\text{eff}} \leq D + 2$ the former setting of d_σ is recovered. For a smaller ratio of μ_{eff}/λ or for $\mu_{\text{eff}} > D + 2$, the new setting allows larger (i.e. faster) changes of σ . Here, $\mu_{\text{eff}} = 1$. For $\lambda \geq 3$, the new setting might be harmful in a noisy or too rugged landscape. Finally, the step-size multiplier was clamped from above at $\exp(1)$, while we do not believe this had any effect in the presented experiments. Each initial solution \mathbf{X}_0 was uniformly sampled in $[-4, 4]^D$ and the step-size σ_0 was initialized to 2. The source code used for the experiments is available at².

As the same parameter setting has been used for all test functions, the crafting effort CrE of all algorithms is 0.

3. CPU TIMING EXPERIMENTS

For the timing experiment, all three algorithms were run on f_8 with a maximum of $10^4 D$ function evaluations and restarted until at least 30 seconds have passed (according to Figure 2 in [6]). The experiments have been conducted with an 8 core Intel Xeon E5520 machine with 2.27 GHz under Ubuntu 9.1 linux and Matlab R2008a. The time per function evaluation was 3.3; 3.3; 3.0; 3.1; 3.4; 4.0 times 10^{-4} seconds for (1,4)-CMA-ES, 3.1; 3.0; 3.0; 3.2; 3.4; 4.0 times 10^{-4} seconds for (1,4_m)-CMA-ES, and 7.1; 7.3; 7.7; 8.1; 7.1; 8.1 times 10^{-4} seconds for (1,4_m^s)-CMA-ES in dimensions 2; 3; 5; 10; 20; 40 respectively. Note that MATLAB distributes the computations over all 8 cores only for 20D and 40D.

4. RESULTS

In this section, experiments according to [6] on the benchmark functions given in [4, 7] are presented. The **expected running time (ERT)**, used in the figures and table, depends on a given target function value, $f_t = f_{\text{opt}} + \Delta f$, and is computed over all relevant trials as the number of function evaluations executed during each trial while the best function value did not reach f_t , summed over all trials and divided by the number of trials that actually reached f_t [6, 11]. **Statistical significance** is tested with the rank-sum test for a given target Δf_t using, for each trial, either the number of needed function evaluations to reach Δf_t (inverted

and multiplied by -1), or, if the target was not reached, the best Δf -value achieved, measured only up to the smallest number of overall function evaluations for any unsuccessful trial under consideration.

4.1 Comparing (1,4_m)- With (1,4)-CMA-ES

The (1,4_m)-CMA-ES is compared with the (1,4)-CMA-ES in Figures 1 and 2 and in Table 1. Mirroring within the (1,4_m)-CMA-ES seems to have an important positive impact compared to the baseline (1,4)-CMA-ES. In 20D, we observe a slight worsening only on the Gallagher functions f_{21} and f_{22} which are not statistically significant. On the other hand, the (1,4_m)-CMA-ES outperforms the (1,4)-CMA-ES on 10 functions of which 4 show statistically significant differences (in 20D and for a target of 10^{-7}): on the sphere function (f_1), the expected running time is 15% lower, on the separable ellipsoid (f_2) it is 18% lower, on the non-separable ellipsoid (f_{10}) 19% lower and on the attractive sector function (f_6) the expected running times even differ by a factor of about 3. These differences are less significant in 5D (Table 1), but Fig. 1 shows similar improvements also for 10D and smaller dimensions.

4.2 Comparing (1,4_m^s)- With (1,4_m)-CMA-ES

Results comparing the (1,4_m)-CMA-ES and the (1,4_m^s)-CMA-ES are presented in Figures 3 and 4 and in Table 2.

The results show that the (1,4_m^s)-CMA-ES is statistically significantly better than the (1,4_m)-CMA-ES on five functions in 20D with expected running times that are about 26% smaller on f_1 and f_2 , 25% smaller on f_{10} , 20% smaller on f_{11} , and about 20% smaller on f_{14} than the ones for the (1,4_m)-CMA-ES. Contrary, there are only two functions where the (1,4_m^s)-CMA-ES is worse than the (1,4_m)-CMA-ES (13% on f_6 and about 70% on f_{22}) but these differences are not statistically significant. To conclude, the (1,4_m^s)-CMA-ES should be clearly preferred over the (1,4_m)-CMA-ES on the BBOB-2010 testbed and, according to Sec. 4.1, also over the (1,4)-CMA-ES. Worth to mention is also the fact, that the (1,4_m^s)-CMA-ES shows better performance than the best algorithm from the BBOB-2009 benchmarking on the attractive sector function (f_6) in 5D by about 30% and a slightly better performance on the ellipsoid (f_{10}) in 20D, see [2].

4.3 Comparing (1,4_m^s)- With (1,4)-CMA-ES

Regarding this comparison, we refrain from showing the plots and tables due to space limitations. Not surprisingly when considering Sec. 4.2, the results show a statistically significant improvement for the (1,4_m^s)-CMA-ES over the (1,4)-CMA-ES on 7 functions in 20D: f_1 (37% faster), f_2 (39% faster), f_5 (faster by a factor of about 2), f_6 (about 3 times faster), f_{10} (39% faster), f_{11} (25% faster), and f_{14} (about 30% faster). Only on the Gallagher function with 21 peaks (f_{22}), an increase of the expected running time by a factor of about 2 can be observed which is, however, not statistically significant due to the low number of successful runs.

5. CONCLUSIONS

The idea behind derandomization by means of mirroring introduced in [1] is to use only one random sample from a multivariate normal distribution to create two (negatively correlated or *mirrored*) offspring. Thereby, one offspring is generated by adding a random sample to the parent solution

¹cmaes.m, version 3.41.beta, from http://www.lri.fr/~hansen/cmaes_inmatlab.html

²<http://coco.gforge.inria.fr/doku.php?id=bbob-2010-results>

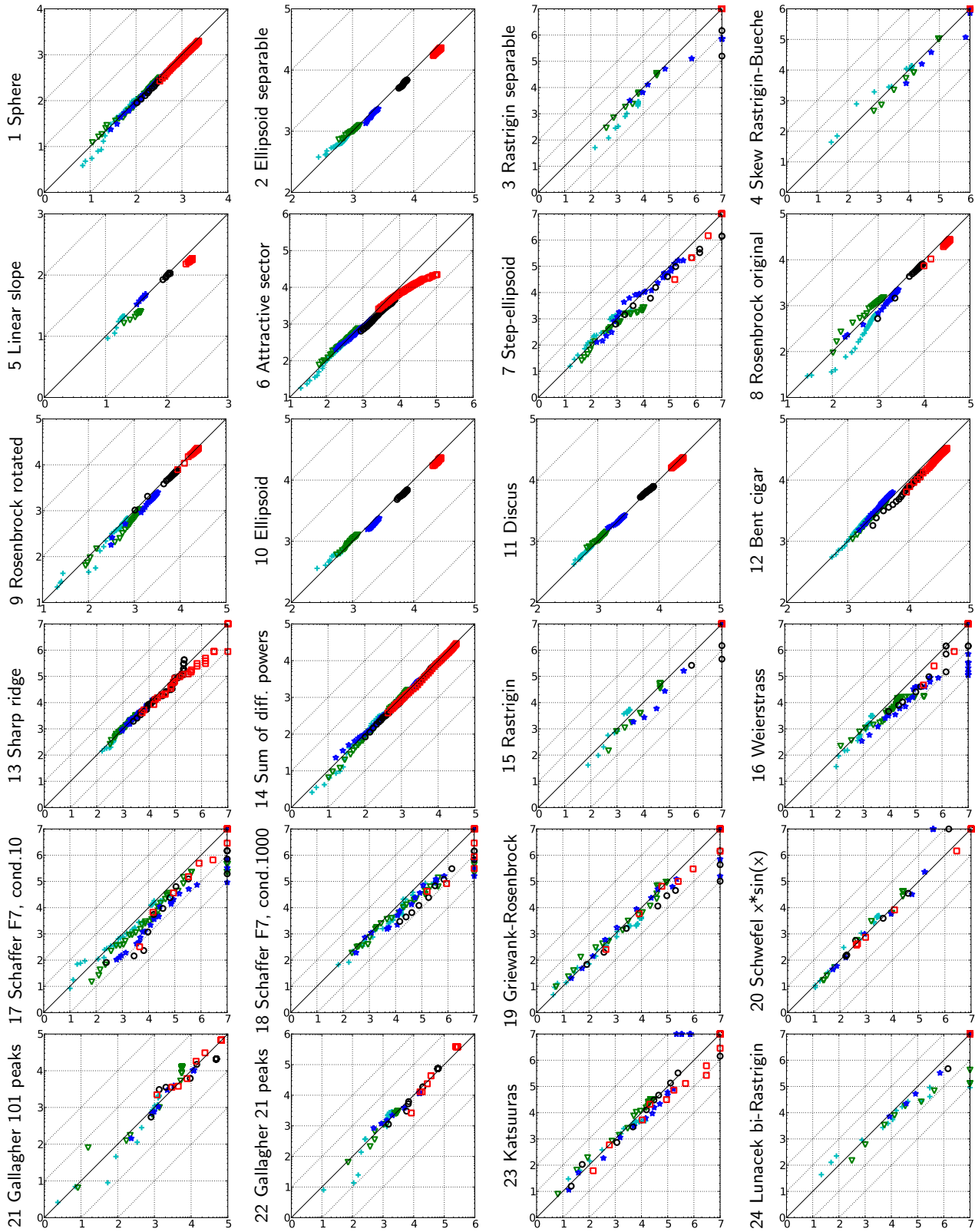


Figure 1: Expected running time (ERT in log10 of number of function evaluations) of $(1,4_m)$ -CMA-ES versus $(1,4)$ -CMA-ES for 46 target values $\Delta f \in [10^{-8}, 10]$ in each dimension for functions f_1 - f_{24} . Markers on the upper or right edge indicate that the target value was never reached by $(1,4_m)$ -CMA-ES or $(1,4)$ -CMA-ES respectively. Markers represent dimension: 2:+, 3:∇, 5:*, 10:o, 20:□.

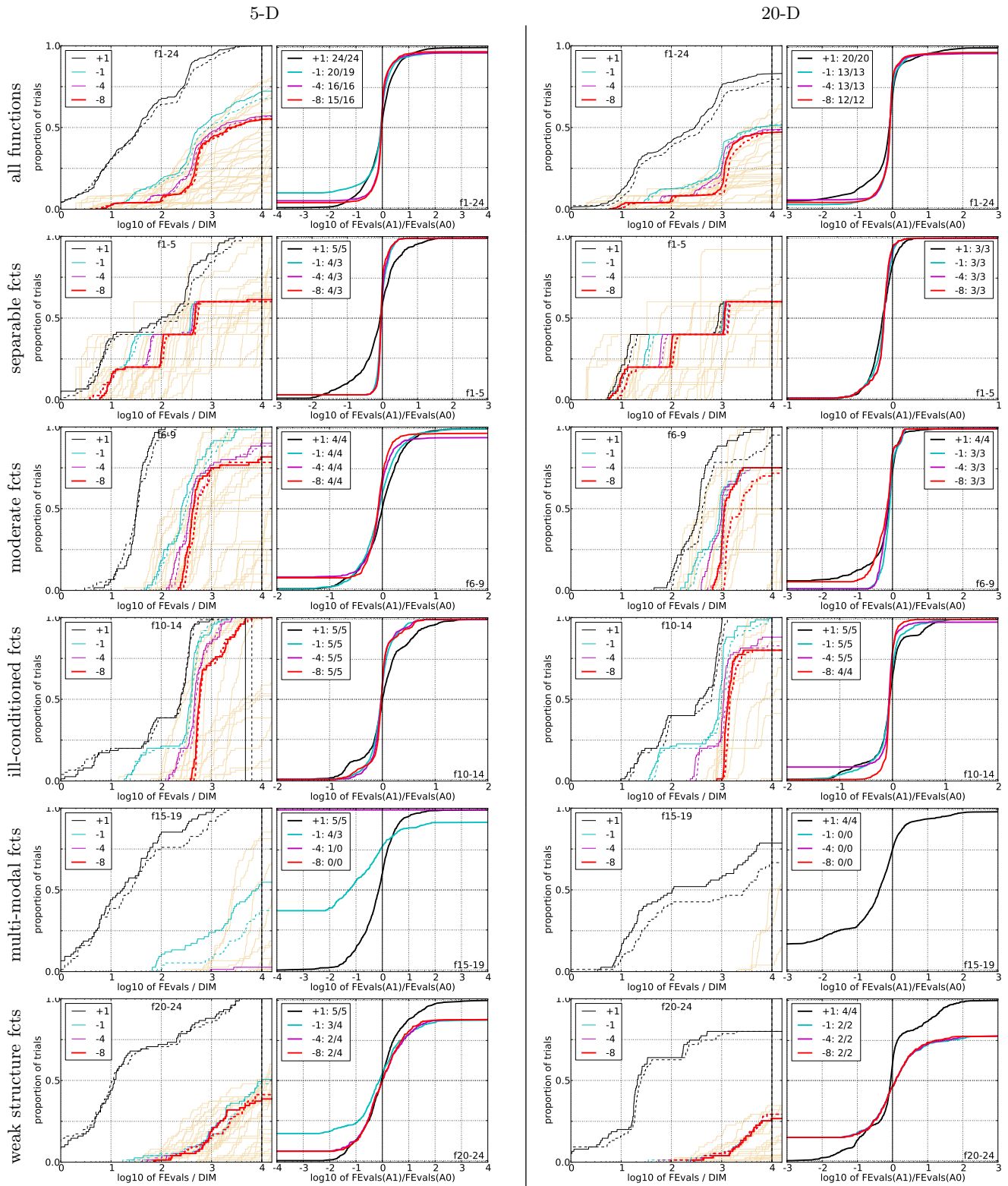


Figure 2: Empirical cumulative distributions (ECDF) of run lengths and speed-up ratios in 5-D (left) and 20-D (right). Left sub-columns: ECDF of the number of function evaluations divided by dimension D (FEvals/ D) to reach a target value $f_{\text{opt}} + \Delta f$ with $\Delta f = 10^k$, where $k \in \{1, -1, -4, -8\}$ is given by the first value in the legend, for $(1,4_m)$ -CMA-ES (solid) and $(1,4)$ -CMA-ES (dashed). Light beige lines show the ECDF of FEvals for target value $\Delta f = 10^{-8}$ of algorithms benchmarked during BBOB-2009. Right sub-columns: ECDF of FEval ratios of $(1,4_m)$ -CMA-ES divided by $(1,4)$ -CMA-ES, all trial pairs for each function. Pairs where both trials failed are disregarded, pairs where one trial failed are visible in the limits being > 0 or < 1 . The legends indicate the number of functions that were solved in at least one trial ($(1,4_m)$ -CMA-ES first).

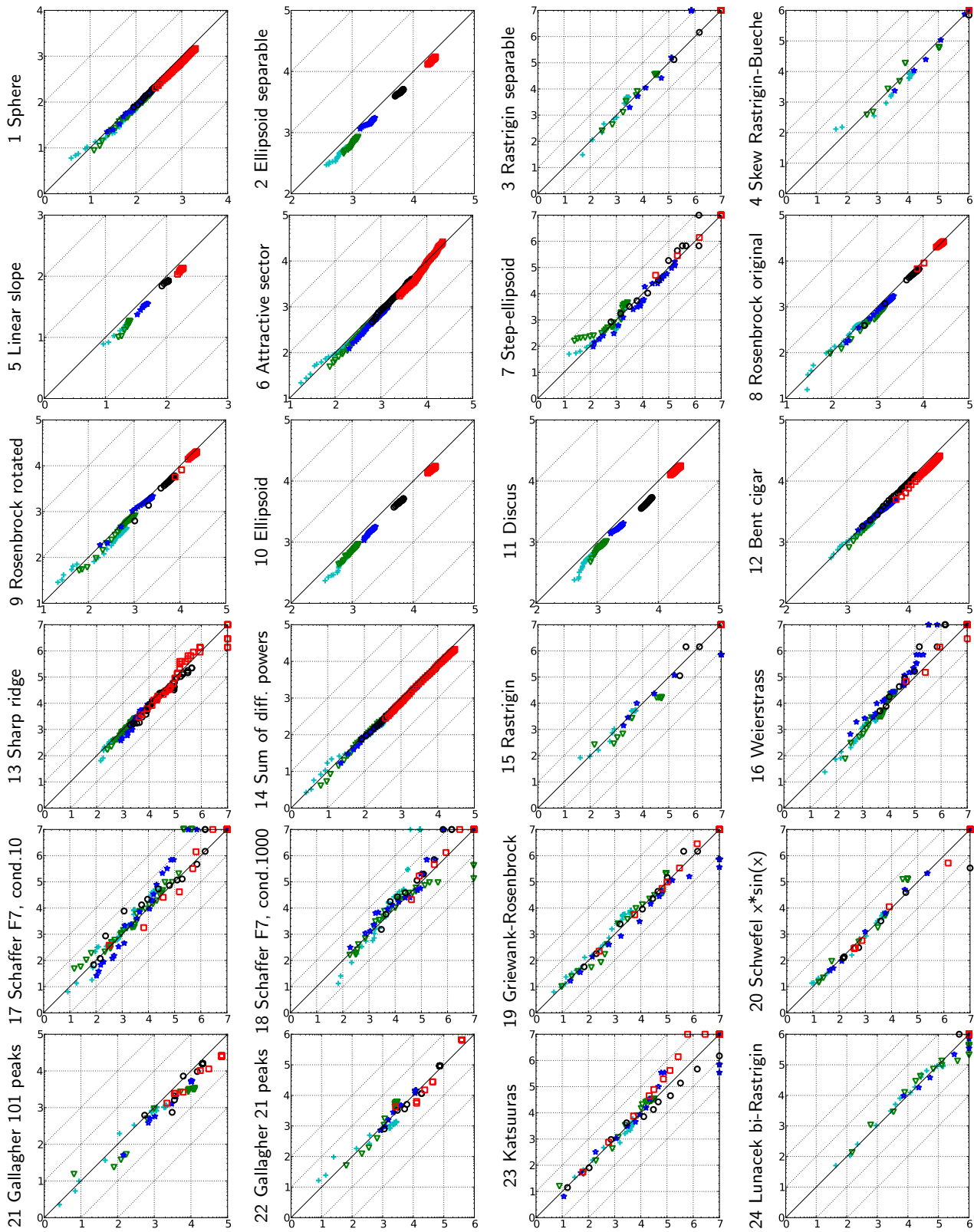


Figure 3: Expected running time (ERT in log10 of number of function evaluations) of $(1,4_m^s)$ -CMA-ES versus $(1,4_m)$ -CMA-ES for 46 target values $\Delta f \in [10^{-8}, 10]$ in each dimension for functions f_1 – f_{24} . Markers on the upper or right edge indicate that the target value was never reached by $(1,4_m^s)$ -CMA-ES or $(1,4_m)$ -CMA-ES respectively. Markers represent dimension: 2:+, 3:∇, 5:*, 10:○, 20:□.

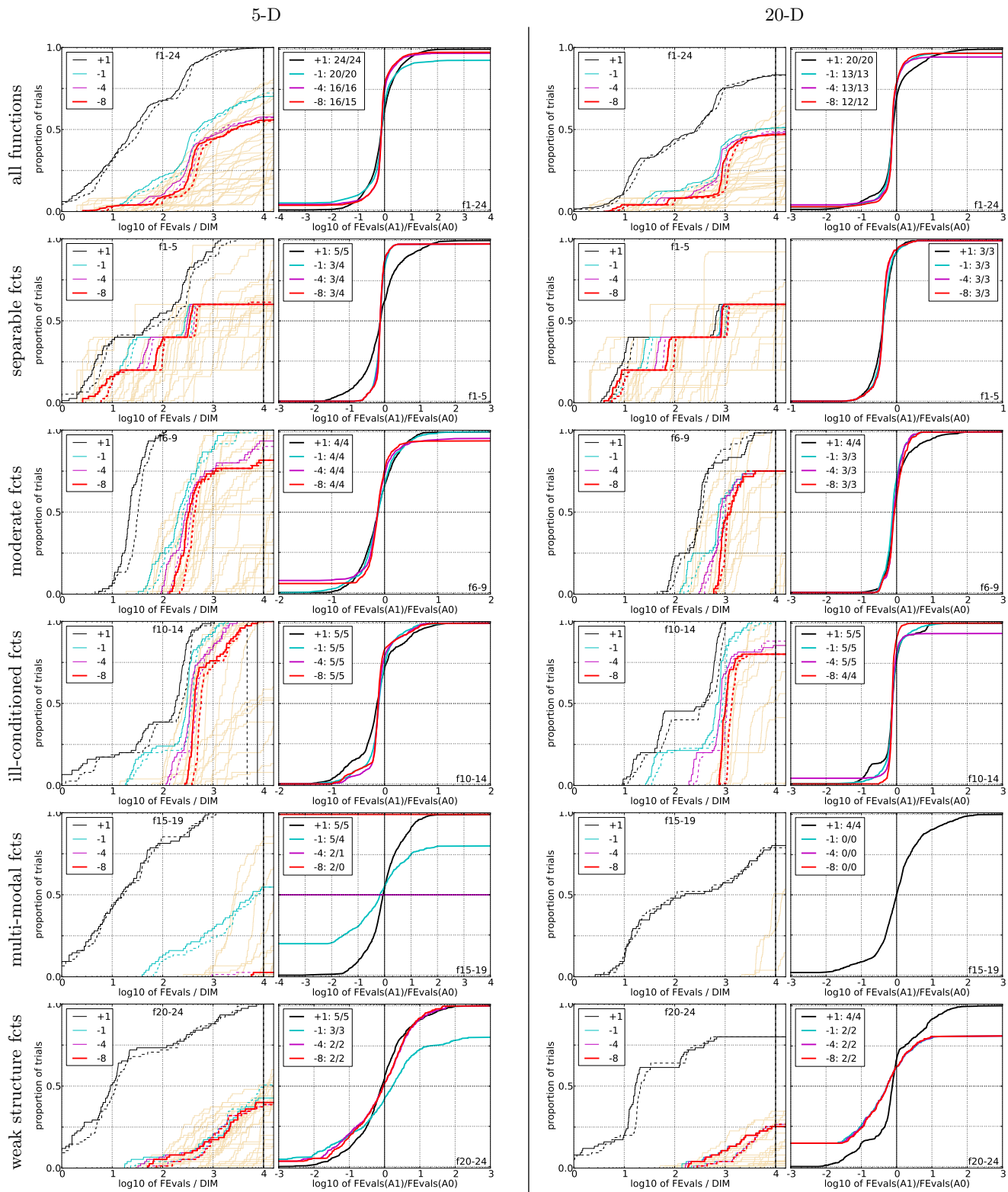


Figure 4: Empirical cumulative distributions of run lengths of $(1,4_m^s)$ -CMA-ES (solid) and $(1,4_m)$ -CMA-ES (dashed) and speed-up ratios of the former divided by the latter in 5-D (left) and 20-D (right) as in Fig. 2.

[2] A. Auger, S. Finck, N. Hansen, and R. Ros. BBOB 2009: Comparison tables of all algorithms on all noiseless functions. Technical Report RT-0383, INRIA, April 2010.
 [3] A. Auger and N. Hansen. Performance evaluation of an

advanced local search evolutionary algorithm. In *Proceedings of the IEEE Congress on Evolutionary Computation (CEC 2005)*, pages 1777–1784, 2005.
 [4] S. Finck, N. Hansen, R. Ros, and A. Auger. Real-parameter

