



**HAL**  
open science

## Benchmarking MO-CMA-ES and COMO-CMA-ES on the Bi-objective bbob-biobj Testbed

Paul Dufossé, Cheikh Touré

► **To cite this version:**

Paul Dufossé, Cheikh Touré. Benchmarking MO-CMA-ES and COMO-CMA-ES on the Bi-objective bbob-biobj Testbed. GECCO 2019 Companion - The Genetic and Evolutionary Computation Conference, Jul 2019, Prague, Czech Republic. 10.1145/3319619.3326892 . hal-02161252

**HAL Id: hal-02161252**

**<https://inria.hal.science/hal-02161252v1>**

Submitted on 20 Jun 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Benchmarking MO-CMA-ES and COMO-CMA-ES on the Bi-objective bbob-biobj Testbed

Paul Dufossé

Thales DMS France, Inria, CMAP, Ecole Polytechnique, IP  
Paris  
paul.dufosse@inria.fr

Cheikh Touré

Inria, CMAP, Ecole Polytechnique, IP Paris  
cheikh.toure@polytechnique.edu

## ABSTRACT

In this paper, we propose a comparative benchmark of MO-CMA-ES, COMO-CMA-ES (recently introduced in [12]) and NSGA-II, using the COCO framework for performance assessment and the Bi-objective test suite bbob-biobj. For a fixed number of points  $p$ , COMO-CMA-ES approximates an optimal  $p$ -distribution of the Hypervolume Indicator. While not designed to perform on archive-based assessment, *i.e.* with respect to all points evaluated so far by the algorithm, COMO-CMA-ES behaves well on the COCO platform. The experiments are done in a true Black-Blox spirit by using a minimal setting relative to the information shared by the 55 problems of the bbob-biobj Testbed.

## CCS CONCEPTS

•Computing methodologies → Continuous space search;

## KEYWORDS

Benchmarking, Black-box optimization, Bi-objective optimization

### ACM Reference format:

Paul Dufossé and Cheikh Touré. 2019. Benchmarking MO-CMA-ES and COMO-CMA-ES on the Bi-objective bbob-biobj Testbed. In *Proceedings of Genetic and Evolutionary Computation Conference Companion, Prague, Czech Republic, July 13–17, 2019 (GECCO '19 Companion)*, 8 pages. DOI: 10.1145/3319619.3326892

## 1 INTRODUCTION

As the COCO platform assesses bi-objective algorithms by computing the covered Hypervolume of evaluated points in a black-box optimization framework, it is natural that well-performing benchmarked algorithms on COCO are designed to tackle these two issues (covering the front and black-box paradigm). Hence [9] uses a hybrid algorithm and parameter tuning to handle the black-box difficulty, and [11] develops an unbounded population size variant of MO-CMA-ES [8, 14] for the Pareto front covering issue.

MO-CMA-ES and COMO-CMA-ES, however, are non-hybrid algorithms with a population size  $p$  (number of kernels for COMO-CMA-ES) fixed *a priori*.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

GECCO '19 Companion, Prague, Czech Republic

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM. 978-1-4503-6748-6/19/07...\$15.00

DOI: 10.1145/3319619.3326892

## 2 ALGORITHMS DESCRIPTION

NSGA-II is a popular choice among Evolutionary Multi-Objective (EMO) algorithms [3], based on non dominated sorting of candidates and crowding distance comparisons. It has already been benchmarked on the bbob-biobj Testbed [1]. The archive data, *i.e.* the data for all points evaluated so far, is shown under the name NSGA-II along with the best2016 reference, as a baseline for the reader.

The elitist Multiobjective CMA-ES, MO-CMA-ES [8], is based on a two-way ranking: the Pareto ranking and the hypervolume contribution among individuals with the same Pareto rank. Each parent is not only a point in the search space but a 5-tuple of data representing a (1+1)-CMA-ES, which is a Single Objective (SO) optimizer defined in [8]. We use here the improved step-size adaptation designed in [14]. The algorithm also uses greedy mating, which means that parents are selected only among non-dominated solutions. We choose the default settings for the  $\lambda_{MO}$ -(1+1)-CMA-ES presented in [8]. In this paper, the term population size refers to the MO population size, namely the parameter  $\lambda_{MO}$ .

A recent EMO algorithm using the non-elitist CMA-ES [4] is the Comma Multiobjective CMA-ES (COMO-CMA-ES) defined in [12], and is designed to approximate the optimal  $p$ -distribution of the Hypervolume indicator. It is the instantiation of a wider framework called Sofomore, on the non-elitist CMA-ES [12]. We also take the default setup for the standard CMA-ES presented in [4]. The *Single-objective Optimization FOR Optimizing Multiobjective Optimization pRobLEms* framework (Sofomore) finds  $p$  points approximating the optimal  $p$ -distribution of the Hypervolume indicator, by iteratively maximizing an extended notion of the hypervolume contribution (hypervolume improvement) called Uncrowded Hypervolume Improvement UHVI [12]. As the two-way ranking in [8], UHVI unflattens the hypervolume improvement's level sets in dominated regions, but the novelty is inherent to the fact that this unflattening operation still preserves diversity among the solutions, by ensuring them to be uncrowded.

## 3 IMPLEMENTATION AND EXPERIMENTAL PROCEDURE

The code for MO-CMA-ES was written in Matlab, 2014, using routines from the Shark library for hypervolume computation. COMO-CMA-ES is implemented in Python and uses NumPy. Note that the hypervolume is computed in pure Python and does not call any NumPy method. Implementation details for NSGA-II can be found in [1].

For MO-CMA-ES, the algorithm starts with random starting point uniformly sampled in  $I = [-5, 5]^n$  with an initial stepsize  $\sigma_0$

Algorithm	$p$	2-D	3-D	5-D	10-D	20-D
MO-	10	2.8	3.3	3.0	3.6	3.3
CMA-ES	32	2.7	2.9	2.8	2.9	2.9
	100	3.0	3.2	3.5	3.4	3.0
COMO-	10	5.9	6.0	5.0	4.7	4.8
CMA-ES	32	5.8	5.2	5.8	4.1	4.9
	100	12	8.9	7.7	7.5	4.3

**Table 1: CPU times per function evaluation of MO-CMA-ES and COMO-CMA-ES (in  $10^{-4}$  seconds) for varying dimensions.  $p$  is the population size (or number of kernels).**

equals to 1. From this starting point a population of  $p$  (1+1)-CMA-ES is evolved, with  $p$  the population size of the algorithm. Any (1+1)-CMA-ES can send a warning, meaning that it has somehow terminated. When any warning is met, the MO algorithm stops. Then a restart is performed with a random starting point in  $I$ , with the same settings. The population size is a core parameter fixed *a priori* for each run (including restarts), which is why MO-CMA-ES is benchmarked with different population sizes, namely 10, 32 and 100. The reference point used for computing the contributing hypervolume is iteratively chosen as the nadir point among the current population, and adding 1 to each coordinate.

COMO-CMA-ES also starts with a random point sampled in  $I$ , with an initial stepsize  $\sigma_0 = \sqrt{n}$ . No restart is performed and no parameter tuning is done on purpose to observe the behaviour of the algorithm with default settings. In the same spirit, the only stopping criteria is given by the allocated budget of function evaluations. As COMO-CMA-ES is designed to approximate the optimal  $p$ -distribution of the Hypervolume indicator, it is then necessary to fix one and for all a reference point for each run; which is done by setting an attribute of each COCO problem called `largest_fvalues_of_interest` [7] as reference point.

We benchmark MO-CMA-ES (with population sizes equal to 10, 32 and 100), COMO-CMA-ES (with 3, 10, 32, 100, 316 and 1000 kernels), and NSGA-II. The latter is run with a population size of 100 as in [1]. COCO is run with the common settings for the bbob-biobj test suite; 10 instances with dimension  $n \in \{2, 3, 5, 10, 20\}$ . The allocated budget is  $10^4 n$  function evaluations for each MO-CMA-ES, and  $10^5 n$  for NSGA-II and each COMO-CMA-ES.

## 4 CPU TIMING

In order to evaluate the duration of each algorithm, we have reported here the timing results when running the full budget experiment. Both codes were run on a linux machine with 64 cores, Intel®Xeon®E7 to E3 v4 processors. We are interested in CPU time per function evaluation for varying dimensions and population sizes. To account for Matlab internal parallelization we use the function `cpitime` and not real timing. The results can be found in Table 1.

We are benchmarking algorithms written in different languages, therefore comparisons should be made carefully. COMO-CMA-ES computation per function evaluation takes 1.5 to 4 times longer than MO-CMA-ES, and this ratio is less important when the dimension increases.

## 5 RESULTS

Results from experiments according to [7], [5] and [2] on the benchmark functions given in [13] for the three algorithms called MO-\*, NSGA-II or COMO-\* (with \* being the fixed population size) are presented in Figures 1, 2, 3 and 4. In Tables 2 and 3 we display results for each algorithm only with a population size of 100.

The **average runtime (aRT)** used in the tables depends on a given quality indicator value,  $I_{\text{target}} = I_{\text{ref}} + \Delta I_{\text{HV}}^{\text{COCO}}$ , and is computed over all relevant trials as the number of function evaluations executed during each trial while the best indicator value did not reach  $I_{\text{target}}$ , summed over all trials and divided by the number of trials that actually reached  $I_{\text{target}}$  [7]. **Statistical significance** is tested with the rank-sum test for a given target  $I_{\text{target}}$  using, for each trial, either the number of needed function evaluations to reach  $I_{\text{target}}$  (inverted and multiplied by  $-1$ ), or, if the target was not reached, the best  $\Delta I_{\text{HV}}^{\text{COCO}}$ -value achieved, measured only up to the smallest number of overall function evaluations for any unsuccessful trial under consideration.

The experiments were performed with COCO [6], version 2.2, the plots were produced with version 2.3.1.

For either COMO-CMA-ES or MO-CMA-ES, one can compare the shapes of the empirical cumulative distribution functions (ECDFs) for different population sizes. Note that for each algorithm and each COCO problem, the function is evaluated first with a zero-vector, since we know this produces a point with better f-values than a random sampled point. This allows rigorous comparison of the different algorithms ECDFs.

A glance on the data allows to surmise the following: the larger the population size, the better the reached target precisions are in the long run. However for small budgets, algorithms with smaller population size reach better target precisions. For example in Figure 1 on function **f28**, MO-10 is the first to solve 35% of the targets, MO-32 is the first to solve 45% for a larger budget, and finally MO-100 has the best overall performance with more than 50% of the targets solved. This result is not new and a MO-CMA-ES adapting the population size has been studied for example in [10]. It is still interesting to note that this tendency appears for almost all functions in the test suite.

The same comment can be made on the same function **f28** for COMO-CMA-ES with 3, 10 and 32 kernels: among the COMO-CMA-ES variants, COMO-3 is the first one to solve 40% of the targets, COMO-10 the first to solve 55%, and COMO-32 the first to solve 70%. Note that COMO-32 performs well at a budget of  $10^4 n$ , compared to the other COMO-CMA-ES. And for a budget of  $10^5 n$ , COMO-100 and COMO-316 are systematically better. We can expect the variants with more kernels (1000 is proposed here) to solve more targets for longer runs, where the ones with smaller number of kernels will stagnate sooner.

There are some functions where MO-CMA-ES outperforms the best 2016 data for small budgets, say up to  $10^2 n$ , see **f7**, **f17**, **f27**, **f32** on Figure 1. In contrast the COMO-CMA-ES' ECDFs dynamics are different: compared to the MO-CMA-ES with the same population size, it solves less targets on small budgets; there is a longer-lasting starting phase. This is partly due to the elitist nature of the (1+1)-CMA-ES used by MO-CMA-ES and the non-elitist nature of the

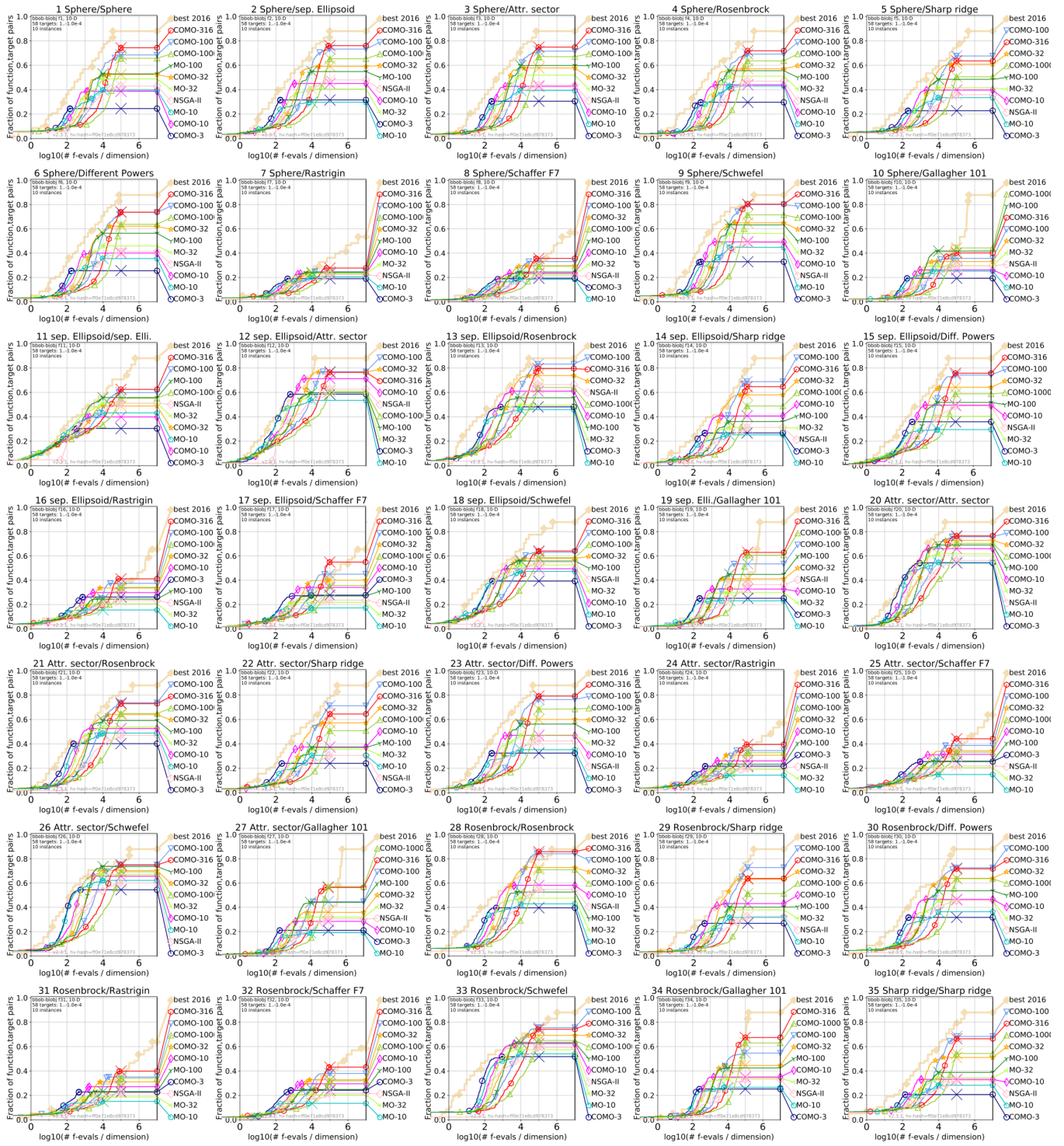
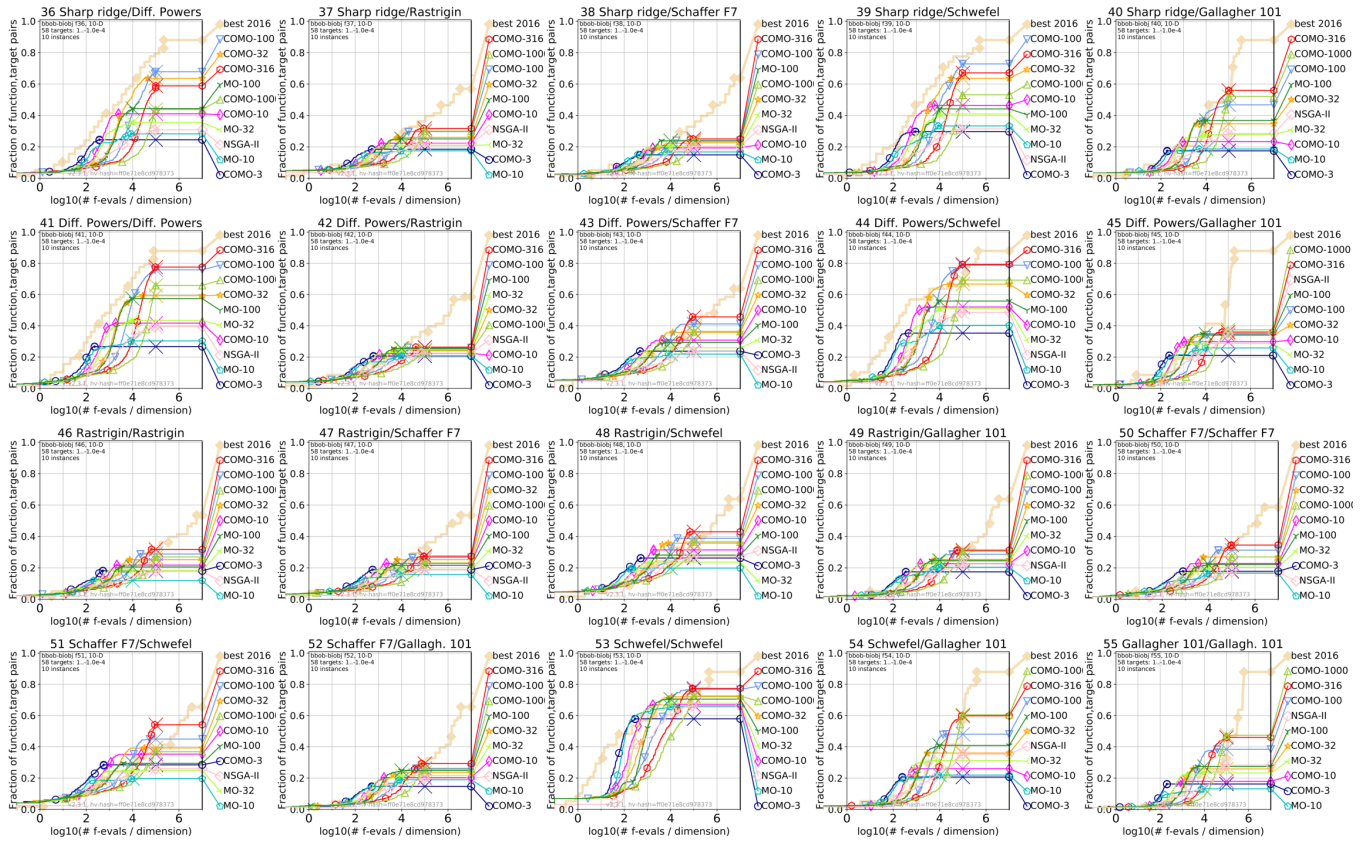


Figure 1: Empirical cumulative distribution of simulated (bootstrapped) runtimes, measured in number of objective function evaluations, divided by dimension (FEvals/DIM) for the 58 targets  $\{-10^{-4}, -10^{-4.2}, -10^{-4.4}, -10^{-4.6}, -10^{-4.8}, -10^{-5}, 0, 10^{-5}, 10^{-4.9}, 10^{-4.8}, \dots, 10^{-0.1}, 10^0\}$  in dimension 10.





**Figure 2: Bootstrapped empirical cumulative distribution of the number of objective function evaluations divided by dimension (FEvals/DIM) as in Fig. 1 but for functions  $f_{36}$  to  $f_{55}$  in 10-D.**

standard CMA-ES used by COMO-CMA-ES. However COMO-CMA-ES performs better in the long run, for instance with a budget of  $10^4 n$  on functions **f2**, **f4**, **f19**, **f53** in Figures 1, 2. Note that this effect (starting phase duration and fraction of targets solved) increases with the number of kernels. Remark that the MO-CMA-ES algorithms perform particularly well on problems where one of the functions is the Gallagher 101 (**f40**, **f45**, **f49**, **f52**, **f54**, **f55**) where a MO-\* systematically outperforms the best2016 reference.

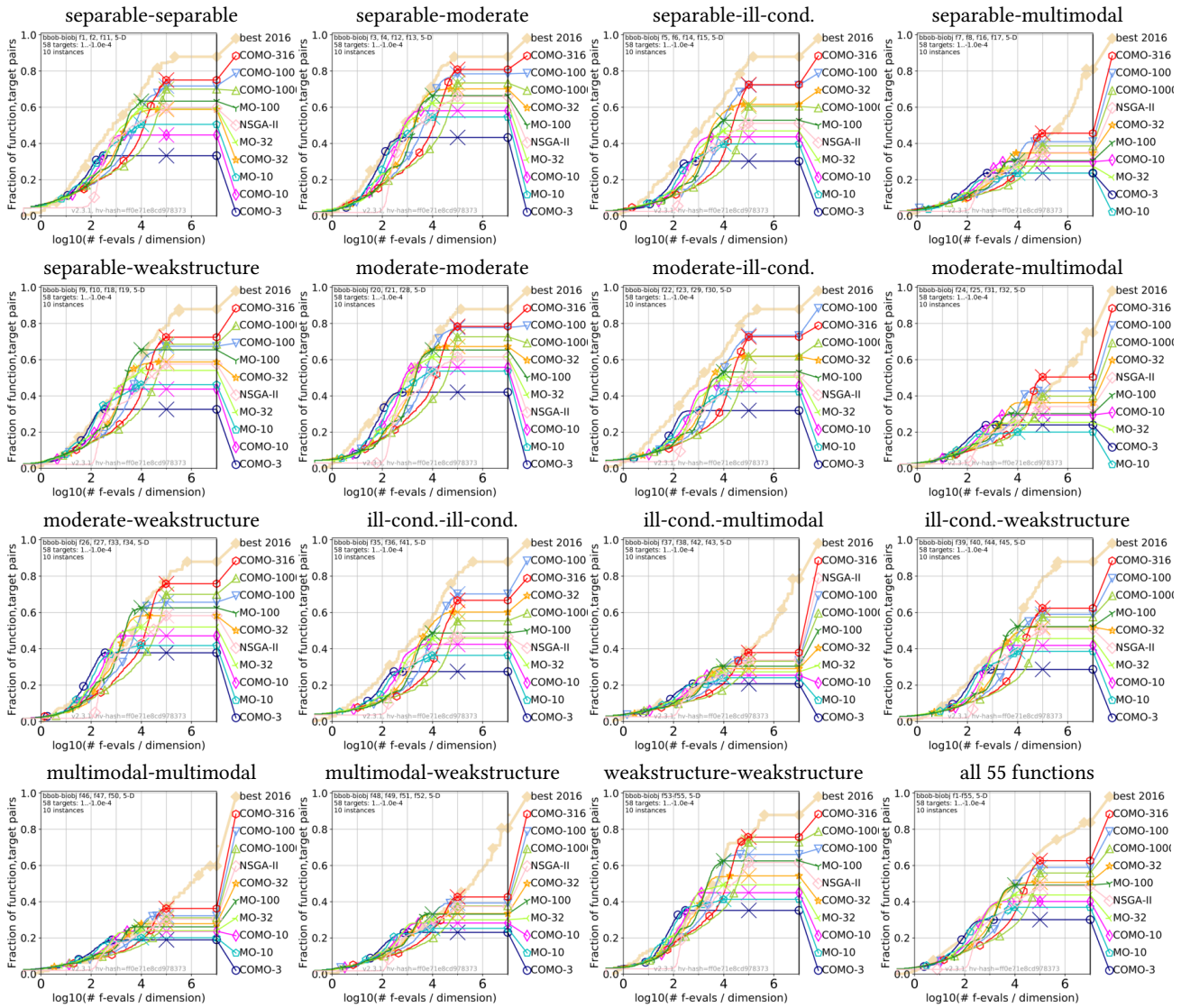
In 5-D on Figure 3, if we look at the fraction of targets found for a budget of  $10^4 n$ , a global statement is that NSGA-II is ranking below either MO-100 or COMO-100 (which have the same population size), but it can be above the variants with a population size far from 100. In 20-D, as we can see on Figure 4, for a budget of  $10^4 n$  and considering all subgroups of functions, COMO-100 solves more targets than MO-100 which in turn performs better than NSGA-II. As an example, on **ill-cond**, **ill-cond** subgroup, COMO-100 solves 50% of the target precisions versus roughly 40% for MO-100 and 20% for the NSGA-II. This kind of gap appears in most subgroups with one ill conditioned objective function. COMO-3 and COMO-10 also perform better than best2016 on subgroups with one multimodal objective function for small budgets, between  $10^2 n$  and  $10^4 n$ .

## 6 DISCUSSION / CONCLUSION

We have benchmarked COMO-CMA-ES and MO-CMA-ES with different population sizes to particularly test the influence of this parameter on the performances. We also used NSGA-II as a baseline. COMO-CMA-ES performs well although it is *a priori* designed to approximate the optimal  $p$ -distribution of the Hypervolume Indicator, for  $p$  the population size (number of kernels). As observed in [12], this performance is mainly due to the large stationary variance obtained with non-elitist (comma) evolution strategies and with the use of UHVI which guides the evolution towards the uncrowded space in the non-dominated region. The overall results for MO-CMA-ES and COMO-CMA-ES show that a smaller population size performs better for smaller budget, while a larger population size end up performing better for a budget large enough. Hence as done in [11] with MO-CMA-ES, a population size adaptation for COMO-CMA-ES should guarantee better performance on COCO.

## REFERENCES

- [1] Anne Auger, Dimo Brockhoff, Nikolaus Hansen, Dejan Tušar, Tea Tušar, and Tobias Wagner. 2016. Benchmarking MATLAB’s gamultiobj (NSGA-II) on the Bi-objective BBOB-2016 Test Suite. In *GECCO 2016 - Genetic and Evolutionary Computation Conference*. ACM, Denver, CO, United States, 1233–1239. DOI: <http://dx.doi.org/10.1145/2908961.2931706>
- [2] D. Brockhoff, T. Tušar, D. Tušar, T. Wagner, N. Hansen, and A. Auger. 2016. Biobjective Performance Assessment with the COCO Platform. *ArXiv e-prints* arXiv:1605.01746 (2016).



**Figure 3:** Bootstrapped empirical cumulative distribution of the number of objective function evaluations divided by dimension (FEvals/DIM) for 58 targets with target precision in  $\{-10^{-4}, -10^{-4.2}, -10^{-4.4}, -10^{-4.6}, -10^{-4.8}, -10^{-5}, 0, 10^{-5}, 10^{-4.9}, 10^{-4.8}, \dots, 10^{-0.1}, 10^0\}$  for all functions and subgroups in 5-D. As reference algorithm, the best algorithm from BBOB 2016 is shown as light thick line with diamond markers.

[3] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. 2002. A Fast and Elitist Multiobjective Genetic Algorithm: NSGA-II. *Trans. Evol. Comp.* 6, 2 (April 2002), 182–197. DOI: <http://dx.doi.org/10.1109/4235.996017>

[4] Nikolaus Hansen. 2016. The CMA evolution strategy: A tutorial. *arXiv preprint arXiv:1604.00772* (2016).

[5] N. Hansen, A. Auger, D. Brockhoff, D. Tušar, and T. Tušar. 2016. COCO: Performance Assessment. *ArXiv e-prints arXiv:1605.03560* (2016).

[6] N. Hansen, A. Auger, O. Mersmann, T. Tušar, and D. Brockhoff. 2016. COCO: A Platform for Comparing Continuous Optimizers in a Black-Box Setting. *ArXiv e-prints arXiv:1603.08785* (2016).

[7] N. Hansen, T. Tušar, O. Mersmann, A. Auger, and D. Brockhoff. 2016. COCO: The Experimental Procedure. *ArXiv e-prints arXiv:1603.08776* (2016).

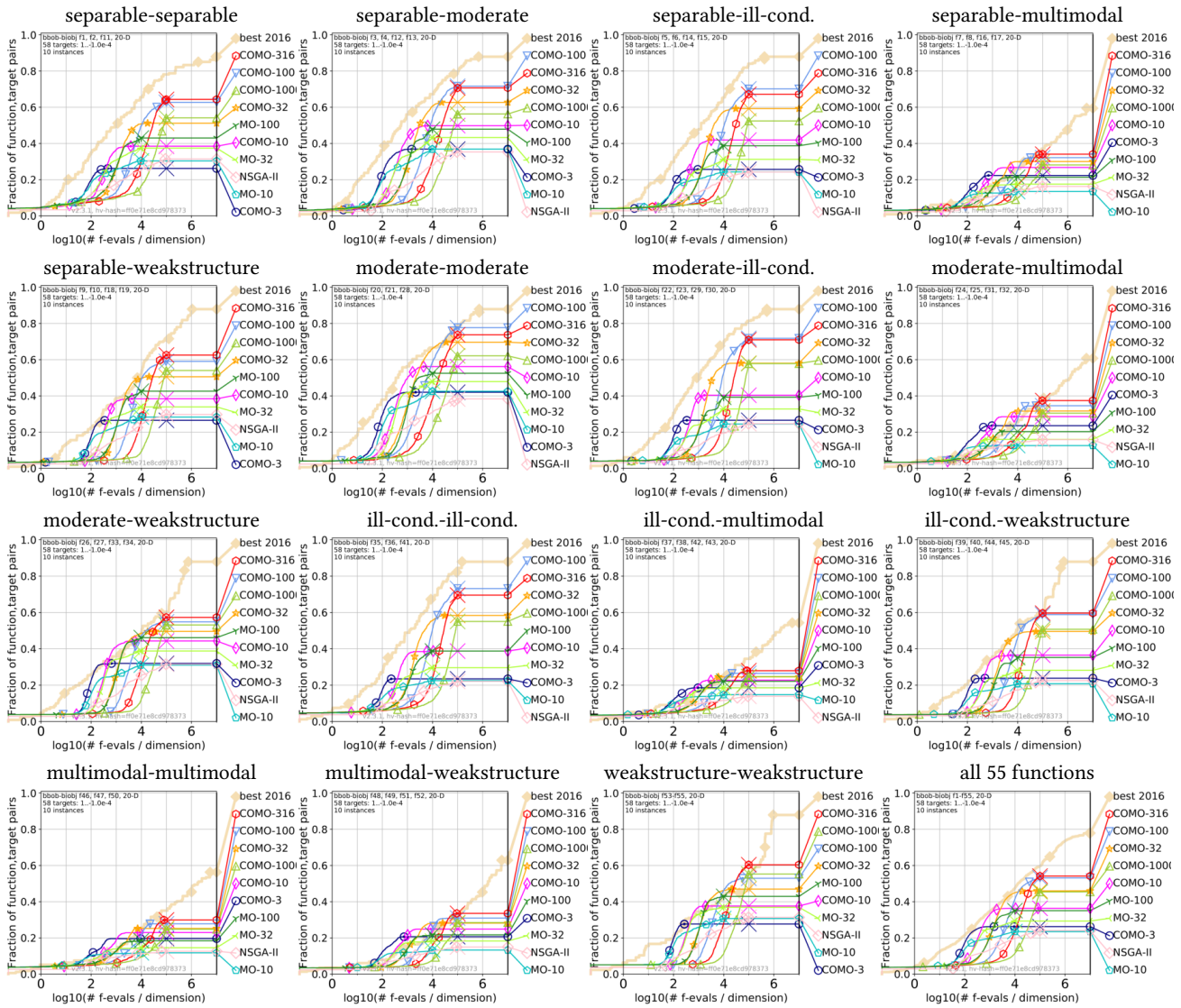
[8] Christian Igel, Nikolaus Hansen, and Stefan Roth. 2007. Covariance Matrix Adaptation for Multi-objective Optimization. *Evolutionary Computation* 15, 1 (2007), 1–28. DOI: <http://dx.doi.org/10.1162/evco.2007.15.1> arXiv:<https://doi.org/10.1162/evco.2007.15.1> PMID: 17388777.

[9] Oswin Krause, Tobias Glasmachers, Nikolaus Hansen, and Christian Igel. 2016. Unbounded population MO-CMA-ES for the bi-objective BBOB test suite. In *Proceedings of the 2016 on Genetic and Evolutionary Computation Conference Companion*. ACM, 1177–1184.

[10] Steffen Limmer and Dietmar Fey. 2016. Investigation of strategies for an increasing population size in multi-objective CMA-ES. In *2016 IEEE Congress on Evolutionary Computation (CEC)*. IEEE, Vancouver, BC, Canada, 476–483. DOI: <http://dx.doi.org/10.1109/CEC.2016.7743832>

[11] Ilya Loshchilov and Tobias Glasmachers. 2016. Anytime bi-objective optimization with a hybrid multi-objective CMA-ES (HMO-CMA-ES). In *Proceedings of the 2016 on Genetic and Evolutionary Computation Conference Companion*. ACM, 1169–1176.

[12] Cheikh Toure, Nikolaus Hansen, Anne Auger, and Dimo Brockhoff. 2019. Uncrowded Hypervolume Improvement: COMO-CMA-ES and the Sofomore framework. In *Proceedings of the Genetic and Evolutionary Computation Conference*. ACM, To appear.



**Figure 4:** Bootstrapped empirical cumulative distribution of the number of objective function evaluations divided by dimension (FEvals/DIM) for 58 targets with target precision in  $\{-10^{-4}, -10^{-4.2}, -10^{-4.4}, -10^{-4.6}, -10^{-4.8}, -10^{-5}, 0, 10^{-5}, 10^{-4.9}, 10^{-4.8}, \dots, 10^{-0.1}, 10^0\}$  for all functions and subgroups in 20-D. As reference algorithm, the best algorithm from BBOB 2016 is shown as light thick line with diamond markers.

[13] T. Tušar, D. Brockhoff, N. Hansen, and A. Auger. 2016. COCO: The Bi-objective Black-Box Optimization Benchmarking (bbob-biobj) Test Suite. *ArXiv e-prints* arXiv:1604.00359 (2016).

[14] Thomas Voß, Nikolaus Hansen, and Christian Igel. 2010. Improved Step Size Adaptation for the MO-CMA-ES. In *Genetic And Evolutionary Computation Conference*. ACM, Portland, United States, 487–494. DOI: <http://dx.doi.org/10.1145/1830483.1830573>

$\Delta f$	1e0	1e-1	1e-2	1e-3	$\#succ$
f1	1	75	584	3660	10/10
f2	5.0	105	601	3715	10/10
f3	3.0	115	665	5170	10/10
f4	2.0	109	571	2669	10/10
f5	2.0	120	1277	21889	10/10
f6	3.0	73	688	3976	10/10
f7	2.0	2763	1.2e5	3.5e6	10/10
f8	3.0	2167	1.6e5	2.1e6	10/10
f9	4.0	96	521	1986	10/10
f10	4.0	323	9839	52107	10/10
f11	5.0	56	436	9189	10/10
f12	5.0	44	641	3991	10/10
f13	7.0	60	560	5582	10/10
f14	5.0	247	1713	12801	10/10
f15	6.0	74	677	6559	10/10
f16	3.0	2810	2.0e5	4.2e6	10/10
f17	29	2935	48624	2.4e6	10/10
f18	2.0	56	557	6912	10/10
f19	9.0	1292	6164	88464	10/10
f20	4.0	70	1162	5724	10/10
f21	5.0	86	3249	9924	10/10
f22	3.0	97	1168	12608	10/10
f23	1	59	618	4410	10/10
f24	5.0	2347	1.8e5	4.1e6	10/10
f25	9.0	3143	1.2e5	2.5e6	10/10
f26	7.0	59	2182	13673	10/10
f27	6.0	2631	21971	45576	10/10
f28	2.0	20	145	1230	10/10
f29	3.0	114	1413	13660	10/10
f30	1	33	366	2294	10/10
f31	3.0	2166	50028	3.2e6	10/10
f32	3.0	2131	1.1e5	2.4e6	10/10
f33	6.0	527	1214	3730	10/10
f34	9.0	1376	15575	54195	10/10
f35	1	141	2268	51388	10/10
f36	2.0	204	4640	41237	10/10
f37	2.0	3956	1.5e5	2.5e6	10/10
f38	4.0	4366	2.7e5	5.4e6	10/10
f39	2.0	215	1160	47472	8/10
f40	2.0	998	34442	2.0e5	8/10
f41	2.0	48	708	6343	10/10
f42	2.0	2525	3.4e5	4.3e6	10/10
f43	4.0	2468	1.5e5	3.8e6	2/10
f44	6.0	86	513	1853	10/10
f45	4.0	816	5730	53653	10/10
f46	4.0	11925	4.9e5	5.6e7	10/10
f47	3.0	4172	2.7e5	4.4e6	10/10
f48	9.0	2160	1.1e5	2.1e6	2/10
f49	9.0	3737	2.0e5	1.3e6	1/10
f50	6.0	3913	2.2e5	2.6e6	1/10
f51	8.0	1251	32465	2.5e6	1/10
f52	2.0	3027	1.4e5	2.1e6	1/10
f53	2.0	13	42	1443	2/10
f54	98	1514	12856	45502	10/10
f55	3.0	1415	18086	51245	9/10

Table 2: Average runtime (aRT in number of function evaluations) divided by the respective best aRT measured during BBOB-2016 in dimension 5. This aRT ratio and, in braces as dispersion measure, the half difference between 10 and 90%-tile of bootstrapped run lengths appear for each algorithm and target, the corresponding reference aRT in the first row. The different target  $\Delta I_{HV}^{COCO}$ -values are shown in the top row. #succ is the number of trials that reached the (final) target  $I_{ref} + 10^{-5}$ . The median number of conducted function evaluations is additionally given in *italics*, if the target in the last column was never reached. Entries, succeeded by a star, are statistically significantly better (according to the rank-sum test) when compared to all other algorithms of the table, with  $p = 0.05$  or  $p = 10^{-k}$  when the number  $k$  following the star is larger than 1, with Bonferroni correction by the number of functions (55). A  $\downarrow$  indicates the target against the best algorithm from BBOB 2016. Best results are printed in bold.

