

AMaLGaM IDEAs in Noiseless Black-Box Optimization Benchmarking

Peter A.N. Bosman
Centre for Mathematics and
Computer Science
P.O. Box 94079
1090 GB Amsterdam
The Netherlands
Peter.Bosman@cwi.nl

Jörn Grahl
Johannes Gutenberg
University Mainz
Dept. of Information Systems
& Business Administration
Jakob Welder-Weg 9
D-55128 Mainz, Germany
grahl@uni-mainz.de

Dirk Thierens
Utrecht University
Dept. of Information and
Computing Sciences
P.O. Box 80089
3508 TB Utrecht
The Netherlands
Dirk.Thierens@cs.uu.nl

ABSTRACT

This paper describes the application of a Gaussian Estimation-of-Distribution (EDA) for real-valued optimization to the noiseless part of a benchmark introduced in 2009 called BBOB (Black-Box Optimization Benchmarking). Specifically, the EDA considered here is the recently introduced parameter-free version of the Adapted Maximum-Likelihood Gaussian Model Iterated Density-Estimation Evolutionary Algorithm (AMaLGaM-IDEA). Also the version with incremental model building (iAMaLGaM-IDEA) is considered.

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

Keywords

Benchmarking, Black-box optimization, Evolutionary computation

1. METHOD

Estimation-of-distribution algorithms (EDAs) [7, 8] are an important strand of research on black-box optimization (BBO). EDAs attempt to automatically exploit features of a problem's structure by probabilistically modeling the search space based on previously evaluated solutions and generating new solutions by sampling the probabilistic model.

The general EDA procedure is as follows. A population \mathcal{P} of n solutions is maintained. Through selection, a vector \mathcal{S} is selected from \mathcal{P} . A probability distribution over the solution space is then estimated using \mathcal{S} as a data set. New

solutions are generated by sampling the estimated probability distribution. Finally, the newly generated samples are incorporated into the population and the process repeats until a termination criterion has been satisfied.

The EDA considered here is the Adapted Maximum-Likelihood Gaussian Model Iterated Density-Estimation Evolutionary Algorithm (AMaLGaM-IDEA, or AMaLGaM for short). In AMaLGaM, the probability distribution used is the normal, also known as the Gaussian, distribution. This EDA uses maximum-likelihood estimates for the mean and the covariance matrix, estimated from the selected solutions. It has a mechanism that scales up the covariance matrix when required to prevent premature convergence on slopes. It furthermore has a mechanism that anticipates the mean shift in the next generation to speed up descent (in case of minimization) along slopes. For a more extensive description, we refer the interested reader to the literature [1].

In addition to the above base procedure, recently a parameter-free version of AMaLGaM was introduced [3]. After experimental analysis, settings were proposed for all parameters. Guidelines were developed for the minimally required population size that allows unimodal problems to be solved. On multimodal problems a restart mechanism is required to increase the probability of success. The specific restart scheme considered increases the number of solutions upon each restart by alternating between two approaches: a single run with a larger population and more parallel runs. To maximize the joint global effect of the parallel runs, their locality is increased by starting them in separate regions that are obtained from clustering the search space first. When increasing the number of parallel runs, the subpopulation size is also increased slightly so as to increase the robustness of the more localized searches.

Distribution estimation in AMaLGaM is done anew from scratch each generation. Subsequent iterations however have much in common and therefore the required population size can be reduced by incremental learning, i.e. combining the distribution estimated from \mathcal{S} with the distribution used in the previous generation. In iAMaLGaM a memory-decay approach is taken to this end. On unimodal problems the required population size was found to indeed be significantly reduced while at the same time requiring less function evaluations to reach the same solution quality. Results on multimodal landscapes indicated however that if memory resources are not very important, a larger base-population size helps

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. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

GECCO '09, July 8–12, 2009, Montréal Québec, Canada.
Copyright 2009 ACM 978-1-60558-505-5/09/07 ...\$5.00.

in optimizing multimodal problems, thus favoring the non-incremental approach. For this reason we tested both AMaL-GaM and iAMaL-GaM on the BBOB benchmark.

Next to the full covariance matrix, two other versions of AMaL-GaM exist that reduce the number of distribution parameters to be estimated. One version uses Bayesian factorizations to select only the most important covariances while another version allows only variances. If only a few dependencies between problem variables exist, these methods outperform the use of the full covariance matrix in asymptotic complexity for the scalability in terms of required function evaluations and required time. These restrictions however also render the EDA non-rotationally invariant and therefore less generally applicable. For this reason and for the sake of space, we do not submit these variants to the BBOB benchmark here. A closer look at the differences with the full covariance matrix can be found in [3]; BBOB benchmarks for additional variants are given in [2].

For technical completeness, pseudo-code is presented below. A note on the pseudo-code: in iAMaL-GaM, for $\hat{\Sigma}(0)$ a matrix with the ML variances on the diagonal and zeros off the diagonal is used. Also, $\hat{\mu}^{\text{Shift}}(t)$ is non-existent for $t = 0$ and for $t = 1$ it is $\hat{\mu}(1) - \hat{\mu}(0)$. SDR stands for standard-deviation ratio, NIS stands for no-improvement stretch.

```
(i)AMaL-GaM-Free
1  $s \leftarrow 0; n^{\text{Base}} \leftarrow 17 + 3D^{1.5}$  (iAMaL-GaM:  $n^{\text{Base}} \leftarrow 10D^{0.5}$ )
2 do
3   if  $(s \bmod 2) = 0$  then  $n \leftarrow (1 + s/2)n^{\text{Base}}; p \leftarrow 2^{s/2}$ 
4   else  $n \leftarrow 2^{1+s/2}n^{\text{Base}}; p \leftarrow 1$ 
5   Run (i)AMaL-GaM with population size  $n$  and  $p$  parallel runs,
   starting from the clustering of  $np$  randomly generated solutions
   into  $p$  clusters and using  $\eta^{\text{DEC}} \leftarrow 0.9; \eta^{\text{INC}} \leftarrow 1/\eta^{\text{DEC}}; \theta^{\text{SDR}} \leftarrow 1;$ 
    $\tau \leftarrow 0.35; \alpha^{\text{AMS}} \leftarrow \frac{1}{2}\tau(n/(n-1)); \delta^{\text{AMS}} \leftarrow 2; \text{NIS}^{\text{MAX}} \leftarrow 25 + D$ 
6    $s \leftarrow s + 1$ 
7   while optimum not found and max. eval. not reached
```

```
(i)AMaL-GaM
1  $\eta^{\Sigma} \leftarrow 1; \eta^{\text{Shift}} \leftarrow 1$ 
   (iAMaL-GaM:  $\eta^{\Sigma} \leftarrow 1 - e^{-1.1 \lfloor \tau n \rfloor^{1.2} / D^{1.6}}; \eta^{\text{Shift}} \leftarrow 1 - e^{-1.2 \lfloor \tau n \rfloor^{0.31} / D^{0.50}}$ )
2  $c^{\text{Multiplier}} \leftarrow 1; n^{\text{AMS}} \leftarrow \alpha^{\text{AMS}}(n-1); \text{NIS} \leftarrow 0; t \leftarrow 0$ 
3 do
4    $\mathcal{S} \leftarrow$  the best  $\lfloor \tau n \rfloor$  solutions in  $\mathcal{P}$  (truncation selection)
5    $\hat{\mu}(t) \leftarrow \frac{1}{|\mathcal{S}|} \sum_{i=0}^{|\mathcal{S}|-1} \mathcal{S}_i$ 
6    $\hat{\Sigma}(t) \leftarrow (1 - \eta^{\Sigma})\hat{\Sigma}(t-1) + \eta^{\Sigma} \frac{1}{|\mathcal{S}|} \sum_{i=0}^{|\mathcal{S}|-1} (\mathcal{S}_i - \hat{\mu}(t)) (\mathcal{S}_i - \hat{\mu}(t))^T$ 
7    $\hat{\mu}^{\text{Shift}}(t) \leftarrow (1 - \eta^{\text{Shift}})\hat{\mu}^{\text{Shift}}(t-1) + \eta^{\text{Shift}} (\hat{\mu}(t) - \hat{\mu}(t-1))$ 
8    $\hat{\mu} \leftarrow \hat{\mu}(t); \hat{\Sigma} \leftarrow c^{\text{Multiplier}}\hat{\Sigma}(t); \mathbf{LL}^* \leftarrow$  Cholesky decomp. of  $\hat{\Sigma}$ 
9    $\mathcal{P}_0 \leftarrow$  the best solution in  $\mathcal{S}$ 
10   $\mathcal{P}_{1..n-1} \leftarrow n-1$  samples from  $\mathcal{N}(\hat{\mu}, \hat{\Sigma}) = \hat{\mu} + \mathbf{L}\mathbf{N}(\mathbf{0}, \mathbf{I})$ 
11  for  $n^{\text{AMS}}$  random solutions  $\mathcal{P}_j$  ( $1 \leq j \leq n-1$ )
12  do  $\mathcal{P}_j \leftarrow \mathcal{P}_j + \delta^{\text{AMS}} c^{\text{Multiplier}} \hat{\mu}^{\text{Shift}}(t)$ 
13  if any  $\mathcal{P}_i$  better than  $\mathcal{P}_0$  ( $1 \leq i \leq n-1$ )
14  then
15     $\text{NIS} \leftarrow 0$ 
16    if  $c^{\text{Multiplier}} < 1$  then  $c^{\text{Multiplier}} \leftarrow 1$ 
17     $\mathbf{x}^{\text{avg-imp}} \leftarrow$  average of all  $\mathcal{P}_i$  better than  $\mathcal{P}_0$  ( $1 \leq i \leq n-1$ )
18     $\text{SDR} \leftarrow \max_{0 \leq i \leq D-1} \{ |(\mathbf{L}^{-1}(\mathbf{x}^{\text{avg-imp}} - \hat{\mu}))_i| \}$ 
19    if  $\text{SDR} > \theta^{\text{SDR}}$  then  $c^{\text{Multiplier}} \leftarrow \eta^{\text{INC}} c^{\text{Multiplier}}$ 
20  else
21    if  $c^{\text{Multiplier}} \leq 1$  then  $\text{NIS} \leftarrow \text{NIS} + 1$ 
22    if  $(c^{\text{Multiplier}} > 1)$  or  $(\text{NIS} \geq \text{NIS}^{\text{MAX}})$ 
23    then  $c^{\text{Multiplier}} \leftarrow \eta^{\text{DEC}} c^{\text{Multiplier}}$ 
24    if  $(c^{\text{Multiplier}} < 1)$  and  $(\text{NIS} < \text{NIS}^{\text{MAX}})$  then  $c^{\text{Multiplier}} \leftarrow 1$ 
25     $t \leftarrow t + 1$ 
26  while opt. not found, max. eval. not reached and  $c^{\text{Multiplier}} \geq 10^{-10}$ 
```

2. PARAMETERS AND OTHER SETTINGS

For initialization, a uniform sampling in $[-5, 5]^D$ was used, where D denotes the dimension of the search space. The experiments according to [5] on the benchmark functions given in [4, 6] have been conducted using the provided C-code.

The AMaL-GaM implementation used is also in C. A maximum of $10^6 D$ function evaluations is allowed. No changes were made to parameter-free AMaL-GaM as described in [3] and as outlined above. Therefore no parameter tuning was required and the crafting effort CrE [5] is zero.

3. CPU TIMING EXPERIMENT

For the timing experiment the full covariance matrix variant for both AMaL-GaM and iAMaL-GaM were run with a maximum of $10^6 D$ function evaluations and restarted until 30 seconds had passed (according to Figure 2 in [5]). The experiments have been conducted on an Intel Q6600 Core2Quad 2.4 GHz processor under Fedora Linux release 10 (Cambridge). In 2, 3, 5, 10, 20 and 40 dimensions, the time in 10^{-7} seconds per function evaluation was as follows:

	2	3	5	10	20	40
AMaL-GaM	1.9	2.2	3.0	5.0	10	24
iAMaL-GaM	1.9	2.3	3.0	5.3	11	29

4. RESULTS AND CONCLUSION

Results from experiments according to [5] on the benchmark functions given in [4, 6] are presented in Figures 1 and 2 and in Table 1 for AMaL-GaM and in Figures 3 and 4 and in Table 2 for iAMaL-GaM.

Problems with weak structure appear to be the hardest for (i)AMaL-GaM. Even within $10^6 D$ evaluations the optimum cannot be found within a desirable precision, especially for larger D . The difference between AMaL-GaM and iAMaL-GaM is not large which supports the design of the population-size reducing incremental-learning approach used. Consistent with earlier findings, the incremental approach is better on unimodal functions, whereas the non-incremental approach is (slightly) better on multimodal functions, most likely due to the larger base population-size.

5. REFERENCES

- [1] P. A. N. Bosman, J. Grahl, and D. Thierens. Enhancing the performance of maximum-likelihood Gaussian EDAs using anticipated mean shift. In G. Rudolph et al., editors, *Parallel Problem Solving from Nature — PPSN X*, pages 133–134, Berlin, 2008. Springer-Verlag.
- [2] P. A. N. Bosman, J. Grahl, and D. Thierens. A parameter-free Gaussian EDA called AMaL-GaM-IDEA: algorithms and benchmarks. CWI technical report (*To Appear*), 2009.
- [3] P.A.N. Bosman. On empirical memory design, faster selection of Bayesian factorizations and parameter-free Gaussian EDAs. In G. Raidl et al., editors, *Proc. of the Genetic and Evolutionary Computation Conference — GECCO-2009*, New York, New York, 2009. ACM Press. (*To Appear*).
- [4] S. Finck, N. Hansen, R. Ros, and A. Auger. Real-parameter black-box optimization benchmarking 2009: Presentation of the noiseless functions. Technical Report 2009/20, Research Center PPE, 2009.
- [5] N. Hansen, A. Auger, S. Finck, and R. Ros. Real-parameter black-box optimization benchmarking 2009: Experimental setup. Technical Report RR-6828, INRIA, 2009.
- [6] N. Hansen, S. Finck, R. Ros, and A. Auger. Real-parameter black-box optimization benchmarking 2009: Noiseless functions definitions. Technical Report RR-6829, INRIA, 2009.
- [7] J. A. Lozano, P. Larrañaga, I. Inza, and E. Bengoetxea. *Towards a New Evolutionary Computation. Advances in Estimation of Distribution Algorithms*. Springer-Verlag, Berlin, 2006.
- [8] M. Pelikan, K. Sastry, and E. Cantú-Paz. *Scalable Optimization via Probabilistic Modeling: From Algorithms to Applications*. Springer-Verlag, Berlin, 2006.

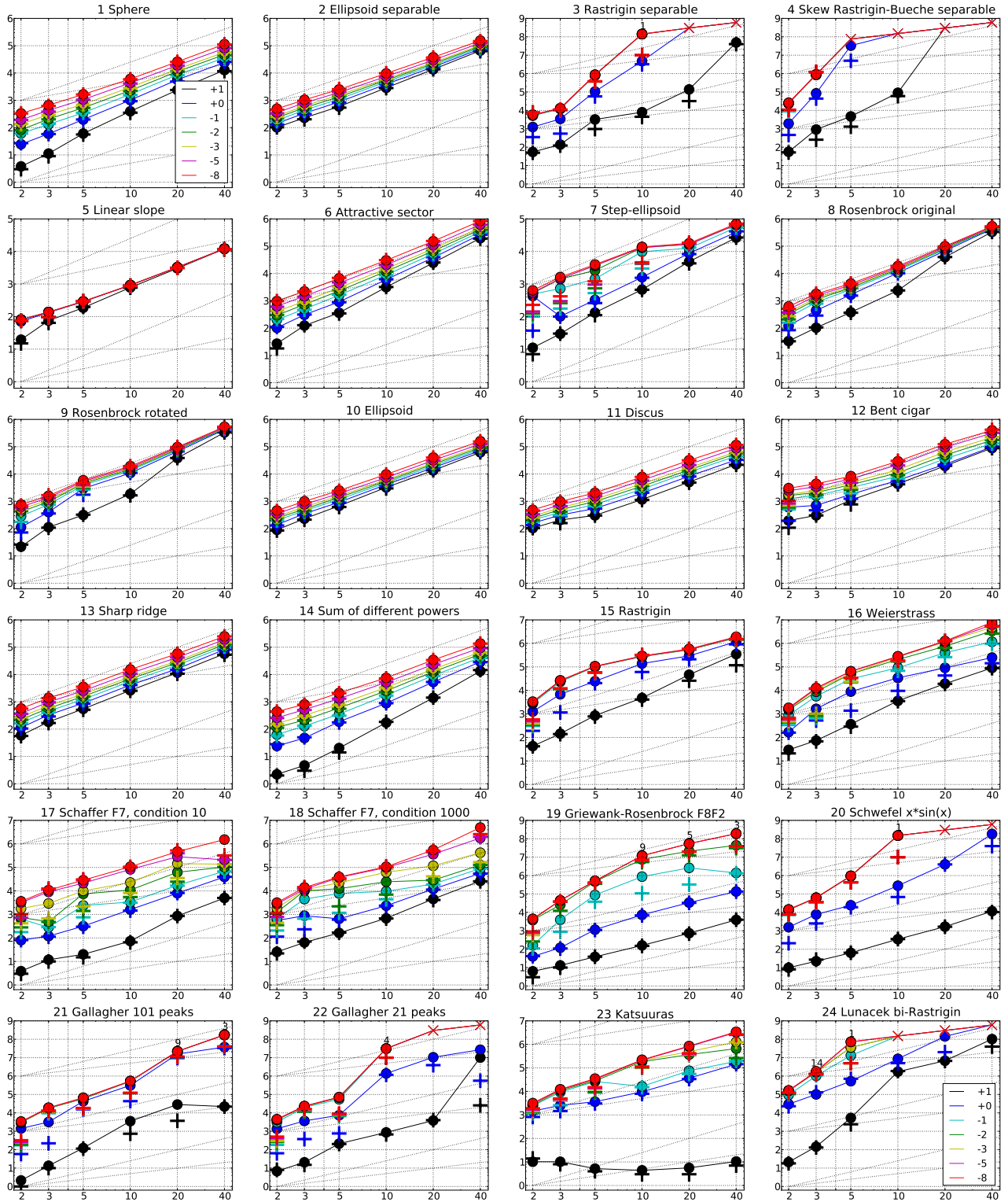


Figure 1: AMaLGaM: Expected Running Time (ERT, \bullet) to reach $f_{\text{opt}} + \Delta f$ and median number of function evaluations of successful trials (+), shown for $\Delta f = 10, 1, 10^{-1}, 10^{-2}, 10^{-3}, 10^{-5}, 10^{-8}$ (the exponent is given in the legend of f_1 and f_{24}) versus dimension in log-log presentation. The ERT(Δf) equals to $\#FES(\Delta f)$ divided by the number of successful trials, where a trial is successful if $f_{\text{opt}} + \Delta f$ was surpassed during the trial. The $\#FES(\Delta f)$ are the total number of function evaluations while $f_{\text{opt}} + \Delta f$ was not surpassed during the trial from all respective trials (successful and unsuccessful), and f_{opt} denotes the optimal function value. Crosses (\times) indicate the total number of function evaluations $\#FES(-\infty)$. Numbers above ERT-symbols indicate the number of successful trials. Annotated numbers on the ordinate are decimal logarithms. Additional grid lines show linear and quadratic scaling.

f_1 in 5-D, N=15, mFE=2108				f_1 in 20-D, N=15, mFE=32945				f_2 in 5-D, N=15, mFE=2941				f_2 in 20-D, N=15, mFE=46861										
Δf	#	ERT	90%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}					
10	15	6.1e1	5.2e1	6.9e1	6.1e1	15	2.4e3	2.3e3	2.5e3	2.4e3	10	15	5.9e2	5.5e2	6.4e2	5.9e2	15	1.4e4	1.3e4	1.4e4	1.4e4	
1	15	2.0e2	1.8e2	2.1e2	2.0e2	15	5.7e3	5.2e3	6.2e3	5.7e3	1	15	8.7e2	8.0e2	9.5e2	8.7e2	15	1.7e4	1.6e4	1.7e4	1.7e4	
1e-1	15	3.5e2	3.3e2	3.7e2	3.5e2	15	8.5e3	7.9e3	9.2e3	8.5e3	1e-1	15	1.1e3	1.0e3	1.2e3	1.1e3	15	1.9e4	1.8e4	2.0e4	1.9e4	
1e-3	15	7.1e2	6.7e2	7.5e2	7.1e2	15	1.4e4	1.3e4	1.5e4	1.4e4	1e-3	15	1.6e3	1.5e3	1.7e3	1.6e3	15	2.4e4	2.3e4	2.5e4	2.4e4	
1e-5	15	1.1e3	1.0e3	1.1e3	1.1e3	15	1.9e4	1.8e4	2.0e4	1.9e4	1e-5	15	1.9e3	1.8e3	2.0e3	1.9e3	15	3.0e4	2.9e4	3.1e4	3.0e4	
1e-8	15	1.6e3	1.6e3	1.7e3	1.6e3	15	2.6e4	2.5e4	2.7e4	2.6e4	1e-8	15	2.5e3	2.4e3	2.6e3	2.5e3	15	3.8e4	3.6e4	3.9e4	3.8e4	
f_3 in 5-D, N=15, mFE=2.66e6				f_3 in 20-D, N=15, mFE=2.00e7				f_4 in 5-D, N=15, mFE=5.01e6				f_4 in 20-D, N=15, mFE=2.00e7										
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}					
10	15	3.2e3	9.4e2	5.5e3	3.2e3	15	1.4e5	9.4e4	1.8e5	1.4e5	10	15	4.7e3	2.3e3	7.1e3	4.7e3	0	<i>14e+0</i>	<i>13e+0</i>	<i>15e+0</i>	4.5e6	
1	15	1.1e5	7.7e4	1.4e5	1.1e5	0	<i>40e-1</i>	<i>30e-1</i>	<i>50e-1</i>	5.6e5	1	2	3.3e7	1.8e7	>7e7	5.0e6						
1e-1	15	7.9e5	5.2e5	1.1e6	7.9e5						1e-1	0	<i>20e-1</i>	<i>99e-2</i>	<i>30e-1</i>	8.9e5						
1e-3	15	8.5e5	5.6e5	1.1e6	8.5e5						1e-3											
1e-5	15	8.6e5	5.6e5	1.2e6	8.6e5						1e-5											
1e-8	15	8.7e5	5.8e5	1.2e6	8.7e5						1e-8											
f_5 in 5-D, N=15, mFE=491				f_5 in 20-D, N=15, mFE=4545				f_6 in 5-D, N=15, mFE=7498				f_6 in 20-D, N=15, mFE=168697										
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}					
10	15	1.9e2	1.8e2	2.0e2	1.9e2	15	3.0e3	2.8e3	3.2e3	3.0e3	10	15	3.6e2	3.2e2	4.0e2	3.6e2	15	2.5e4	2.4e4	2.6e4	2.5e4	
1	15	2.8e2	2.6e2	3.1e2	2.8e2	15	3.2e3	3.0e3	3.4e3	3.2e3	1	15	9.1e2	8.3e2	1.0e3	9.1e2	15	3.8e4	3.7e4	4.0e4	3.8e4	
1e-1	15	2.9e2	2.7e2	3.1e2	2.9e2	15	3.3e3	3.1e3	3.5e3	3.3e3	1e-1	15	1.6e3	1.5e3	1.8e3	1.6e3	15	5.2e4	5.0e4	5.4e4	5.2e4	
1e-3	15	2.9e2	2.7e2	3.1e2	2.9e2	15	3.3e3	3.0e3	3.5e3	3.3e3	1e-3	15	3.0e3	2.8e3	3.3e3	3.0e3	15	8.1e4	7.9e4	8.3e4	8.1e4	
1e-5	15	2.9e2	2.7e2	3.1e2	2.9e2	15	3.3e3	3.1e3	3.5e3	3.3e3	1e-5	15	4.4e3	4.2e3	4.7e3	4.4e3	15	1.1e5	1.1e5	1.1e5	1.1e5	
1e-8	15	2.9e2	2.7e2	3.1e2	2.9e2	15	3.3e3	3.0e3	3.5e3	3.3e3	1e-8	15	6.7e3	6.4e3	6.9e3	6.7e3	15	1.5e5	1.5e5	1.6e5	1.5e5	
f_7 in 5-D, N=15, mFE=14818				f_7 in 20-D, N=15, mFE=21017				f_8 in 5-D, N=15, mFE=5440				f_8 in 20-D, N=15, mFE=116157										
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}					
10	15	1.3e2	1.0e2	1.6e2	1.3e2	15	4.9e3	4.4e3	5.4e3	4.9e3	10	15	3.8e2	3.4e2	4.1e2	3.8e2	15	4.0e4	3.9e4	4.1e4	4.0e4	
1	15	3.2e2	2.8e2	3.8e2	3.2e2	15	8.8e3	8.2e3	9.4e3	8.8e3	1	15	1.7e3	1.5e3	1.8e3	1.7e3	15	6.8e4	6.6e4	7.0e4	6.8e4	
1e-1	15	1.4e3	5.3e2	2.9e3	1.4e3	15	1.3e4	1.2e4	1.3e4	1.3e4	1e-1	15	2.6e3	2.4e3	2.8e3	2.6e3	15	7.7e4	7.5e4	7.8e4	7.7e4	
1e-3	15	3.7e3	1.9e3	5.5e3	3.7e3	15	1.7e4	1.6e4	1.7e4	1.7e4	1e-3	15	3.3e3	3.1e3	3.6e3	3.3e3	15	8.6e4	8.5e4	8.8e4	8.6e4	
1e-5	15	3.7e3	1.9e3	5.5e3	3.7e3	15	1.7e4	1.6e4	1.7e4	1.7e4	1e-5	15	3.8e3	3.5e3	4.0e3	3.8e3	15	9.2e4	9.0e4	9.4e4	9.2e4	
1e-8	15	4.0e3	2.2e3	5.7e3	4.0e3	15	1.8e4	1.7e4	1.9e4	1.8e4	1e-8	15	4.3e3	4.0e3	4.5e3	4.3e3	15	1.0e5	9.7e4	1.0e5	1.0e5	
f_9 in 5-D, N=15, mFE=25268				f_9 in 20-D, N=15, mFE=112465				f_{10} in 5-D, N=15, mFE=3382				f_{10} in 20-D, N=15, mFE=46293										
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}					
10	15	3.2e2	2.9e2	3.4e2	3.2e2	15	3.8e4	3.8e4	3.9e4	3.8e4	10	15	6.8e2	6.1e2	7.6e2	6.8e2	15	1.3e4	1.3e4	1.4e4	1.3e4	
1	15	2.9e3	1.6e3	4.3e3	2.9e3	15	6.7e4	6.6e4	6.8e4	6.7e4	1	15	9.2e2	8.3e2	1.0e3	9.2e2	15	1.7e4	1.6e4	1.8e4	1.7e4	
1e-1	15	3.9e3	2.5e3	5.4e3	3.9e3	15	7.5e4	7.4e4	7.6e4	7.5e4	1e-1	15	1.2e3	1.1e3	1.3e3	1.2e3	15	2.0e4	1.9e4	2.1e4	2.0e4	
1e-3	15	4.8e3	3.4e3	6.4e3	4.8e3	15	8.4e4	8.2e4	8.5e4	8.4e4	1e-3	15	1.6e3	1.4e3	1.7e3	1.6e3	15	2.6e4	2.5e4	2.7e4	2.6e4	
1e-5	15	5.3e3	3.8e3	6.9e3	5.3e3	15	8.9e4	8.8e4	9.1e4	8.9e4	1e-5	15	1.9e3	1.8e3	2.0e3	1.9e3	15	3.2e4	3.0e4	3.3e4	3.2e4	
1e-8	15	5.9e3	4.4e3	7.5e3	5.9e3	15	9.7e4	9.5e4	9.9e4	9.7e4	1e-8	15	2.5e3	2.3e3	2.6e3	2.5e3	15	4.0e4	3.9e4	4.2e4	4.0e4	
f_{11} in 5-D, N=15, mFE=2549				f_{11} in 20-D, N=15, mFE=40045				f_{12} in 5-D, N=15, mFE=19209				f_{12} in 20-D, N=15, mFE=144557										
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}					
10	15	3.0e2	2.8e2	3.3e2	3.0e2	15	5.1e3	4.7e3	5.4e3	5.1e3	10	15	1.1e3	8.3e2	1.4e3	1.1e3	15	2.0e4	1.9e4	2.1e4	2.0e4	
1	15	5.5e2	4.9e2	6.1e2	5.5e2	15	8.3e3	7.8e3	8.8e3	8.3e3	1	15	1.8e3	1.2e3	2.5e3	1.8e3	15	3.3e4	2.1e4	4.4e4	2.3e4	
1e-1	15	7.6e2	6.9e2	8.4e2	7.6e2	15	1.2e4	1.1e4	1.2e4	1.2e4	1e-1	15	3.2e3	2.3e3	4.1e3	3.2e3	15	2.6e4	3.3e4	3.9e4	3.6e4	
1e-3	15	1.2e3	1.1e3	1.3e3	1.2e3	15	1.7e4	1.6e4	1.8e4	1.7e4	1e-3	15	4.9e3	3.8e3	6.0e3	4.9e3	15	6.4e4	6.0e4	6.7e4	6.4e4	
1e-5	15	1.5e3	1.4e3	1.6e3	1.5e3	15	2.3e4	2.1e4	2.4e4	2.3e4	1e-5	15	6.6e3	5.3e3	8.0e3	6.6e3	15	9.5e4	9.1e4	9.8e4	9.5e4	
1e-8	15	2.0e3	1.9e3	2.2e3	2.0e3	15	3.2e4	3.0e4	3.3e4	3.2e4	1e-8	15	8.3e3	6.8e3	9.9e3	8.3e3	15	1.2e5	1.2e5	1.3e5	1.2e5	
f_{13} in 5-D, N=15, mFE=4019				f_{13} in 20-D, N=15, mFE=70149				f_{14} in 5-D, N=15, mFE=2892				f_{14} in 20-D, N=15, mFE=42885										
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}					
10	15	5.5e2	4.9e2	6.1e2	5.5e2	15	1.2e4	1.1e4	1.3e4	1.2e4	10	15	2.0e1	1.5e1	2.6e1	2.0e1	15	1.4e3	1.3e3	1.5e3	1.4e3	
1	15	8.9e2	8.2e2	9.7e2	8.9e2	15	1.6e4	1.5e4	1.7e4	1.6e4	1	15	1.9e2	1.7e2	2.1e2	1.9e2	15	5.3e3	4.8e3	5.8e3	5.3e3	
1e-1	15	1.2e3	1.1e3	1.3e3	1.2e3	15	2.2e4	2.1e4	2.3e4	2.2e4	1e-1	15	3.5e2	3.3e2	3.8e2	3.5e2	15	9.1e3	8.4e3	9.6e3	9.1e3	
1e-3	15	1.9e3	1.8e3	2.0e3	1.9e3	15	3.1e4	3.0e4	3.2e4	3.1e4	1e-3	15	8.1e2	7.5e2	8.7e2	8.1e2	15	1.6e4	1.5e4	1.6e4	1.6e4	
1e-5	15	2.5e3	2.4e3	2.7e3	2.5e3	15	4.2e4	4.2e4	4.3e4	4.2e4	1e-5	15	1.3e3	1.2e3	1.4e3	1.3e3	15	2.3e4	2.2e4	2.5e4	2.3e4	
1e-8	15	3.5e3	3.4e3	3.6e3	3.5e3	15	5.8e4	5.7e4	5.9e4	5.8e4	1e-8	15	2.1e3	1.9e3	2.2e3	2.1e3	15	3.4e4	3.3e4	3.6e4	3.4e4	
f_{15} in 5-D, N=15, mFE=622081				f_{15} in 20-D, N=15, mFE=1.12e6				f_{16} in 5-D, N=15, mFE=255129				f_{16} in 20-D, N=15, mFE=2.48e6										
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}	#	ERT	10%	90%	90%	RT _{succ}					
10	15	8.5e2	7.7e2	9.4e2	8.5e2	15	4.7e4	2.7e4	6.6e4	4.7e4	10	15	3.7e2	2.8e2	4.6e2	3.7e2	15	1.9e4	1.8e4	2.1e4	1.9e4	
1	15	2.5e4	2.1e4	2.9e4	2.5e4	15	3.0e5	2.5e5	3.4e5	3.0e5</												

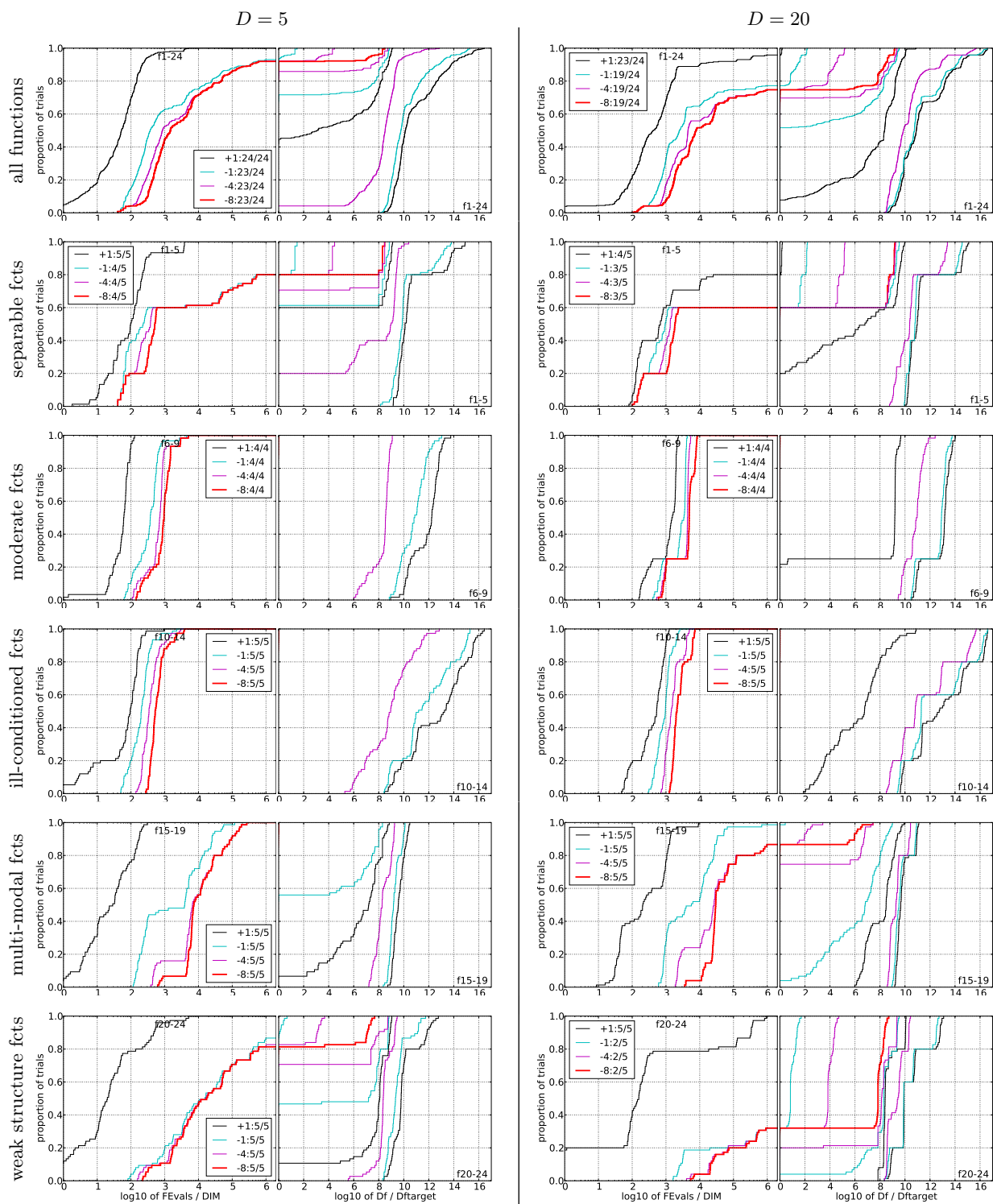


Figure 2: AMaLGaM: Empirical cumulative distribution functions (ECDFs), plotting the fraction of trials versus running time (left) or Δf . Left subplots: ECDF of the running time (number of function evaluations), divided by search space dimension D , to fall below $f_{\text{opt}} + \Delta f$ with $\Delta f = 10^k$, where k is the first value in the legend. Right subplots: ECDF of the best achieved Δf divided by 10^k (upper left lines in continuation of the left subplot), and best achieved Δf divided by 10^{-8} for running times of $D, 10D, 100D \dots$ function evaluations (from right to left cycling black-cyan-magenta). Top row: all results from all functions; second row: separable functions; third row: misc. moderate functions; fourth row: ill-conditioned functions; fifth row: multi-modal functions with adequate structure; last row: multi-modal functions with weak structure. The legends indicate the number of functions that were solved in at least one trial. FEvals denotes number of function evaluations, D and DIM denote search space dimension, and Δf and Df denote the difference to the optimal function value.

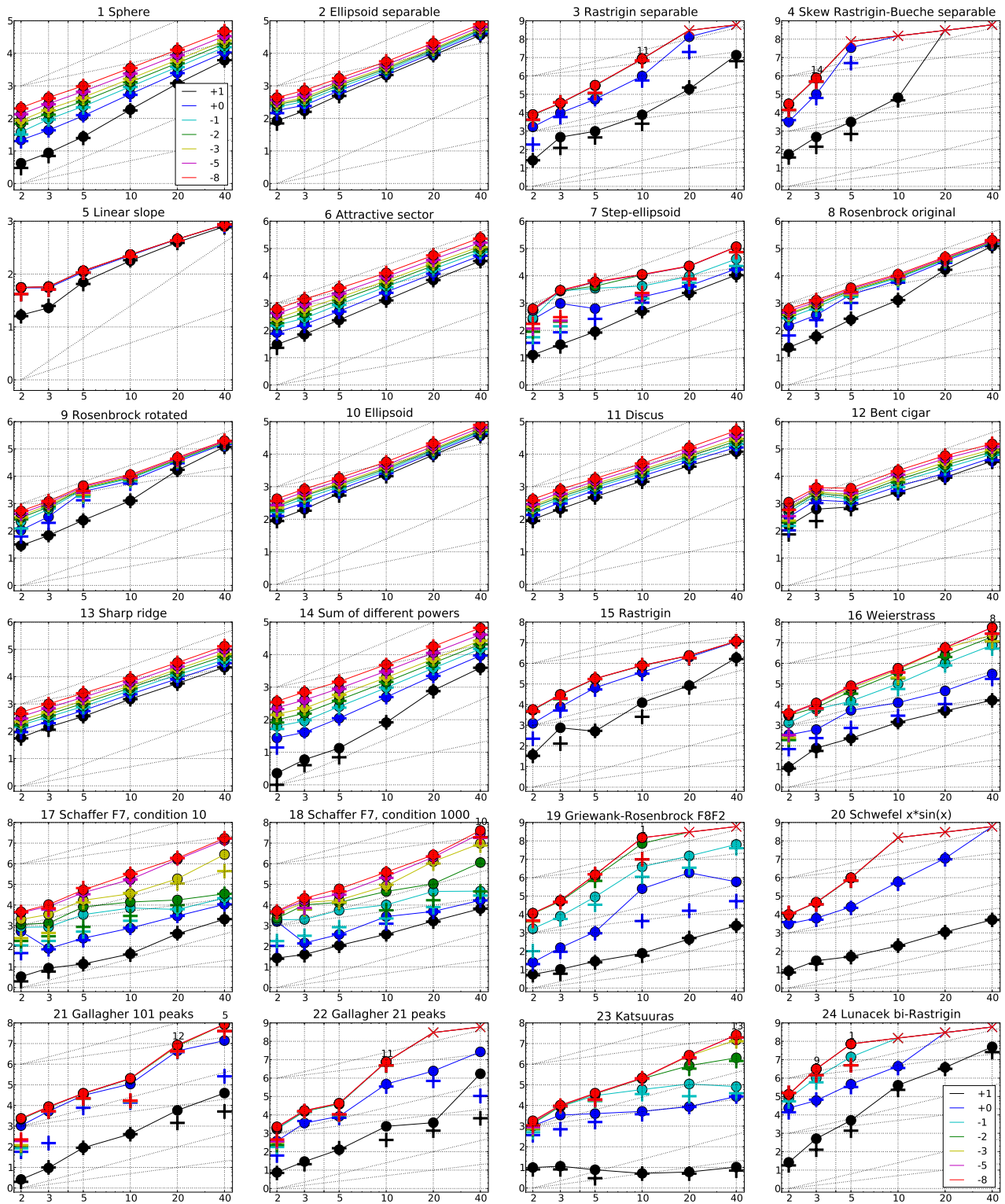


Figure 3: iAMaLGaM: Expected Running Time (ERT, \bullet) to reach $f_{\text{opt}} + \Delta f$ and median number of function evaluations of successful trials (+), shown for $\Delta f = 10, 1, 10^{-1}, 10^{-2}, 10^{-3}, 10^{-5}, 10^{-8}$ (the exponent is given in the legend of f_1 and f_{24}) versus dimension in log-log presentation. The $\text{ERT}(\Delta f)$ equals to $\#\text{FEs}(\Delta f)$ divided by the number of successful trials, where a trial is successful if $f_{\text{opt}} + \Delta f$ was surpassed during the trial. The $\#\text{FEs}(\Delta f)$ are the total number of function evaluations while $f_{\text{opt}} + \Delta f$ was not surpassed during the trial from all respective trials (successful and unsuccessful), and f_{opt} denotes the optimal function value. Crosses (\times) indicate the total number of function evaluations $\#\text{FEs}(-\infty)$. Numbers above ERT-symbols indicate the number of successful trials. Annotated numbers on the ordinate are decimal logarithms. Additional grid lines show linear and quadratic scaling.

f_1 in 5-D, N=15, mFE=1198					f_1 in 20-D, N=15, mFE=13503					f_2 in 5-D, N=15, mFE=2206					f_2 in 20-D, N=15, mFE=22791						
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	
10	15	2.7e1	2.2e1	3.3e1	2.7e1	15	1.2e3	1.1e3	1.2e3	1.2e3	15	5.2e2	4.7e2	5.8e2	5.2e2	15	8.6e3	8.3e3	8.8e3	8.6e3	
1	15	1.2e2	1.1e2	1.3e2	1.2e2	15	2.5e3	2.4e3	2.5e3	2.5e3	10	7.1e2	6.5e2	7.7e2	7.1e2	15	1.0e4	1.0e4	1.1e4	1.0e4	
1e-1	15	2.3e2	2.2e2	2.5e2	2.3e2	15	3.8e3	3.8e3	3.8e3	3.8e3	1e-1	15	8.8e2	8.1e2	9.5e2	8.8e2	15	1.2e4	1.1e4	1.2e4	1.2e4
1e-3	15	4.4e2	4.2e2	4.7e2	4.4e2	15	6.3e3	6.2e3	6.3e3	6.3e3	1e-3	15	1.2e3	1.1e3	1.2e3	1.2e3	15	1.4e4	1.4e4	1.4e4	1.4e4
1e-5	15	6.8e2	6.5e2	7.0e2	6.8e2	15	8.8e3	8.8e3	8.9e3	8.8e3	1e-5	15	1.4e3	1.3e3	1.5e3	1.4e3	15	1.7e4	1.6e4	1.7e4	1.7e4
1e-8	15	1.0e3	9.7e2	1.0e3	1.0e3	15	1.3e4	1.3e4	1.3e4	1.3e4	1e-8	15	1.7e3	1.6e3	1.8e3	1.7e3	15	2.1e4	2.0e4	2.1e4	2.1e4
f_3 in 5-D, N=15, mFE=2.13e6					f_3 in 20-D, N=15, mFE=2.00e7					f_4 in 5-D, N=15, mFE=5.01e6					f_4 in 20-D, N=15, mFE=2.00e7						
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	
10	15	9.8e2	4.4e2	1.5e3	9.8e2	15	1.9e5	1.5e5	2.3e5	1.9e5	10	15	3.1e3	1.9e3	4.4e3	3.1e3	0	<i>13e+0</i>	<i>12e+0</i>	<i>15e+0</i>	7.9e6
1	15	5.4e4	3.6e4	7.3e4	5.4e4	2	1.3e8	7.1e7	>3e8	2.0e7	1	2	3.4e7	1.7e7	>7e7	2.8e6					
1e-1	15	2.9e5	1.5e5	4.5e5	2.9e5	0	<i>20e-1</i>	<i>99e-2</i>	<i>40e-1</i>	6.3e6	1e-1	0	<i>20e-1</i>	<i>99e-2</i>	<i>20e-1</i>	7.1e5					
1e-3	15	3.0e5	1.5e5	4.6e5	3.0e5						1e-3										
1e-5	15	3.1e5	1.6e5	4.7e5	3.1e5						1e-5										
1e-8	15	3.1e5	1.6e5	4.9e5	3.1e5						1e-8										
f_5 in 5-D, N=15, mFE=232					f_5 in 20-D, N=15, mFE=689					f_6 in 5-D, N=15, mFE=4600					f_6 in 20-D, N=15, mFE=68285						
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	
10	15	7.1e1	6.2e1	7.9e1	7.1e1	15	4.0e2	3.7e2	4.3e2	4.0e2	10	15	2.4e2	2.1e2	2.7e2	2.4e2	15	7.0e3	6.8e3	7.2e3	7.0e3
1	15	1.1e2	9.4e1	1.2e2	1.1e2	15	4.6e2	4.3e2	4.8e2	4.6e2	1	15	5.0e2	4.5e2	5.5e2	5.0e2	15	1.2e4	1.2e4	1.2e4	1.2e4
1e-1	15	1.2e2	1.0e2	1.3e2	1.2e2	15	4.6e2	4.4e2	4.8e2	4.6e2	1e-1	15	8.9e2	7.9e2	9.9e2	8.9e2	15	1.8e4	1.7e4	1.8e4	1.8e4
1e-3	15	1.2e2	1.0e2	1.3e2	1.2e2	15	4.6e2	4.4e2	4.8e2	4.6e2	1e-3	15	1.6e3	1.5e3	1.7e3	1.6e3	15	2.9e4	2.7e4	3.0e4	2.9e4
1e-5	15	1.2e2	1.0e2	1.3e2	1.2e2	15	4.6e2	4.4e2	4.8e2	4.6e2	1e-5	15	2.3e3	2.1e3	2.4e3	2.3e3	15	4.0e4	3.8e4	4.2e4	4.0e4
1e-8	15	1.2e2	1.0e2	1.3e2	1.2e2	15	4.6e2	4.4e2	4.8e2	4.6e2	1e-8	15	3.5e3	3.3e3	3.6e3	3.5e3	15	5.6e4	5.4e4	5.8e4	5.6e4
f_7 in 5-D, N=15, mFE=20656					f_7 in 20-D, N=15, mFE=1517					f_8 in 5-D, N=15, mFE=11429					f_8 in 20-D, N=15, mFE=116850						
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	
10	15	8.9e1	7.6e1	1.0e2	8.9e1	15	2.3e3	2.3e3	2.4e3	2.3e3	10	15	2.5e2	2.3e2	2.7e2	2.5e2	15	1.7e4	1.6e4	1.7e4	1.7e4
1	15	6.3e2	2.6e2	1.0e3	6.3e2	15	4.3e3	4.0e3	4.6e3	4.3e3	1	15	2.1e3	1.1e3	3.0e3	2.1e3	15	3.4e4	2.9e4	4.0e4	3.4e4
1e-1	15	3.5e3	2.7e3	4.3e3	3.5e3	15	9.5e3	6.8e3	1.2e4	9.5e3	1e-1	15	2.6e3	1.6e3	3.6e3	2.6e3	15	3.8e4	3.3e4	4.4e4	3.8e4
1e-3	15	5.8e3	4.5e3	7.5e3	5.8e3	15	2.2e4	1.5e4	2.9e4	2.2e4	1e-3	15	3.0e3	2.1e3	4.0e3	3.0e3	15	4.2e4	3.7e4	4.7e4	4.2e4
1e-5	15	5.8e3	4.5e3	7.5e3	5.8e3	15	2.2e4	1.5e4	2.9e4	2.2e4	1e-5	15	3.3e3	2.4e3	4.4e3	3.3e3	15	4.5e4	3.9e4	5.0e4	4.5e4
1e-8	15	6.1e3	4.7e3	7.8e3	6.1e3	15	2.3e4	1.6e4	3.0e4	2.3e4	1e-8	15	3.7e3	2.7e3	4.7e3	3.7e3	15	4.9e4	4.3e4	5.4e4	4.9e4
f_9 in 5-D, N=15, mFE=13398					f_9 in 20-D, N=15, mFE=121138					f_{10} in 5-D, N=15, mFE=2458					f_{10} in 20-D, N=15, mFE=23694						
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	
10	15	2.4e2	2.3e2	2.6e2	2.4e2	15	1.7e4	1.6e4	1.8e4	1.7e4	10	15	6.3e2	5.6e2	7.1e2	6.3e2	15	9.5e3	9.1e3	9.8e3	9.5e3
1	15	2.8e3	1.7e3	3.9e3	2.8e3	15	3.4e4	2.9e4	4.0e4	3.4e4	1	15	8.1e2	7.2e2	9.1e2	8.1e2	15	1.1e4	1.1e4	1.1e4	1.1e4
1e-1	15	3.3e3	2.1e3	4.5e3	3.3e3	15	3.8e4	3.3e4	4.4e4	3.8e4	1e-1	15	1.0e3	8.9e2	1.1e3	1.0e3	15	1.2e4	1.2e4	1.3e4	1.2e4
1e-3	15	3.8e3	2.6e3	5.1e3	3.8e3	15	4.2e4	3.7e4	4.8e4	4.2e4	1e-3	15	1.3e3	1.2e3	1.4e3	1.3e3	15	1.5e4	1.5e4	1.5e4	1.5e4
1e-5	15	4.1e3	2.9e3	5.4e3	4.1e3	15	4.5e4	3.9e4	5.1e4	4.5e4	1e-5	15	1.5e3	1.4e3	1.6e3	1.5e3	15	1.8e4	1.7e4	1.8e4	1.8e4
1e-8	15	4.5e3	3.3e3	5.8e3	4.5e3	15	4.9e4	4.3e4	5.5e4	4.9e4	1e-8	15	1.9e3	1.8e3	2.0e3	1.9e3	15	2.1e4	2.1e4	2.2e4	2.1e4
f_{11} in 5-D, N=15, mFE=2248					f_{11} in 20-D, N=15, mFE=18491					f_{12} in 5-D, N=15, mFE=7372					f_{12} in 20-D, N=15, mFE=67898						
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	
10	15	4.9e2	4.3e2	5.5e2	4.9e2	15	4.4e3	4.2e3	4.7e3	4.4e3	10	15	7.5e2	6.6e2	8.4e2	7.5e2	15	9.0e3	8.9e3	9.2e3	9.0e3
1	15	7.8e2	6.9e2	8.6e2	7.8e2	15	6.0e3	5.7e3	6.3e3	6.0e3	1	15	1.1e3	1.0e3	1.3e3	1.1e3	15	1.3e4	1.0e4	1.5e4	1.3e4
1e-1	15	9.4e2	8.6e2	1.0e3	9.4e2	15	7.2e3	6.9e3	7.5e3	7.2e3	1e-1	15	1.5e3	1.3e3	1.6e3	1.5e3	15	2.0e4	1.7e4	2.2e4	2.0e4
1e-3	15	1.2e3	1.1e3	1.3e3	1.2e3	15	9.8e3	9.5e3	1.0e4	9.8e3	1e-3	15	2.2e3	1.9e3	2.4e3	2.2e3	15	3.2e4	3.0e4	3.5e4	3.2e4
1e-5	15	1.5e3	1.4e3	1.6e3	1.5e3	15	1.2e4	1.2e4	1.3e4	1.2e4	1e-5	15	2.8e3	2.4e3	3.1e3	2.8e3	15	4.4e4	4.2e4	4.7e4	4.4e4
1e-8	15	1.8e3	1.7e3	1.9e3	1.8e3	15	1.6e4	1.6e4	1.6e4	1.6e4	1e-8	15	3.6e3	3.2e3	4.1e3	3.6e3	15	5.7e4	5.5e4	5.9e4	5.7e4
f_{13} in 5-D, N=15, mFE=3025					f_{13} in 20-D, N=15, mFE=40034					f_{14} in 5-D, N=15, mFE=1975					f_{14} in 20-D, N=15, mFE=20254						
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	
10	15	3.5e2	3.3e2	3.6e2	3.5e2	15	5.6e3	5.6e3	5.7e3	5.6e3	10	15	1.3e1	9.9e0	1.7e1	1.3e1	15	7.9e2	7.5e2	8.4e2	7.9e2
1	15	5.8e2	5.5e2	6.0e2	5.8e2	15	8.4e3	8.2e3	8.5e3	8.4e3	1	15	1.1e2	9.7e1	1.2e2	1.1e2	15	2.2e3	2.2e3	2.3e3	2.2e3
1e-1	15	8.1e2	7.8e2	8.5e2	8.1e2	15	1.2e4	1.1e4	1.2e4	1.2e4	1e-1	15	2.5e2	2.3e2	2.7e2	2.5e2	15	3.8e3	3.8e3	3.9e3	3.8e3
1e-3	15	1.3e3	1.2e3	1.4e3	1.3e3	15	1.9e4	1.8e4	2.0e4	1.9e4	1e-3	15	5.7e2	5.4e2	6.1e2	5.7e2	15	7.3e3	7.2e3	7.4e3	7.3e3
1e-5	15	1.8e3	1.7e3	1.8e3	1.8e3	15	2.4e4	2.4e4	2.5e4	2.4e4	1e-5	15	9.1e2	8.6e2	9.6e2	9.1e2	15	1.1e4	1.1e4	1.2e4	1.1e4
1e-8	15	2.5e3	2.4e3	2.6e3	2.5e3	15	3.3e4	3.2e4	3.4e4	3.3e4	1e-8	15	1.5e3	1.4e3	1.6e3	1.5e3	15	1.8e4	1.7e4	1.8e4	1.8e4
f_{15} in 5-D, N=15, mFE=485369					f_{15} in 20-D, N=15, mFE=4.85e6					f_{16} in 5-D, N=15, mFE=223478					f_{16} in 20-D, N=15, mFE=1.58e7						
Δf	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	#	ERT	10%	90%	RT _{succ}	
10	15	5.1e2	4.6e2	5.6e2	5.1e2	15	8.4e4	6.3e4	1.1e5	8.4e4	10	15	2.3e2	1.7e2							

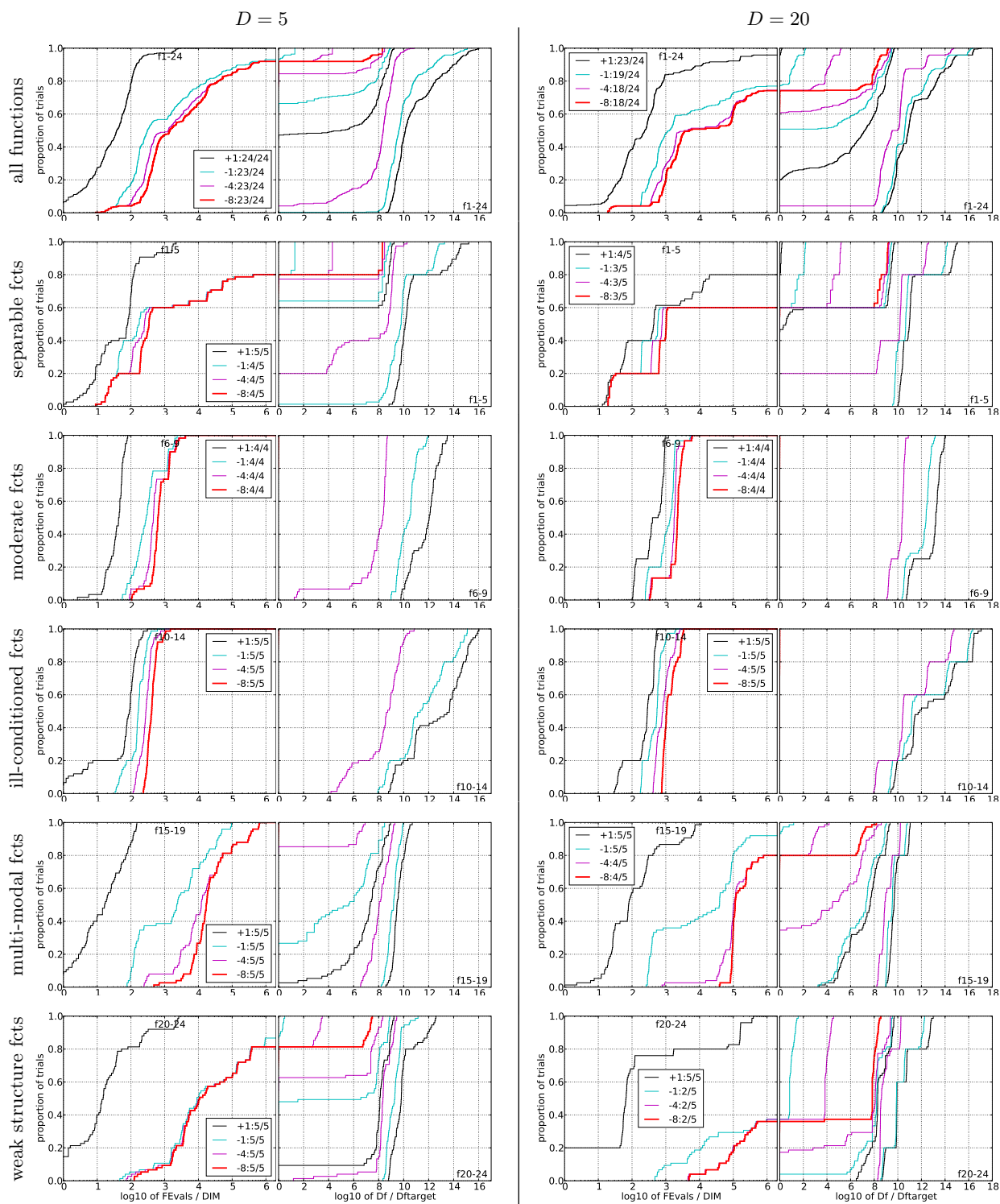


Figure 4: iAMaLGaM: Empirical cumulative distribution functions (ECDFs), plotting the fraction of trials versus running time (left) or Δf . Left subplots: ECDF of the running time (number of function evaluations), divided by search space dimension D , to fall below $f_{\text{opt}} + \Delta f$ with $\Delta f = 10^k$, where k is the first value in the legend. Right subplots: ECDF of the best achieved Δf divided by 10^k (upper left lines in continuation of the left subplot), and best achieved Δf divided by 10^{-8} for running times of $D, 10D, 100D \dots$ function evaluations (from right to left cycling black-cyan-magenta). Top row: all results from all functions; second row: separable functions; third row: misc. moderate functions; fourth row: ill-conditioned functions; fifth row: multi-modal functions with adequate structure; last row: multi-modal functions with weak structure. The legends indicate the number of functions that were solved in at least one trial. FEvals denotes number of function evaluations, D and DIM denote search space dimension, and Δf and Df denote the difference to the optimal function value.