



## Assignment of master's thesis

<b>Title:</b>	Accelerating evolutionary algorithms by means of neural networks
<b>Student:</b>	Bc. Jaroslav Langer
<b>Supervisor:</b>	prof. Ing. RNDr. Martin Holeňa, CSc.
<b>Study program:</b>	Informatics
<b>Branch / specialization:</b>	Knowledge Engineering
<b>Department:</b>	Department of Applied Mathematics
<b>Validity:</b>	until the end of summer semester 2024/2025

### Instructions

1. Familiarize yourself with the most successful method for continuous black-box optimization, the CMA-ES evolutionary algorithm (see Hansen, 2006).
2. Get acquainted with the surrogate modeling approach for black-box optimization, particularly regarding the CMA-ES algorithm (see, for example, Bajer et al., 2019).
3. Learn about the Comparing Continuous Optimizers (COCO) tool, developed at Inria Saclay, which is used to test black-box optimization methods and their surrogate models.
4. Explore the following four types of modern artificial neural networks:
  - a) Multi-layer perceptron with kernel parameterization (see, for example, Paria et al., 2022)
  - b) Variational autoencoder (see, for example, Kim et al., 2023)
  - c) Generative adversarial networks (see, for example, Lu et al., 2022)
  - d) Networks for learning the prior distribution of data (see Müller et al., 2023)
5. Based on the knowledge gained in points 1-4, your professional interests, and considering the implementation challenges of integrating CMA-ES with the selected network, select one or two of these networks to evaluate their effectiveness as surrogate models for CMA-ES.
6. Using a suitable implementation of CMA-ES, implement a new version of CMA-ES that uses your chosen network(s) as surrogate model(s).
7. Test the implemented version of CMA-ES using the COCO tool, which will compare it with several dozen other versions of CMA-ES, including those using surrogate models (and, if you are interested, with other methods for black-box optimization).



## References

- > L. Bajer, Z. Pitra, J. Repický, and M. Holeňa. Gaussian Process Surrogate Models for the CMA Evolution Strategy. *Evolutionary Computation*, 27:665–697, 2019.
- > N. Hansen. The CMA Evolution Strategy: A Comparing Review. In *Towards a New Evolutionary Computation*, pages 75–102. Springer, 2006.
- > S. Kim, PY. Lu, C. Lob, J. Smith, J. Snoek, et al. Deep Learning for Bayesian Optimization of Scientific Problems with High-Dimensional Structure. *Transactions on Machine Learning Research*, 1:openreview tPMQ6Je2rB, 2023.
- > M. Lu, S. Ning, S. Liu, F. Sun, B. Zhang, et al. OPT-GAN: A Broad-Spectrum Global Optimizer for Black-box Problems by Learning Distribution. *Arxiv 2102.03888v5*, 2022.
- > S. Müller, M. Feurer, N. Hollmann, F. Hutter. PFNs4BO: In-Context Learning for Bayesian Optimization. *ICML*, 2023.
- > B. Paria, B. Póczos, K. Ravikumar, J. Schneider, A.S. Suggala, et al. Be Greedy – a Simple Algorithm for Blackbox Optimization using Neural Networks. In *ICML Workshop on Adaptive Experimental Design and Active Learning in the Real World*, pages 1–27, 2022.

Master's thesis

**ACCELERATING  
EVOLUTIONARY  
ALGORITHMS BY  
MEANS OF NEURAL  
NETWORKS**

**Bc. Jaroslav Langer**

Faculty of Information Technology  
Department of Applied Mathematics  
Supervisor: prof. Ing. RNDr. Martin Holeňa, CSc.  
May 9, 2024

Czech Technical University in Prague  
Faculty of Information Technology

© 2024 Bc. Jaroslav Langer. All rights reserved.

*This thesis is school work as defined by Copyright Act of the Czech Republic. It has been submitted at Czech Technical University in Prague, Faculty of Information Technology. The thesis is protected by the Copyright Act and its usage without author's permission is prohibited (with exceptions defined by the Copyright Act).*

Citation of this thesis: Langer Jaroslav. *Accelerating Evolutionary Algorithms by Means of Neural Networks*. Master's thesis. Czech Technical University in Prague, Faculty of Information Technology, 2024.

## Contents

<b>Acknowledgments</b>	<b>vi</b>
<b>Declaration</b>	<b>vii</b>
<b>Abstract</b>	<b>viii</b>
<b>Introduction</b>	<b>1</b>
0.1 Derivative-free Optimization . . . . .	1
0.1.1 Evolutionary Algorithms (EAs) . . . . .	1
0.2 CMA-ES . . . . .	2
0.3 Surrogate Modeling . . . . .	3
0.4 Benchmarking . . . . .	4
0.4.1 BBOB Noiseless Functions . . . . .	4
0.4.2 Performance Measurement . . . . .	6
<b>1 Related Work</b>	<b>8</b>
1.1 Neural Network Assisted CMA-ES . . . . .	8
1.2 DTS-CMA-ES . . . . .	9
1.3 lq-CMA-ES . . . . .	9
1.4 Surrogate Model - Evolution Control Interplay . . . . .	10
<b>2 Choosing the Approach</b>	<b>11</b>
2.1 Proposed Neural Networks . . . . .	11
2.1.1 VAE . . . . .	11
2.1.2 GAN . . . . .	11
2.1.3 PFNs . . . . .	12
2.1.4 NTK . . . . .	12
2.2 Findings, Motives, and Decisions . . . . .	12
2.2.1 Ensemble of Activation Functions . . . . .	14
2.2.2 Anchored Ensemble and RAF . . . . .	15
<b>3 Methodology</b>	<b>16</b>
3.1 CMA-ES Variant . . . . .	17
3.2 Surrogate . . . . .	17
3.2.1 Subset Method . . . . .	17
3.2.2 Feature Transformation . . . . .	18
3.2.3 Target Variable Transformation . . . . .	18

3.2.4	Weighing the Train Points . . . . .	19
3.2.5	The Model . . . . .	19
3.3	Evolution Control . . . . .	20
3.4	Implementation . . . . .	20
3.5	Benchmarking Platforms . . . . .	22
<b>4</b>	<b>Results</b>	<b>23</b>
4.1	Preliminary - Dimension 2 . . . . .	24
4.2	Benchmarking RAF Dimensions 2,3,5,10,20 . . . . .	31
<b>5</b>	<b>Conclusions</b>	<b>33</b>
5.1	Future Work . . . . .	33
<b>A</b>	<b>Test 1 and 2 ERT Multiples for BBOB2009 Targets</b>	<b>35</b>
	List of abbreviations	46
	Attachment Contents	47

## List of Figures

1	CMA-ES Step Example . . . . .	2
2	Separable Functions . . . . .	5
3	Functions with Low or Moderate Conditioning . . . . .	5
4	Unimodal Functions with High Conditioning . . . . .	6
5	Multi-modal Functions with Global Structure . . . . .	6
6	Multi-modal Functions with Weak Global Structure . . . . .	6
2.1	Example of NTK-GP Standard Deviations over Sphere Optimization . . . . .	13
2.2	Example of sine modeling and AF effect . . . . .	14
3.1	Example of Search Points and The Archive . . . . .	18
3.2	Example Correlations of NNs with different AF . . . . .	19
3.3	Example Predictions of NNs With Different AFs . . . . .	20
4.1	Runtime ECDF for Function Groups Dim 2 Preliminary . . . . .	28
4.2	ECDF Functions 1-14 Budget Based Targets BBOB2009 . . . . .	29
4.3	ECDF Functions 15-24 Budget Based Targets BBOB2009 . . . . .	30
4.4	Runtime ECDF for Function Groups Dim 2 . . . . .	31
4.5	Expected running time (ERT) divided by dimension versus dimension in log-log presentation . . . . .	32

## List of Tables

4.1	Experiment $10^{-8}$ -target hits . . . . .	25
4.2	Experiment 50 E/D-target hits . . . . .	27
A.1	ERT Multiples for Budget Based Targets BBOB2009 Dimension 2 . . . . .	36
A.2	ERT Multiples for Budget Based Targets BBOB2009 Dimension 3 . . . . .	37
A.3	ERT Multiples for Budget Based Targets BBOB2009 Dimension 5 . . . . .	38
A.4	ERT Multiples for Budget Based Targets BBOB2009 Dimension 10 . . . . .	39

A.5 ERT Multiples for Budget Based Targets BBOB2009 Dimension	
20 . . . . .	40

## List of code listings

3.1 Function <code>seek_minimum</code> Signature . . . . .	20
3.2 SurrogateCallable Protocol . . . . .	21
3.3 EvolutionControl Protocol . . . . .	21



*Computational resources were provided by the e-INFRA CZ project (ID:90254), supported by the Ministry of Education, Youth and Sports of the Czech Republic.*

*I want to thank my supervisor, Prof. Ing. RNDr. Martin Holeňa, CSc.*

*I also want to thank my family: Iva and Jiří & Jaroslav and Nikol & Karolína and Jaroslav, as well as Miloslava and Josef & Emílie and Jaroslav.*

*I am also grateful to my friends, whom I won't list exhaustively due to space constraints, but the first few who come to mind are: Kamil, Radomil, Jan, Jan, Filip, Ondřej & Daniel, Adam, Vojtěch & Matúš, Barbora, Valeriia, Dominika, Lucie, Lukáš & David and Petr & Jana and Kateřina & individually Natálie, Martina, and Martin.*

*A great source of motivation came from the group of people around the university sports volleyball and yoga so I thank Jana, Klára, Antonín, Radim, Lenka, Radka, Lucie, Michal, Štěpán, René, and Matěj.*

*I would like to thank all the teachers who were accommodating in accepting my assignments after the deadline with some penalties, or in allowing me to take examinations on the last days of the examination period, sometimes even on the very last day. Their support made my studies much more manageable and less painful.*

*Special thanks go to Kamil Dedecius for showing that teaching can come from within and that a teacher can truly care about students despite all the duties the teacher might have.*

*Another thanks goes to Lukáš Krejčí for being very tolerant when negotiating the allocation of my time until the end of my studies.*

*The most special thanks go to my beloved Bedřiška.*

## Declaration

I hereby declare that the presented thesis is my own work and that I have cited all sources of information in accordance with the Guideline for adhering to ethical principles when elaborating an academic final thesis.

I acknowledge that my thesis is subject to the rights and obligations stipulated by the Act No. 121/2000 Coll., the Copyright Act, as amended. In accordance with Section 2373(2) of Act No. 89/2012 Coll., the Civil Code, as amended, I hereby grant a non-exclusive authorization (licence) to utilize this thesis, including all computer programs that are part of it or attached to it and all documentation thereof (hereinafter collectively referred to as the "Work"), to any and all persons who wish to use the Work. Such persons are entitled to use the Work in any manner that does not diminish the value of the Work and for any purpose (including use for profit). This authorisation is unlimited in time, territory and quantity.

In Prague on May 9, 2024

## Abstract

The Covariance Matrix Adaptation Evolution Strategy (CMA-ES) was combined with ensembles of networks with different activation functions as a surrogate model. Over a year of rigorous experimentation, these ensembles demonstrated an ability to estimate uncertainty. However, while they showed potential, they did not prove to be a universal solution, yielding mixed results across different functions and dimensions.

Significant focus was placed on dimension 2, where enhancements in performance were observed, albeit with varying success in other dimensions where results often deteriorated. This investigation posits that, with an investment of time comparable to that required for Gaussian Process (GP)-based surrogates, similar or better outcomes might be achieved. Yet, the practicality of such an investment is questioned, especially considering the existing efficacy of advanced methods like DTS-CMA-ES and lq-CMA-ES.

A principal contribution of this thesis is the development of a mini-framework for surrogate testing, designed to lower the barriers to entry for researchers and practitioners interested in applying surrogate models to CMA-ES.

**Keywords** derivative-free optimization, evolution strategy, ES, surrogate-modelling, NN, RAF, evolution control, EC

## Abstrakt

Metoda CMA-ES byla zkombinována s ansámby sítí s různými aktivačními funkcemi jako náhradní model. Během ročního důkladného experimentování tyto ansámby prokázaly schopnost odhadovat nejistotu. Ačkoli ukázaly potenciál, neprokázaly se jako univerzální řešení, přinášely smíšené výsledky napříč různými funkcemi a dimenzemi.

Značná pozornost byla věnována dimenzi 2, kde bylo pozorováno zlepšení výkonu, ačkoli s různým úspěchem v ostatních dimenzích, kde výsledky často zhoršovaly. Toto vyšetřování předpokládá, že s investicí času srovnatelnou

s tou, která je vyžadována pro náhradní modely založené na Gaussovských procesech (GP), by mohly být dosaženy podobné nebo lepší výsledky. Nicméně praktičnost takové investice je zpochybňována, zejména vzhledem k existující efektivitě pokročilých metod jako jsou DTS-CMA-ES a lq-CMA-ES.

Hlavním přínosem této práce je vývoj mini-frameworku pro testování náhradních modelů, který je navržen tak, aby snížil bariéry pro vstup pro výzkumníky a praktiky zajímající se o aplikaci náhradních modelů na CMA-ES.

**Klíčová slova** Optimalizace bez známých derivací, Evoluční strategie, ES, Náhradní modelování, NN, RAF, Kontrola evoluce, EC

# Introduction

*Often in life, we strive to make the best decisions but cannot foresee the outcomes. In mathematics, this is called derivative-free optimization.*

## 0.1 Derivative-free Optimization

Derivative-free optimization is a field where the only accessible information about the function being optimized is its values evaluated at given points. In this setting, there are no assumptions about the function's shape or structure, and there is no access to its gradients[1]. The discipline of derivative-free optimization has a rich scientific history. It's unsurprising, given its broad range of applications, spanning aerodynamics design [2], hyper-parameter optimization of machine learning algorithms, and materials science, such as alloy design and synthesis of short polymer fibers, as well as bio-engineering tasks like 3D bio-printing and molecule design[3]. The field is also known by other names, including gradient-free optimization, optimization without derivatives, or zeroth-order optimization; a prominent sub-field is black-box optimization[4].

► **Definiton 0.1** (Archive). *The archive is a collection of evaluated points and their function values. We can say the archive grows during the optimization.*

### 0.1.1 Evolutionary Algorithms (EAs)

Evolutionary Algorithms (EAs) are a prominent approach for derivative-free optimization, inspired by biological evolution. One of their strengths is the ability to deal with complex, multi-modal, and noisy functions. This robustness is rooted in directing the search by a population of candidate solutions[5].

EAs operate through cycles called generations. In each generation, a stochastic variation-often small-is introduced to the current population, followed by the selection of the fittest individuals[6]. Variation operators, such as mutation or crossover, help explore the search space while taking into ac-

count the information already gathered about the problem. Meanwhile, the selection process directs the search towards promising regions, effectively balancing the exploration-exploitation trade-off[5].

Despite these advantages, a major limitation of EAs is the high number of function evaluations typically required. This becomes prohibitive in scenarios where evaluations are expensive, a common case in real-world applications.

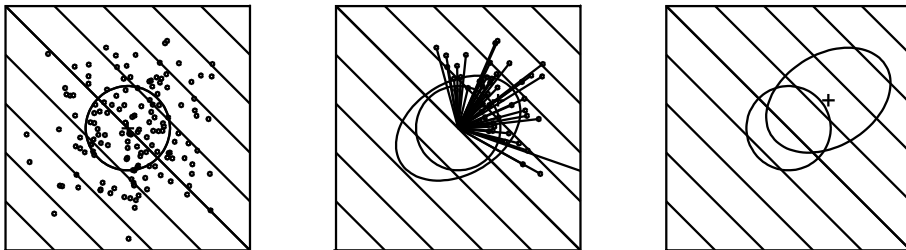
Evolution strategies (ES), as a subset of evolutionary algorithms, date back to the 1960s. They are best suited for black-box optimization in continuous search spaces. Algorithmically speaking, ES methods generate new candidate solutions stochastically, typically sampling from a multivariate normal probability distribution [6].

## 0.2 CMA-ES

Among evolution strategies, the Covariance Matrix Adaptation Evolution Strategy (CMA-ES) is considered one of the state-of-the-art methods for black-box optimization[7]. One of its design goals was to be able to solve ill-conditioned and non-separable problems[8]. In CMA-ES, a population of search points (individuals, offspring)  $x_1, \dots, x_\lambda$  is sampled from the multivariate normal distribution  $\mathcal{N}(m^{(g)}, \mathcal{C}^{(g)})$  where the mean and covariance matrix are adapted based on the function values of the search points.

The new mean,  $m^{(g+1)}$ , is calculated as the weighted average of the fittest subset of the population, with weights assigned such that better points receive larger weights.

Similarly, the new covariance matrix  $\mathcal{C}^{(g+1)}$  is estimated using a weighted subset of the population. Because estimating the covariance from a single population can be unreliable, especially for smaller populations, adaptation procedures: rank- $\mu$ -update and rank-one-update are employed[8].



■ **Figure 1** CMA-ES step example - Left: Search points sample, Middle: Weighing best points, Right: New mean and covariance estimates[8].

**Step-size  $\sigma^{(g)}$  Adaptation** While the Covariance Matrix Adaptation primarily influences the search direction, the step size  $\sigma^{(g)}$  controls the scale of the search. By default, CMA-ES performs a Cumulative Step-size Adaptation (CSA), which utilizes an evolution path, i.e., a cumulative sum of

successive steps. The length of this path is compared to its expected length under random selection, where each step is independent and uncorrelated. If the evolution path is longer than expected, indicating a consistent direction in the steps,  $\sigma$  is increased to allow for larger steps. Conversely, if the steps tend to cancel each other out, suggesting that the step size may be too large,  $\sigma$  is decreased [8]. An alternative is the Two-Point Step-size Adaptation (TPA) [9].

**Active CMA-ES** This variant, also a default setting, leverages information about unsuccessful offspring to actively reduce variances in the mutation distribution along less promising directions of the search space [10].

**Population Size  $\lambda$**  The default population size is set to  $4 + \lfloor 3 \ln n \rfloor$ . Increasing  $\lambda$  generally enhances the global search capabilities but at the cost of more function evaluations[11].

**IPOP-CMA-ES** This is a restart strategy in CMA-ES, where the population size is increased for each restart (IPOP), usually by a factor of 2. As noted above, this adjustment makes the search progressively more global after each restart[12].

**Boundary Handling** In situations where the solution is not expected to be close to the infeasible region, a straightforward approach is to set any value outside the boundaries to the maximum (reasonable) value or to resample the points outside the feasible region. However, these approaches do not work well when the solution lies close to the boundary. For such cases, techniques, termed "repair", exist to address this issue. It is important to note that naive solutions, such as clipping to the boundary, can lead to problems such as divergence or too fast convergence of the step-size. The default configuration for CMA-ES is to operate on the  $\mathbb{R}^D$ [8].

### 0.3 Surrogate Modeling

Due to the nature of evolutionary algorithms, they typically require a large number of function evaluations. However, in many real-world applications, evaluating the optimization function can be time-consuming or costly. In these situations, employing an approximate model, also known as a meta-model or surrogate, to estimate some function values and save on their evaluations can be of great help [13].

The term "surrogate model" is also used in Bayesian optimization (BO) with a slightly different meaning. Bayesian optimization is another recognized approach for derivative-free optimization, particularly well-suited for functions that are expensive to evaluate. It traditionally constructs a probabilistic surrogate model, which places a prior on the function being optimized and continuously updates this model based on new evaluations. The decision of which

points to evaluate next is made by an acquisition function, such as Expected Improvement, that utilizes the posterior distribution modeled by the surrogate[3]. In simple terms, the acquisition function chooses points with a high probability of finding a new optimum.

In this work, the Evolution Strategy (ES) requires values for all search points, not just the most promising ones based on the model. Therefore, the surrogate’s role is not only to identify the most promising points but to estimate all the values, effectively reducing the number of evaluations required. This approach is intuitively based on the premise that the function values of many points can be sufficiently approximated using the data already observed.

## 0.4 Benchmarking

Quantifying and comparing performance is a vital aspect of research on optimization algorithms. Several benchmark suites are available, some of which are listed here[14]. This work specifically targets the Noiseless Functions from the Black-box Optimization Benchmarking (BBOB) test suite, which comprises 24 noise-free real-parameter single-objective functions [15].

### 0.4.1 BBOB Noiseless Functions

According to the authors, the functions were selected to assess the performance of optimization algorithms against typical challenges encountered in continuous domain search. The goal was for the collection to represent a balanced range of difficulties, emphasizing characteristics known to pose challenges in practice, such as ill-conditioning, non-separability, non-convexity, and ruggedness. Easier functions were considered of lesser interest [15, 16, 14].

**Dimensionality  $D$**  All functions are arbitrarily scalable in dimension[17].

Throughout this work,  $D$  without further specification refers to dimension.

**Domain** Each function is defined and evaluable over  $\mathbb{R}^D$ , with the actual search domain specified as  $[-5, 5^D]$ , within which global optimum is always located[18].

**Optimum** For the majority of functions, the global optimum is uniformly distributed over  $[-4, 4]^D$ . The optimum’s value is drawn from a Cauchy distribution, with a zero median and about 50% of the values lying between  $-100$  and  $100$ [18].

**Instances** For most functions, neither the location of the optimum nor the function value is fixed. This variability allows for the generation of instances (potentially random) for each function by shifting both the optimum’s location and value. The underlying assumption for subsequent analysis is that different instances of the same test function present similar



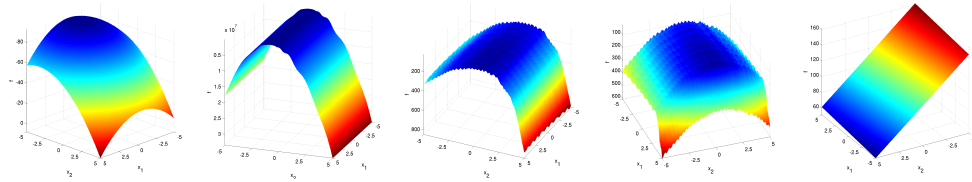
difficulties. This concept of instances enables a fair comparison between deterministic and stochastic solvers and prevents solvers from exploiting specific properties of a particular function instance, such as the position of the optimum[14].

**Problem** A specific instance of a given function is often referred to as a "problem" to distinguish it from the generic function type. For example, a Sphere function may appear in multiple dimensions and instances, whereas a specific Sphere problem can be evaluated or optimized by a solver. The default problems in the BBOB test suite span dimensions (2, 3, 5, 10, 20, 40) and comprise 15 instances each[14].

**Conditioning** Generally, a function is considered ill-conditioned if, for points with similar function values, the minimal displacement in the search space required to achieve a given improvement in function value differs by orders of magnitude[19].

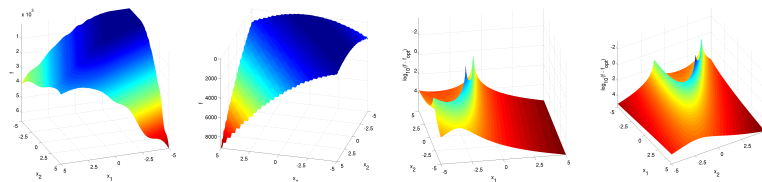
The noiseless functions are divided into 5 groups based on separability, conditioning, modality, and global structure.

### Separable Functions



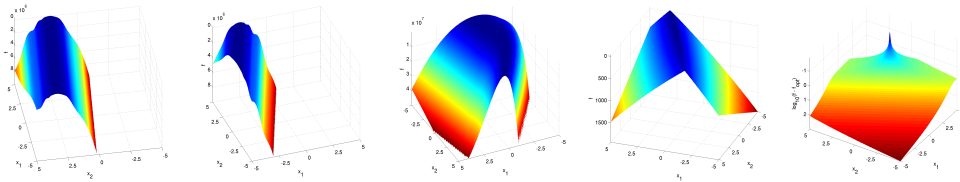
■ **Figure 2** Functions (from the left) 1. Sphere, 2. Ellipsoidal, 3. Rastrigin, 4. Büche-Rastrigin, 5. Linear Slope[18].

### Functions with Low or Moderate Conditioning



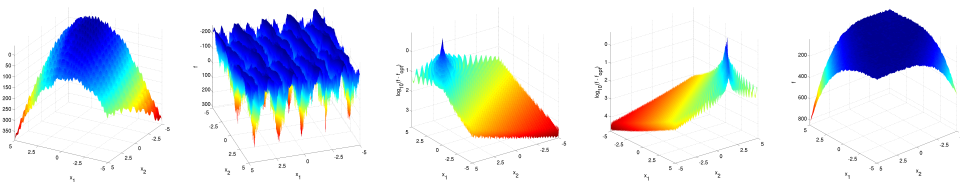
■ **Figure 3** Functions (from the left) 6. Attractive Sector, 7. Step Ellipsoidal, 8. Rosenbrock, 9. Rosenbrock rotated[18].

### Functions with High Conditioning and Unimodal



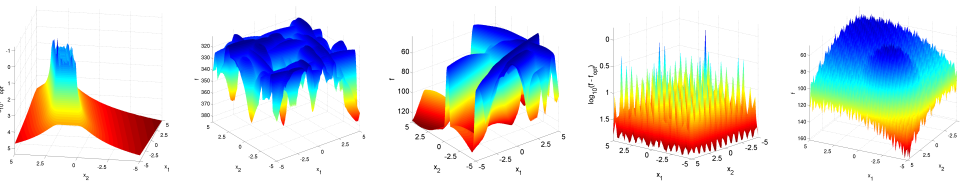
■ **Figure 4** Functions (from the left) 10. Ellipsoidal, 11. Discus, 12. Bent Cigar, 13. Sharp Ridge, 14. Different Powers[18].

### Multi-modal Functions with Adequate Global Structure



■ **Figure 5** Functions (from the left) 15. Rastrigin, 16. Weierstrass, 17. Schaffers F7, 18. Schaffers F7 ill-conditioned, 19. Composite Griewank-Rosenbrock F8F2[18].

### Multi-modal Functions with Weak Global Structure



■ **Figure 6** Functions (from the left) 20. Schwefel, 21. Gallagher's Gaussian 101-me Peaks, 22. Gallagher's Gaussian 21-hi Peaks, 23. Katsuura, 24. Lunacek bi-Rastrigin[18].

## 0.4.2 Performance Measurement

**Budget** The budget refers to the allowed number of function evaluations, typically specified as a multiple of the dimension. There is no inherent predefined budget (see Anytime Assessment Approach). As a result, all results are comparable up to the smallest chosen budget in the solver set.

The larger the budget, the more data are generated for comparison with other algorithms[14].

**Runtime** The runtime of a solver for a target is the number of evaluations performed until the target value is reached or surpassed for the first time[14].

**Expected Runtime (ERT)** Expected Runtime (ERT), or Averaging Runtime for a target, is calculated from a set of trials as “the sum of all evaluations in unsuccessful trials plus the sum of the runtimes in all successful trials, both divided by the number of successful trials”[20].

**Target** A target is a specified function value. Solver performance is measured by the runtimes of these targets. Targets can be defined based on distances to the optimal value, which can be considered precisions, e.g., a problem was solved with a precision of  $10^{-5}$ .

**Anytime Assessment Approach** The performance is measured over the whole run of the solver, not just after reaching a specific runtime, budget, or target [14].

**Budget-based Target Values** In addition to using predefined precisions, targets can be defined based on budgets and a finite set of solvers. For any given budget, the easiest target is chosen such that the expected runtime of every solver for that target exceeds the specified budget. These budget-based or run-length-based targets depend on the solver set and are specific for each function and dimension.

## Related Work

*There is nothing wrong with reinventing the wheel from time to time, but not every day on your way to work.*

Surrogate modeling has proven itself as a valuable improvement to the already powerful CMA-ES algorithm [5]. In this chapter, I first briefly describe the most recent attempt known to me to enhance CMA-ES with Neural Networks. Subsequently, I discuss two state-of-the-art surrogate methods: DTS-CMA-ES and lq-CMA-ES. I present an overview of their main components.

### 1.1 Neural Network Assisted CMA-ES

“Over the past decade, neural networks have emerged as a potential candidate to play the role of surrogate model in black-box optimization”[3].

Multi-layer perceptrons were extensively tested as surrogate models[21]. However, specifically with CMA-ES, to the best of my knowledge, it was only done by Yaochu et al.[2, 22].

“it allows us to take advantage of the self-adaptation of the covariance matrix in online training of the neural network model. The basic idea is that the new data points that lie along the direction in which the evolutionary algorithm proceeds should be given larger weight in online learning”[22]

The MLP network used has one hidden layer with 20 nodes. The evolution control is generation-based, where the ratio of using approximated values is proportional to the mean current error[22].

“The most important question is how many individuals or how many generations should be controlled to guarantee the correct convergence of an evolutionary algorithm when false optima are present in the approximate fitness function. To answer this question, simulations have been carried out on the Ackley function [35] and the Rosenbrock function [36], which have been studied widely in the field of evolutionary computation. Based on these empirical stud-

ies, a framework for model management in generation-based evolution control is proposed.”[22]

## 1.2 DTS-CMA-ES

Is an approach where the surrogate model is a Gaussian Process. In each generation, a GP is trained on a transformed subset of the Archive. The current evaluation ratio  $\alpha^{(g)}$  of points with the greatest acquisition function values are evaluated, and the second model is trained with these additional evaluated points. The predictions of the second model are used as surrogate values, offset by the so far best optimum[23].

- ▶ **Definiton 1.1** (Double Population Bigger Sigma CMA-ES). *In DTS they used initial step-size  $\sigma^{(0)} = \frac{8}{3}$  and initial population size  $\lambda = 8 + \lceil 6 \ln D \rceil$ [23] which (is up to +1) double the CMA-ES default  $4 + \lceil 3 \ln n \rceil$ [8].*
- ▶ **Definiton 1.2** (Closest to Each Test Point Bounded by  $r_{max}^A$  and  $N_{max}$  Subset). *This method maximizes the equality of the number of closest (in Mahalanobis distance) train points for each test point, regarding the maximum distance  $r_{max}^A$  and maximum number of train points allowed[23].*
- ▶ **Definiton 1.3** (CMA-ES Distribution Feature Transformation). *The train points are standardized with regards to the current CMA-ES mean and standard deviation corrected with the step-size[23].*
- ▶ **Definiton 1.4** (Subset Standardization Target Transformation). *The train point values are standardized with mean and standard deviation estimated from the train subset[23].*
- ▶ **Definiton 1.5** (Double Training). *The model is trained twice in the generation, first on the train subset, and second time on the train subset plus some search points that were already evaluated[23].*
- ▶ **Definiton 1.6** (RDE Controlled Ratio EC). *The points to evaluate are chosen as the  $\alpha^{(g)}$  ratio of the points with the greatest acquisition function values. The  $\alpha^{(g)}$  is driven by the smoothed error from the previous generation between the first and second model predictions[23].*

## 1.3 lq-CMA-ES

The surrogate model performs Linear to Full Quadratic Regression based on the number of Archive points available.

“lmm-CMA-ES [ 20 ] uses local meta-models to compute a different surrogate for each solution (offspring). Our approach is different in that i) we maintain only a global model on which all offspring are eventually evaluated;

ii) our regression weights are based on a distance in f-space (i.e. fitness) instead of a distance in x-space; iii) we compare to the true ranking (instead of rank changes from adding data) in order to decide when to accept the model; iv) the internal computational complexity is smaller by the order of  $\lambda/\log(\lambda)$ . Compared to both above approaches, we i) start already after three or four evaluations to build a model; ii) change the model complexity and estimate also linear and diagonal-quadratic models before have gathered  $1.1(n(n+3)/2+1)$  data samples to estimate a full quadratic model; iii) utilize the model optimum to inject in the population in the next iteration, and iv) use CMA-ES with active covariance matrix update”[24]

► **Definiton 1.7** (Last N Sorted Generation Subset). *The lq-CMA-ES subsets the Archive as the last  $n$  points, where  $n$  is driven by the model’s degrees of freedom. The points are sorted inside each generation, so the last point that remains from the previous generation is the one with the best function value compared to the rest of that generation*[24].

► **Definiton 1.8** (Weights 20 to 1). *The regression uses linear weights from 20 to 1, where the greatest weights go to the points at the end, i.e., with the best value*[24].

► **Definiton 1.9** (Train With New Data). *Every time when a prediction is needed, there is a check if the model was trained with the newest evaluated points; if not, then the regression weights are re-estimated*[24].

► **Definiton 1.10** (Smart Offset). *Calculate the difference between the smallest prediction and the point’s function value and offset all predictions by this difference*[24].

► **Definiton 1.11** (Evaluate Until Kendall Under Threshold). *This evaluation control evaluates points sorted from the best prediction to the worst. Calculates Kendall coefficient, if it is greater than 0.85, for the rest of the points the model value is offset by the "smart offset" and used. For details see*[24].

[24]

## 1.4 Surrogate Model - Evolution Control Interplay

In this paper[7], Pitra et al. showed that the three prominent surrogate models can be successfully decoupled from the evolution control they were originally published with.

# Choosing the Approach

In this chapter, I'm going to discuss proposed neural network architectures and their possible application in the domain of surrogate modeling. Also, I show the reasoning behind my decisions and the plans that are connected to them.

## 2.1 Proposed Neural Networks

### 2.1.1 VAE

Variational autoencoder is an NN architecture, built upon the autoencoder with restrictions on the latent space.

“VAE-GP uses a VAE trained ahead of time on an unlabelled dataset representative of  $\mathcal{X}$ . This allows us to encode complex input spaces, such as chemical molecules, into a continuous latent space over which conventional GP-based BO methods can be applied”[25].

This approach could be interesting for surrogate modeling since there were recent attempts to use ANN-GP[26, 27]; it would be interesting to see whether VAE-GP would do better. Especially when “To our knowledge, this is the first time that ANN+GP combinations have been investigated for possible application in surrogate modeling”[27]. I was, explicitly told not to use Gaussian processes, so I chose a different approach to pursue.

### 2.1.2 GAN

GAN showed itself to be a great tool for modeling the distribution of a potential function minimum[28]. It would be interesting to compare it with CMA-ES. However, after discussing it with my supervisor, we didn't see a straightforward way to use the GAN. Having the distribution would be great for choosing which points to evaluate, but the modeled values would still be needed.

### 2.1.3 PFNs

Prior-data Fitted Networks (PFNs)[29] was shown to be a strong surrogate model in terms of Bayesian Optimization. Substituting Gaussian Processes, which would fit this domain neatly. However, I didn't choose this approach for these reasons: In the small dimensions, e.g. 2 with a budget multiplier of 250, at the end of the training, you have a maximum of 500 points. What is even worse, often some of them differ in the function value in the order of  $10^{-7}$  while others in the order of  $10^4$ , so using all the 500 points is optimistic at best. So, my premise was that the model should be learned while the optimization can not be too complex. However, there was an interesting idea: to record data from multiple functions and instances and learn the prior on this data. Unfortunately, I didn't see it doable in the scope of my master's thesis, so I wanted to choose something where my chances of moving forward were bigger.

### 2.1.4 NTK

Currently, the research of correspondence between NN in the infinite width limit and GP induced by Neural Tangent Kernel (NTK) is making a lot of progress, and it is a hot research topic[30, 31]. The advantage of this approach was the similarity to the Gaussian Processes already examined, so it could be a nice comparison of how well the wide neural network substitutes the Gaussian Processes. However, as I was told, this thesis should not be about Gaussian Processes as it's already been thoroughly studied by others, and it was recommended to me that trying to apply a method that should be equivalent, not better, and only in the infinite-width limit doesn't seem like a good idea to start with.

## 2.2 Findings, Motives, and Decisions

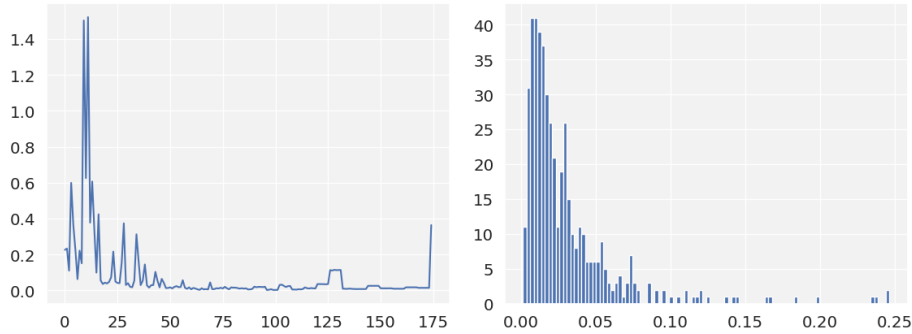
So far, I haven't chosen a way I would consider leading to the completion of these four. So, I started experimenting with plain NNs and Gaussian Processes with NTK parametrization as surrogate models to see where it would take me. At first, I realized that subsetting the train dataset was key for both plain NN and NTK. Standardization is expected by the GPs, but it also helped the NNs a lot.

Also, I started with the (second) easiest function, Sphere, in dimension 2, because the plain CMA-ES finds the optimum in about 200 evaluations  $\pm 100$ , so it was a quick iteration cycle to see if something works or not. In this phase, I realized I like the uncertainty predictions given by the NTK and it looked as if because the values were standardized, I could have a rule as simple as predictions with  $std > 0.025$  or  $std > 0.05$  needs to be evaluated and the certainty of the prediction for the rest is good enough to use the model. Especially, when CMA-ES weights the points only based on the order, not



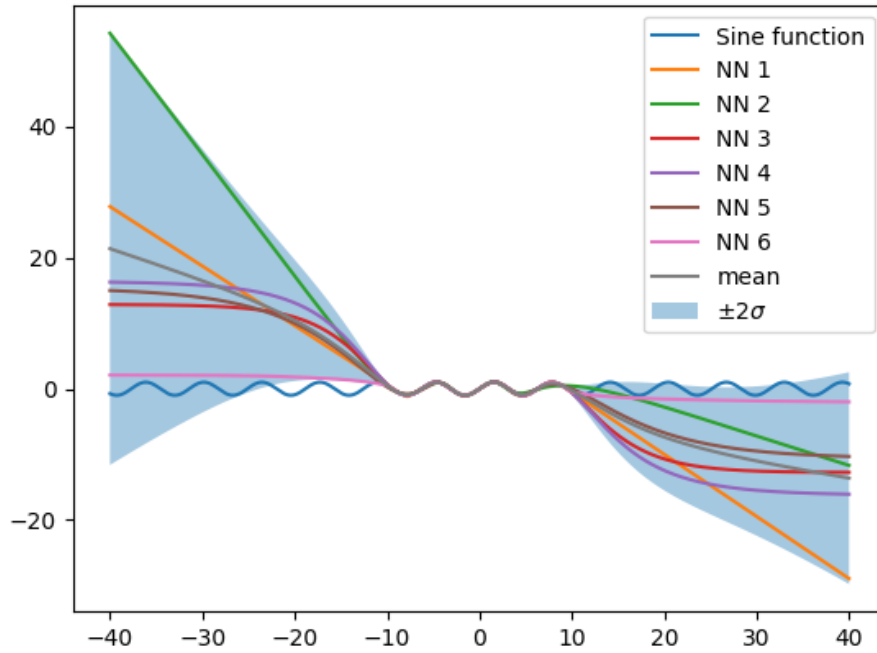
the actual values, so the predictions can be as bad as they keep the right order. I chose the std thresholds based on a histogram of std values which were measured alongside the plain CMA-ES optimization. For one example see figure 2.1.

Even this simple strategy worked well for Sphere, e.g., reducing the number of evaluations to about 60-150.



■ **Figure 2.1** NTK-GP standard deviations over Sphere optimization: evolution (left), histogram (right).

So I choose I want to explore some approach based on predicting also the uncertainty. The NTK was often in papers or supplemental materials, compared to an ensemble of NNs[31]. So I made a dummy dataset of a transformed sine wave to see how the point estimates of the Ensemble of neural networks predict the uncertainty compared to the NTK. Very soon, I noticed it predicts very poorly because the networks, even though randomly initialized, converge to similar solutions, i.e., their outside-of-distribution predictions are usually very degenerated and nothing compared to NTK-GP. I tried it with different activation functions (AF), but the effect was the same. The out 1D out-of-distribution tails degenerate "always" in a similar fashion no matter what AF I used. However, the solutions are unlike each other across the activation functions. For e.g. saturating AFs such as sigmoid, hardsigmoid, hardtanh usually ends up with constant prediction higher or lower than the last point depending on whether the function ended upwards (higher) or downwards (lower). In contrast, activation functions which are at least on one side not constrained e.g. RELU, GELU, SELU, ... tends to predict the last trend, either linear or other depending on the function. The idea was born: let's use an Ensemble of different activation functions for uncertainty predictions.



■ **Figure 2.2** Example of sine modeling using an ensemble of ReLU, ELU, ERF, tanh, sigmoid, and hardtanh AFs

### 2.2.1 Ensemble of Activation Functions

When I saw how much better it works compared to a simple ensemble (extrapolating and interpolating uncertainty), I thought I could not be the first one to notice this. But to my surprise, the topic of the ensemble of activation functions is not a thoroughly researched topic. Ensembles of activation functions were used to improve the performance with CNNs[32]. To the best of my knowledge a little research was done on leveraging it for uncertainty estimation. The following quote illustrates the maturity of this research “to the best of our knowledge, ours is the first work to investigate their [NN ensemble] usefulness for predictive uncertainty estimation”[33]. To emphasize the quote from 2017, it is about ensembles with the same activation function. The only research I found using an ensemble of activation functions for uncertainty estimation is Randomized Activation Functions (RAF)[34], which proposes it as a modification of the Anchored Ensemble[35]. Unfortunately, there is no evaluation of the ensemble without the anchored loss function. Besides these, there are other methods for uncertainty estimation with NNs such as Deep Ensemble or Bayesian Neural Networks; based on the paper, the RAF seems to become one of the state-of-the-art, and for further information, I recommend these[34, 36].

### 2.2.2 Anchored Ensemble and RAF

The Anchored Ensemble[35] is an ensemble with regularization, not about the size of the weights, but about the initial parameters. This direction of research is even more interesting in the times of NTK, where the proofs about the equivalence of the Gaussian process with NTK to the neural network is about the evolution of the network while training about their initial weights. The RAF's contribution is for large number of members to be initialized with random activation function, but for a ensemble smaller than or equal to 5 it is just initialized with given set of five different activation functions.

## Methodology

Given the broad scope of the assignment, I was tasked with choosing an approach and starting from scratch without any existing code base. My work mostly unfolded in three ways.

First, I conducted experiments typically focused on a single function and a single instance with a set seed. This approach allowed me to observe how individual decisions directly affected the outcomes. However, this method can be misleading, as it might reveal effective parameter settings that only work in specific problem instances.

Second, the experiments usually involved either a single function across multiple instances (typically eight to align with running eight parallel processes on a single-threaded Pytorch and NumPy setup) or multiple functions (four to eight) with a single instance. Some results from these experiments were archived, primarily to validate whether specific combinations had been tested or to identify the closest combinations tried.

In this chapter, I detail the various aspects explored during these experiments. In the "Results" section, I will discuss only two groups of experiments: one focused only on dimension 2, and the second benchmarked with all functions and all instances for dimensions 2, 3, 5, 10, and 20. I believe this somewhat aligns with Nikolaus Hansen, who wrote:

We believe, however, that in the process of solver design, a benchmarking framework like COCO has its limitations. During the design phase, usually:

- fewer benchmark functions should be used,
- the functions and measuring tools should be tailored to the given solver and the design question, and
- The overall procedure should be more informal and interactive with rapid iterations.

[14].

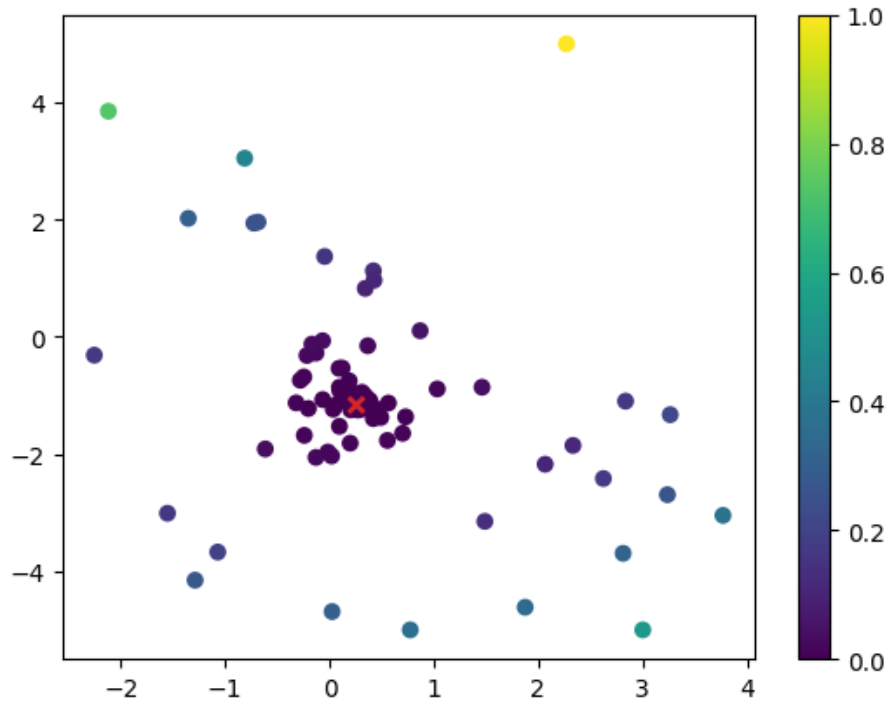
### 3.1 CMA-ES Variant

In terms of CMA-ES and Python language, there are two prominent implementations. First is the one written by the team of Nikolaus Hansen himself CMA-ES/pycma[37]. The second one is ModularCMAES[38] from the IOH-profiler project. I was told the preferred way is to use the ModularCMAES because it's newer, so it could have learned from past mistakes of pycma. My feelings about these implementations are mixed. The ModularCMAES seems better organized (which makes sense being the "second" implementation), but on the other hand, the documentation of ModularCMAES is outdated, so you have to look into the code, with pycma, you also have to look at the code at the end, which is a little more chaotic, but the documentation up to date which I see as a plus. What made the decision for me, was that empirically (not thoroughly tested) the ModularCMAES was underperforming the pycma. I first discovered I need to use the active, IPOP CMA-ES, which is the default in pycma, but opt-in for ModularCMAES. Later I discovered the pycma by default does not handle boundaries in any way, which is beneficial, because it does not distort the population when offspring is sampled outside the boundary. Also, the ModularCMAES (empirically) restarted itself sooner because of being close to "numerical instability," which maybe was, but instead of studying all these tolerances and settings, I chose to use the pycma CMA-ES, especially when I wanted to produce results comparable with others as is.

### 3.2 Surrogate

#### 3.2.1 Subset Method

I already mentioned the importance of subsetting the Archive to train set. I began with the method of taking up to  $20 \times D$  points with the smallest norm to the population center. Later I realized the this is not robust, because sometimes 40 closest points are all close, and sometimes there real outliers, so I added a constrained in any direction the point can not be further form the population mean, than  $10std$  of the population. If 10 standard deviations sound like a lot, with less, there would often be very few points to learn from, e.g., 12 for dimension 2. To illustrate the need for subsetting the archive in the figure3.1 is the x values of the archive and 6 search points for the evaluation, the reason it seems as a single red cross, is that the search points differ in small order of magnitude if the picture was zoomed in a lot, the difference would show off. The figure also shows that when the target variable is linearly scaled into  $[0, 1]$  range, all the values around the search points have the same color, again, subsetting helps but other target variable transformations might be used as well.



■ **Figure 3.1** Example of search points and the archive.

### 3.2.2 Feature Transformation

The feature transformation proved itself very soon; for a long time, I standardized with sample mean and sample standard deviation. Then, I moved to mean and standard deviation (multiplied by the step size) of CMA-ES in the iteration. The values should be similar, but because of the stochastic nature, why use estimates when there is access to the distribution the points are drawn from?

### 3.2.3 Target Variable Transformation

Transformation of the target variable, i.e., the function value, was much trickier. In order to train the model in a generation, we have access to the points we need to model, as well as to the distribution from which the points were drawn, but in terms of target values, there is only the archive. The first try was to use standardization using  $y$  subset mean and std estimates. It didn't work well. Then I tried weighted mean, std by the norm (later Mahalanobis norm) of the points to the sample mean (later CMA-ES mean).

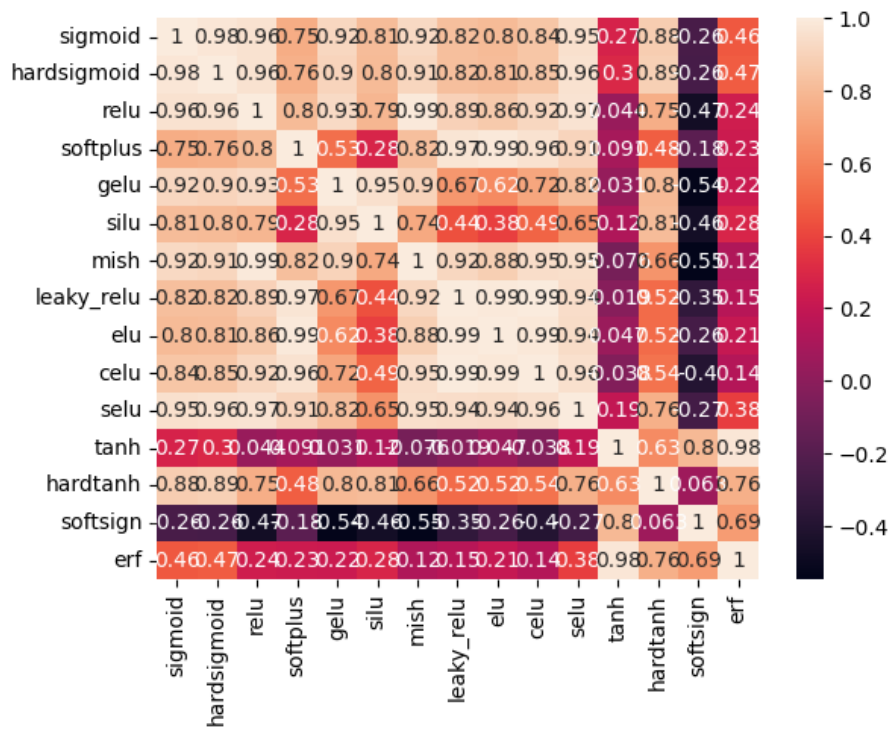
The last transformation I tried was based on experience. Sometimes (often, but I might have chosen the wrong functions to test this, I usually used

functions 1, 6, and 10), you would need to use the third power of the norm so the mean is close to the median, i.e., to prevent the situation where from 30 points there are twenty looking like this: 78.8573XYZ, and your mean is 80. So I tried to the logarithmic transformation to get rid of the outliers affecting the mean and variance so much. And it worked.

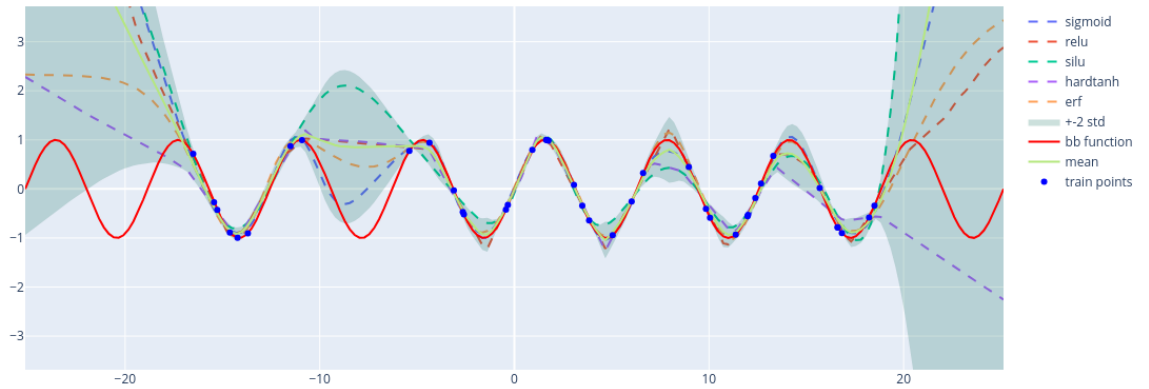
### 3.2.4 Weighing the Train Points

### 3.2.5 The Model

I tested two models, a plain Ensemble of Networks with different activation functions termed EAF and the RAF[34]. To choose which activation functions to use, I modeled a sample of sine values with out-of-bounds test points. For these (with comparable end loss), I looked at correlation coefficients and tried to choose a subset that is the least correlated with each other.



■ **Figure 3.2** Example correlation coefficients of NNs for predictions of NNs with different AF.



■ **Figure 3.3** Example predictions of NNs with sigmoid, relu, silu, hardtanh and erf AFs.

### 3.3 Evolution Control

#### 3.4 Implementation

Since I started with an empty code base, in contrast to enhancing the supervisor's or other's project I needed to make a lot of arbitrary decisions. Such as the version of Python and the libraries I use.

I tried to implement a modular structure, where it is easy, to change any component mentioned above without touching the rest. The top optimization function is called `seek_minimum`. I believe the modularity is best comprehensible from a simple usage example.

■ **Code listing 3.1** Function `seek_minimum` Signature

```
def seek_minimum(
    problem: Problem,
    *,
    es: Cma,
    surrogate_and_ec: Optional[SurrogateAndEc] = None,
    sort_archive=True,
    **details
) -> Archive: ...
```



■ **Code listing 3.2** SurrogateCallable Protocol

```
class SurrogateCallable(Protocol):
    def __call__(
        self,
        *,
        x_train: NDArray[np.float64],
        y_train: NDArray[np.float64],
        x_test: NDArray[np.float64],
        es: Cma,
    ) -> Model: ...
```

■ **Code listing 3.3** EvolutionControl Protocol

```
class EvolutionControl(Protocol):
    def __call__(
        self, *, model: Model, points: XPop, problem: Problem, archive: Archive
    ) -> ValuesAndEvaluatedIdx: ...
```

Another take was to keep the implementation in one file. Generally, I am not a fan of huge-file implementations. While testing it locally (running through the night), visualizing it in Google Colaboratory (enjoying Plotly interactive charts), and benchmarking it at the MetaCentrum server, I found it too convenient to keep it in one simple file such that you can literally copy-paste the whole implementation. If this code has a future, it would be smart to separate it into reasonable modules, especially when the structure (components above) is mostly crystallized.

I took the versions of Python packages from Google Colaboratory.

► **Note 3.1.** The TensorFlow running RAF with these packages raises multiple warnings. I checked if they were produced by the anchoring method (by commenting the part out), and they were not. I also ran the experiment with an older version of TensorFlow 2.9.1, and the performance looked the same; it just took twice the time. The performance of the mean prediction looked comparable to the plain PyTorch ensemble. Also, changing the `data_noise` did affect the surrogate performance (for some functions positively, for some, the opposite). Having said that, I kept this for later. especially If I started utilizing the uncertainty.

### 3.5 Benchmarking Platforms

In order to experiment with the BBOB functions, the most straightforward way is to use COCO[14] the framework the suit was originally proposed with. There is also a newer COCO-inspired alternative IOHprofiler[39]. The BBOB functions are the same, but a few differences are worth mentioning. The IOHprofiler is much easier to install into the Python ecosystem, as it is a classical Python package that you can install with

```
pip install ioh
```

. In contrast, the COCO must be installed from source. It is not a big deal when having full control of the environment, but I didn't succeed with installing it in the Google Colaboratory, and installing it on the MetaCentrum server had its pitfalls. The second difference is between COCO's

```
problem.final_target_hit
```

, which evaluates to true if the problem was evaluated at a maximum  $1e-8$  away from the optimum. The IOHprofiler's

```
problem.state.optimum_found
```

evaluates to true when the optimum was reached exactly (with the floating point rounding).

## Chapter 4

# Results

First things first. All following figures and tables except 4.1 and 4.2 were generated by `cocopp`<sup>1</sup> from experiment data generated by COCO[14] and using their `LATEX`template<sup>2</sup>. The algorithm called 2022/CMA-ES-pycma is from the development branch of `pycma`<sup>3</sup> and this commit<sup>4</sup>, information source<sup>5</sup>. The algorithm called 2023/default-CMA-ES is also from the development branch of `pycma`, commit was not specified<sup>6</sup>, but the COCO Data Archie<sup>7</sup> states it's equivalent to the default CMA-ES from the `pycma` module, version 3.3.0.

► Note 4.1. The IOHAnalyzer<sup>8</sup>, should “visualize solver performance (given in COCO format) interactively in the browser and offers both fixed-target and fixed-budget views”[14]. However, the same data that works with COCO failed to load in the IOHProfiler. This is likely due to `cocopp` recursively and fault-tolerantly using all data in any sub-directory, while IOHProfiler might be more precise about the desired structure. For this reason, I stuck with `cocopp`, and modified the this line<sup>9</sup> to 250.

In this work I focused on Improving the performance of CMA-ES with a small budget. The value 250 E/D in dimensions 2,3,5,10,20. These choices were taken over from this work[23]. Given CMA-ES is not the best method for very contained evaluations: “For budgets below 500D function evaluations, the best performance achieve NEWUOA, MCS and GLOBAL. For larger bud-

---

<sup>1</sup><https://numbbo.github.io/coco-doc/apidocs/cocopp/>

<sup>2</sup><https://github.com/numbbo/coco/tree/master/code-postprocessing/latex-templates>

<sup>3</sup><https://github.com/CMA-ES/pycma/>

<sup>4</sup><https://github.com/CMA-ES/pycma/commit/3bec053cfc2a6a4720ac8a32dcee133778f99608>

<sup>5</sup><https://inria.hal.science/hal-03671431/file/wshp-bbob-article.pdf>

<sup>6</sup><https://hal.science/hal-04089923/document>

<sup>7</sup><https://numbbo.github.io/data-archive/bbob/>

<sup>8</sup><https://iohprofiler.github.io/IOHAnalyzer/>

<sup>9</sup><https://github.com/numbbo/coco/blob/v2.6.3/code-postprocessing/cocopp/genericsettings.py#L55>

gets, BIPOP-CMA-ES and IPOP-SEP-CMA-ES become superior.”[40], it is reasonable to try to push in this direction.

**pre-1** Is ModularCmaes ES with EAF surrogate, and 0.1 fixed evaluation ratio EC sorted by LCB.

**pre-4** Is pycma CMA-ES with EAF surrogate, and 0.1 fixed evaluation ratio EC sorted by LCB.

**pre-4** Is pycma CMA-ES with double initial population, EAF surrogate, and 0.05 fixed evaluation ratio EC sorted by LCB.

**pre-7** Is pycma CMA-ES with EAF surrogate, and 0.75 Kendall  $\tau$  Threshold in current generation evaluations EC.

**1** Is pycma CMA-ES with RAF surrogate and Kendall EC1.11.

**2- $\alpha$**  Is pycma CMA-ES with Log transformed targets shifted by constant from 0, RAF surrogate, and Kendall EC1.11.

**2** Is pycma CMA-ES with Log transformed targets shifted by 0.05 quantile difference of the archived targets, RAF surrogate, and Kendall EC1.11.

## 4.1 Preliminary - Dimension 2

The COCO platform’s default target precisions are 51 evenly log-spaced values between  $10^{-8}$  and  $10^2$ [14].

■ **Table 4.1**  $10^{-8}$ -target hits of 15 instances in  $2D$ .

Test FID	pre-1	pre-4	pre-5	pre-7	1	2- $\alpha$	2
1	15	15	15	15	15	14	15
2	15	15	15	15	15	13	15
3	1	4	3	2	4	1	4
4	0	1	1	0	2	0	3
5	15	15	15	15	15	15	15
6	0	2	0	2	2	1	10
7	11	13	8	12	14	13	14
8	7	12	15	10	11	12	7
9	11	15	15	13	13	11	10
10	0	0	0	11	15	13	8
11	0	0	0	12	13	14	6
12	3	8	7	5	3	7	6
13	0	0	0	4	3	1	2
14	0	0	0	13	12	10	10
15	3	3	1	1	2	1	3
16	8	8	10	9	6	3	1
17	1	1	0	0	1	0	1
18	1	1	0	0	0	0	0
19	1	4	0	2	4	3	1
20	0	1	2	0	1	1	0
21	6	7	3	5	5	3	8
22	5	5	5	8	4	2	7
23	0	6	6	0	0	0	0
24	0	0	0	0	0	0	0

The test two is a perfect showcase for one phenomenon. Because of the training times, and different instances and also the training is randomized, it's not possible to test everything, so when I wanted to improve the test 1, I choose function 6, because from lq-CMA[24] I got the impression, that the function is hard, so improving function 6 will improve a lot of functions. Unfortunately, even though the function 6 was improved immensely, a lot of functions stayed the same or even deteriorated. I believe there are things to be solved, it's

just, finding a way about offsetting the minimum, so it have the right amount of improvement ability helped a lot, but only somewhere.

The COCO platform[14] offers runtime-based targets in the expensive optimization scenario. The artificial best algorithm<sup>10</sup> from BBOB-2009[40] serves as a reference algorithm with either the five budgets of  $0.5D$ ,  $1.2D$ ,  $3D$ ,  $10D$ , and  $50D$  function evaluations, or with 31 targets evenly space on the logarithmic scale between  $0.5D$  and  $50D$ [20].

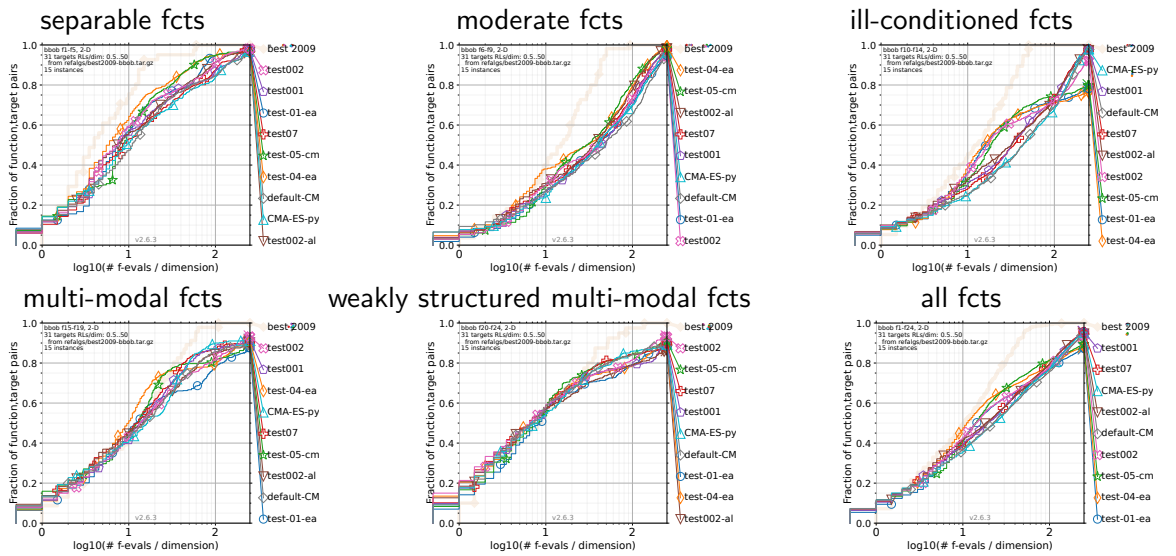
---

<sup>10</sup>The algorithm is artificial because different target values might used runtime results from different algorithms[20].

■ **Table 4.2** Best of BBOB2009 50 E/D-target hits of 15 instances in  $2D$ .

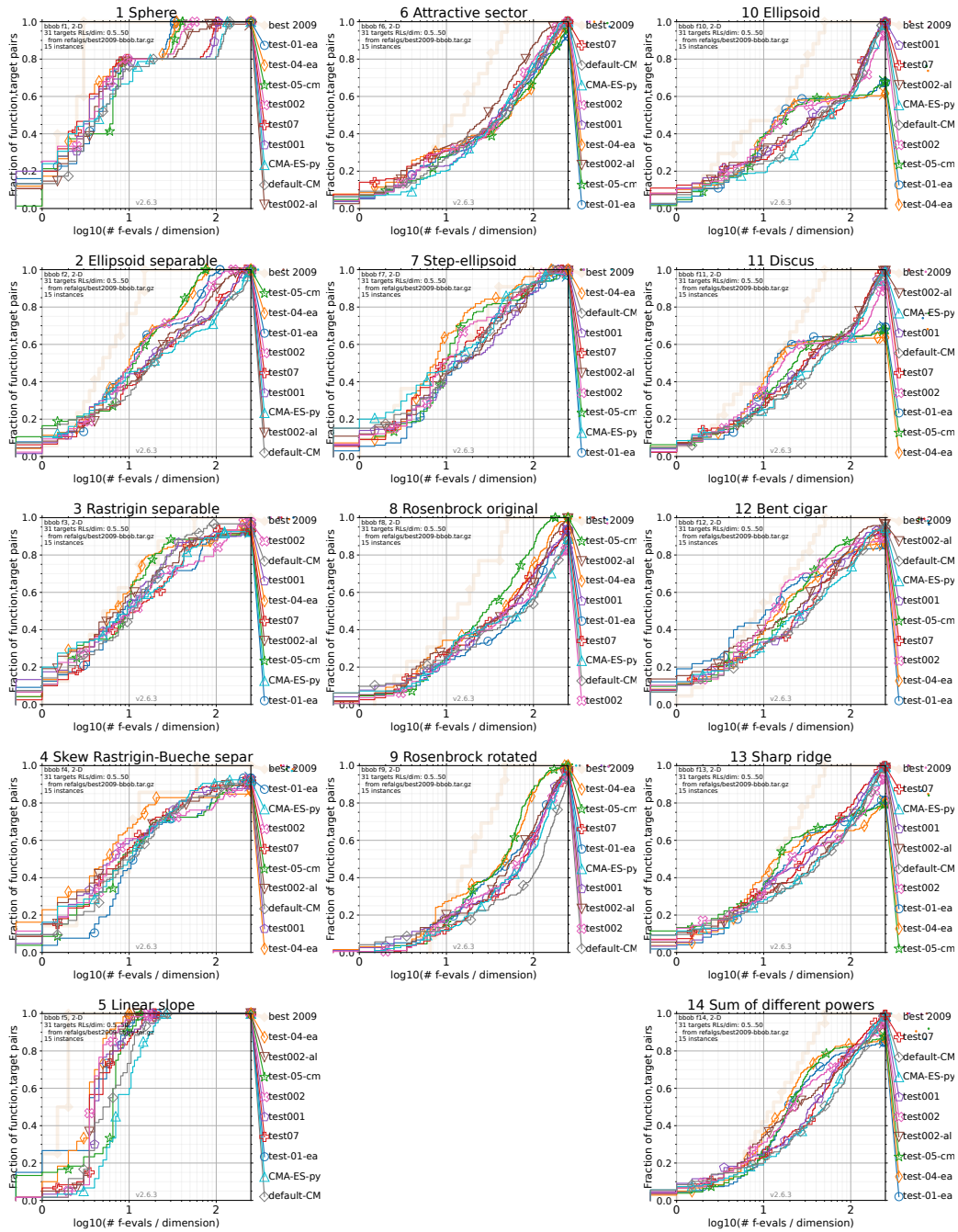
Test FID	pre-1	pre-4	pre-5	pre-7	1	2- $\alpha$	2
1	15	15	15	15	15	14	15
2	15	15	15	15	15	13	15
3	10	12	10	10	12	10	13
4	10	8	9	8	7	8	9
5	15	15	15	15	15	15	15
6	4	11	7	15	14	11	14
7	14	15	14	15	15	14	14
8	9	13	15	12	11	14	9
9	11	15	15	13	13	11	10
10	0	0	0	13	15	13	8
11	0	0	0	12	13	14	7
12	6	9	9	11	8	12	8
13	0	0	0	14	13	10	8
14	0	0	0	13	12	10	10
15	13	15	13	10	11	13	14
16	11	12	13	13	11	8	11
17	11	14	13	15	11	14	13
18	12	12	13	10	12	13	12
19	4	5	1	3	7	6	5
20	6	6	8	5	5	6	8
21	9	10	7	8	11	8	13
22	7	8	9	12	6	7	9
23	13	12	15	14	13	13	13
24	7	8	12	6	7	5	9

I interpret it as the final target  $1e^{-8}$ 4.1 measures whether the algorithm is capable of finding the global optimum, while the 50 E/D expensive targets4.2, are satisfiable with very different minimums as long as they are observed in time. This can be with the fixed number of surrogate evaluations, as in terms of the expensive targets the performance is improved compared to the default CMA-ES, while considering the final targets  $1e^{-8}$ , the performance is worse mostly because convergence to local optimum.

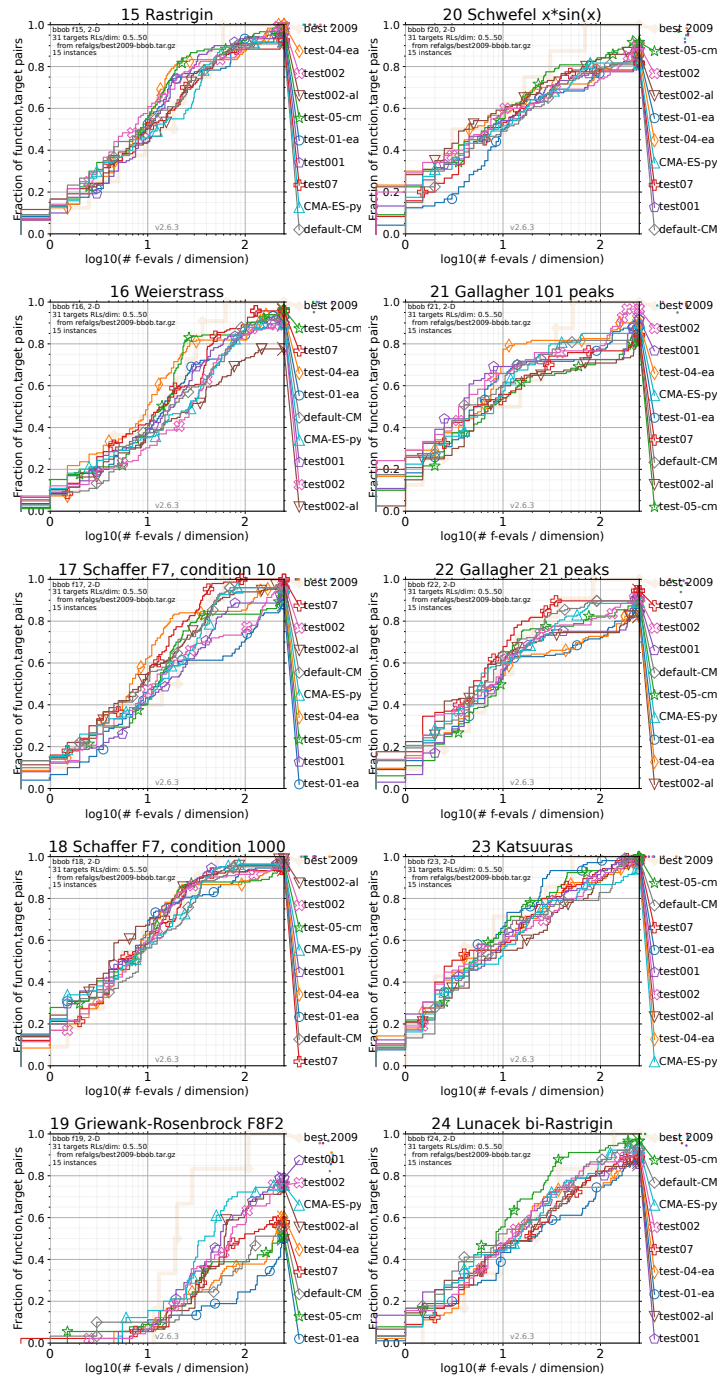


**Figure 4.1** Bootstrapped empirical cumulative distribution of the number of  $f$ -evaluations divided by dimension (FEvals/DIM) for all functions and subgroups in 2-D. The targets are chosen from  $10^{[-8..2]}$  such that the best algorithm from BBOB 2009 just not reached them within a given budget of  $k \times \text{DIM}$ , with 31 different values of  $k$  chosen equidistant in logscale within the interval  $\{0.5, \dots, 50\}$ . As reference algorithm, the best algorithm from BBOB 2009 is shown as light thick line with diamond markers.

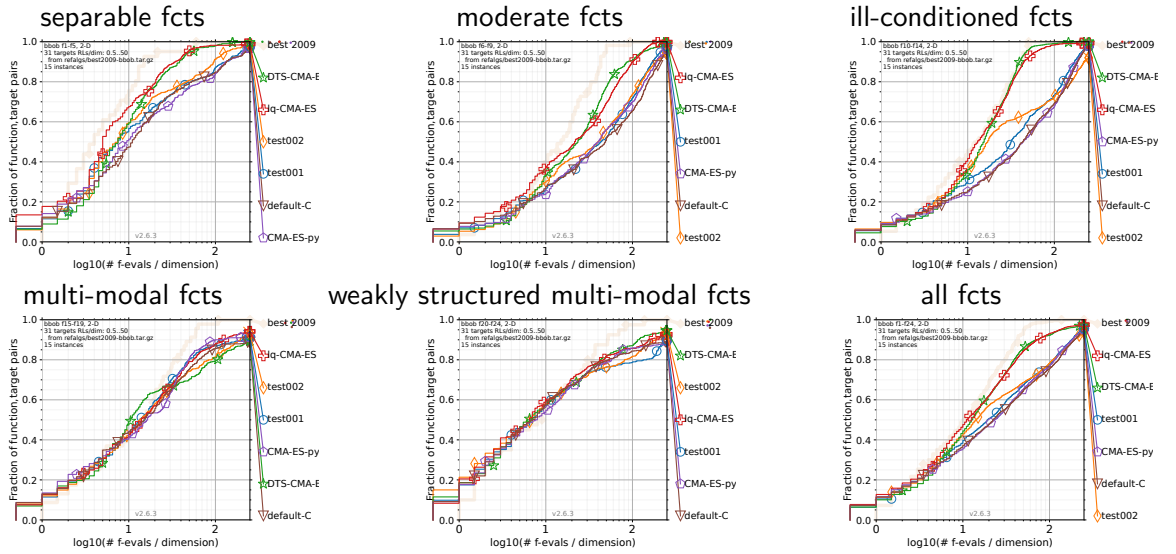




■ **Figure 4.2** Empirical cumulative distribution of simulated (bootstrapped) runtimes, measured in number of  $f$ -evaluations divided by dimension (FEvals/DIM) in dimension 2 and for those targets in  $10^{[-8, .2]}$  that have just not been reached by the best algorithm from BBOB 2009 in a given budget of  $k \times \text{DIM}$ , with 31 different values of  $k$  chosen equidistant in logscale within the interval  $\{0.5, \dots, 50\}$ .

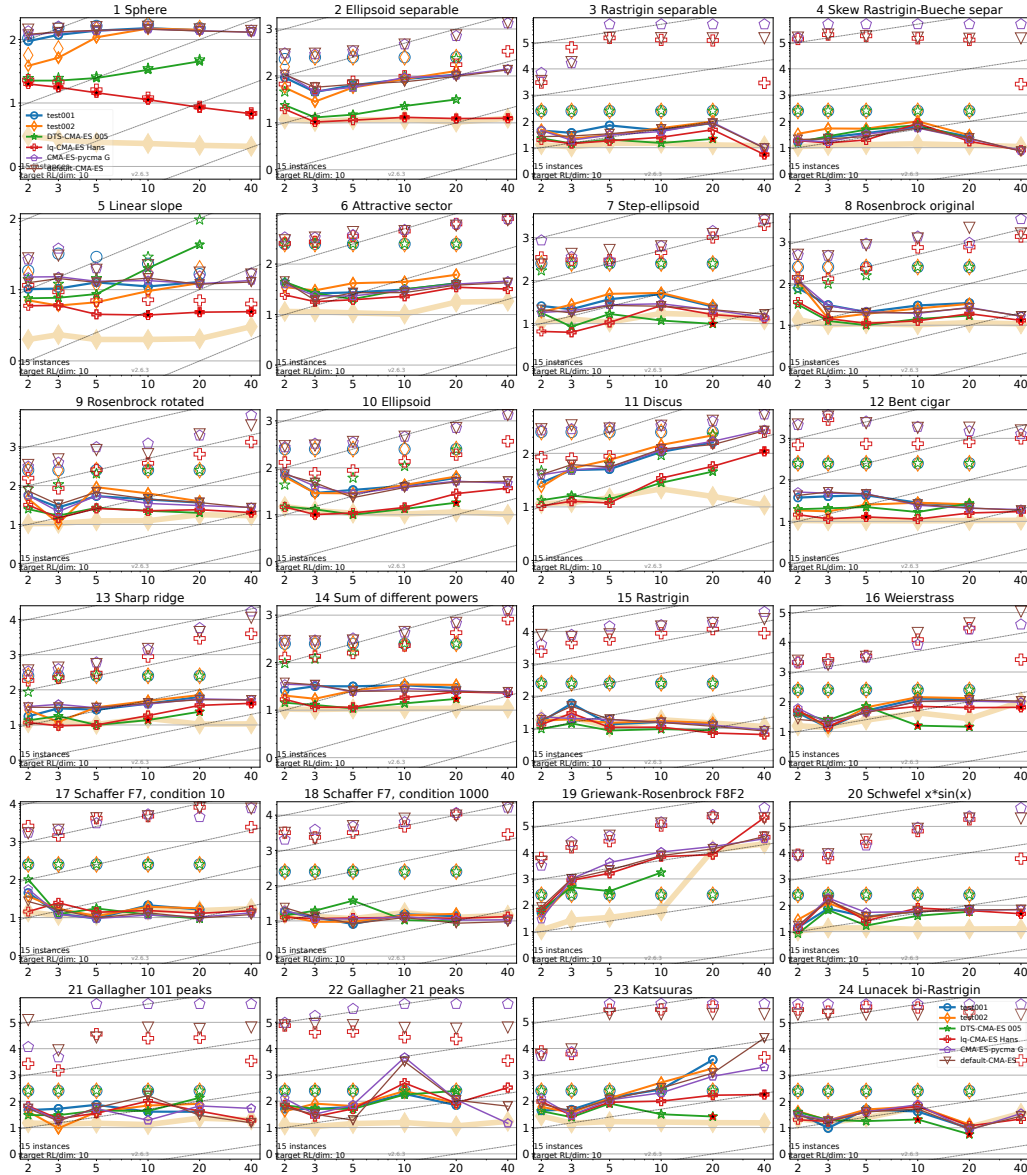


■ **Figure 4.3** Empirical cumulative distribution of simulated (bootstrapped) run-times, measured in number of  $f$ -evaluations divided by dimension (FEvals/DIM) in dimension 2 and for those targets in  $10^{[-8..2]}$  that have just not been reached by the best algorithm from BBOB 2009 in a given budget of  $k \times \text{DIM}$ , with 31 different values of  $k$  chosen equidistant in logscale within the interval  $\{0.5, \dots, 50\}$ .



**Figure 4.4** Bootstrapped empirical cumulative distribution of the number of  $f$ -evaluations divided by dimension (FEvals/DIM) for all functions and subgroups in 2-D. The targets are chosen from  $10^{[-8..2]}$  such that the best algorithm from BBOB 2009 just not reached them within a given budget of  $k \times \text{DIM}$ , with 31 different values of  $k$  chosen equidistant in logscale within the interval  $\{0.5, \dots, 50\}$ . As reference algorithm, the best algorithm from BBOB 2009 is shown as light thick line with diamond markers.

## 4.2 Benchmarking RAF Dimensions 2,3,5,10,20



**Figure 4.5** Expected running time (ERT in number of  $f$ -evaluations as  $\log_{10}$  value) divided by dimension versus dimension. The target function value is chosen such that the best algorithm from BBOB 2009 just failed to achieve an ERT of  $10 \times \text{DIM}$ . Different symbols correspond to different algorithms given in the legend of  $f_1$  and  $f_{24}$ . Light symbols give the maximum number of evaluations from the longest trial divided by dimension. Black stars indicate a statistically better result compared to all other algorithms with  $p < 0.01$  and Bonferroni correction number of dimensions (six). Legend:  $\circ$ : CMA-ES-pycma Gharafi,  $\diamond$ : default-CMA-ES Gissler,  $\star$ : test-01-eaf,  $+$ : test-04-eaf-cma,  $\square$ : test-05-cma-eaf-double-005,  $\nabla$ : test001,  $\times$ : test002,  $\Delta$ : test002-alpha,  $\diamond$ : test07.

## Conclusions

A year-long experimentation with ensembles of wide networks has demonstrated their capability to estimate uncertainty and learn various functions. However, no silver bullet was found. It has been a long process of steady improvements, culminating in slightly improved performance for some functions in the most tested domain, dimension 2, and slightly worse to worse performance for other functions.

I believe that with comparable time investment, similar results to those achieved with GP-based surrogates could be possible. The question remains whether such an investment is worthwhile when one can already use the DTS-CMA-ES or lq-CMA-ES.

The mini-framework for surrogate testing developed in this work I consider to be the main contribution. The individual components are clearly designated and implemented in Python, allowing anyone interested in trying model X as a surrogate to simply write a small wrapper around it and integrate it with some of the prepared subset methods, transformations, evolution controls, etc. The openness of this implementation is also valuable, ensuring that not only the model source code but also all the environment and scripts are available, thereby lowering the entry barrier for any CMA-ES surrogate enthusiast.

### 5.1 Future Work

In terms of Neural Networks in surrogate modeling, I believe there is great potential to leverage uncertainty estimation either directly, such as choosing which points to evaluate based on uncertainty, or in more steps, such as evaluating the points from the most uncertain, and when the error of your prediction for some number of predictions drops below a threshold, using the predictions for the rest of the points.

Similarly, I see a huge potential in using continual learning. The last time Neural Networks were used as CMA-ES surrogates, it was one model con-

tinually learning the functions. I don't think keeping the same model is the right way, but leveraging the evaluated points in the generation perhaps to re-estimate the mean and uncertainty of the predictions, compare them, or, as in the DTS case, measure the error and use it for the evolution control.

Lastly, the Mini-framework developed as a side product of this thesis has a lot of rough edges. The logging is primitive and made ad-hoc to things I needed to measure and, in a way, sufficient to get things done. However improving the measurement capabilities would really help for experiments with CMA-ES surrogates. Because I often needed to run the experiments again when I came up with something I would like to measure, such as Kendall  $\tau$  coefficient first for evaluated points, then for predictions of the whole archive, then for predictions of archive subset, then predictions of individual members of the ensemble, etc. Especially when more people would use it and share the use insights this could be very beneficial for experimenting.

Also, Python is a relatively slow language. One approach to solve this, is to write the code in different language and keep only the interface in Python, e.g. the COCO framework is written in C. I don't know where the Python limits lie, but I would start with profiling, then maybe continue with Numba<sup>1</sup>'s JIT or rewrite the code to JAX<sup>2</sup> where possible and utilize JAX's JIT<sup>3</sup> and `jax.numpy` module<sup>4</sup>. Or to go even further and experiment with MOJO<sup>5</sup>.

---

<sup>1</sup><https://numba.pydata.org/>

<sup>2</sup><https://github.com/google/jax>

<sup>3</sup>[https://jax.readthedocs.io/en/latest/\\_autosummary/jax.jit.html](https://jax.readthedocs.io/en/latest/_autosummary/jax.jit.html)

<sup>4</sup><https://jax.readthedocs.io/en/latest/jax.numpy.html>

<sup>5</sup><https://github.com/modularml/mojo>

## Test 1 and 2 ERT Multiples for BBOB2009 Targets

The tables show: “Expected runtime (ERT in number of  $f$ -evaluations) divided by the respective best ERT measured during BBOB-2009 (when finite) in dimension  $N$ . This ERT ratio and, in braces as dispersion measure, the half difference between 10 and 90%-tile of bootstrapped run lengths appear for each algorithm and run-length based target, the corresponding reference ERT (preceded by the target  $\Delta f$ -value in *italics*) in the first row. #succ is the number of trials that reached the target value of the last column. 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 (24). A  $\uparrow$  signifies the number of trials that were worse than the ERT of the best algorithm from BBOB 2009 shown only when less than 10 percent were worse and the ERT was better. Best results are printed in bold. ”. This description was generated by the COCO platform.













# Bibliography

1. RIOS, Luis Miguel; SAHINIDIS, Nikolaos V. Derivative-free optimization: a review of algorithms and comparison of software implementations. *Journal of Global Optimization*. 2012, vol. 56, no. 3, pp. 1247–1293. ISSN 1573-2916. Available from DOI: [10.1007/s10898-012-9951-y](https://doi.org/10.1007/s10898-012-9951-y).
2. JIN, Yaochu; OLHOFER, M.; SENDHOFF, B. Managing approximate models in evolutionary aerodynamic design optimization. In: *Proceedings of the 2001 Congress on Evolutionary Computation (IEEE Cat. No.01TH8546)*. IEEE, 2001, vol. 1, 592–599 vol. 1. Available from DOI: [10.1109/CEC.2001.934445](https://doi.org/10.1109/CEC.2001.934445).
3. PHAN-TRONG, Dat; TRAN-THE, Hung; GUPTA, Sunil. NeuralBO: A black-box optimization algorithm using deep neural networks. *Neurocomputing*. 2023, vol. 559, p. 126776. ISSN 0925-2312. Available from DOI: [10.1016/j.neucom.2023.126776](https://doi.org/10.1016/j.neucom.2023.126776).
4. LARSON, Jeffrey; MENICKELLY, Matt; WILD, Stefan M. Derivative-free optimization methods. *Acta Numerica*. 2019, vol. 28, pp. 287–404. ISSN 1474-0508. Available from DOI: [10.1017/s0962492919000060](https://doi.org/10.1017/s0962492919000060).
5. LOSHCHILOV, Ilya. *Surrogate-Assisted Evolutionary Algorithms*. 2013. Available also from: <https://theses.hal.science/tel-00823882>. Theses. Université Paris Sud - Paris XI ; Institut national de recherche en informatique et en automatique - INRIA.
6. HANSEN, Nikolaus; ARNOLD, Dirk V.; AUGER, Anne. Evolution Strategies. In: KACPRZYK, Janusz; PEDRYCZ, Witold (eds.). *Handbook of Computational Intelligence*. Springer, 2015. No. 871-898. Available also from: <https://inria.hal.science/hal-01155533>.
7. PITRA, Zbyněk; HANUŠ, Marek; KOZA, Jan; TUMPACH, Jiří; HOLEŇA, Martin. Interaction between model and its evolution control in surrogate-assisted CMA evolution strategy. In: *Proceedings of the Genetic and Evolutionary Computation Conference*. Lille, France: Association for Com-

- puting Machinery, 2021, pp. 528–536. GECCO '21. ISBN 9781450383509. Available from DOI: 10.1145/3449639.3459358.
8. HANSEN, Nikolaus. *The CMA Evolution Strategy: A Tutorial*. arXiv, 2016. Available from DOI: 10.48550/ARXIV.1604.00772.
  9. HANSEN, Nikolaus. *CMA-ES with Two-Point Step-Size Adaptation*. 2008. Research Report. Available also from: <https://inria.hal.science/inria-00276854>.
  10. JASTREBSKI, G.A.; ARNOLD, D.V. Improving Evolution Strategies through Active Covariance Matrix Adaptation. In: *2006 IEEE International Conference on Evolutionary Computation*. 2006, pp. 2814–2821. Available from DOI: 10.1109/CEC.2006.1688662.
  11. HANSEN, Nikolaus. The CMA Evolution Strategy: A Comparing Review. In: *Towards a New Evolutionary Computation: Advances in the Estimation of Distribution Algorithms*. Ed. by LOZANO, Jose A.; LARRAÑAGA, Pedro; INZA, Iñaki; BENGOTXEA, Endika. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, pp. 75–102. ISBN 978-3-540-32494-2. Available from DOI: 10.1007/3-540-32494-1\_4.
  12. AUGER, A.; HANSEN, N. A restart CMA evolution strategy with increasing population size. In: *2005 IEEE Congress on Evolutionary Computation*. 2005, vol. 2, 1769–1776 Vol. 2. Available from DOI: 10.1109/CEC.2005.1554902.
  13. JIN, Yaochu; SENDHOFF, Bernhard. Fitness Approximation In Evolutionary Computation - a Survey. In: LANGDON, William B.; CANTÚPAZ, Erick; MATHIAS, Keith E.; ROY, Rajkumar; DAVIS, David; POLI, Riccardo; BALAKRISHNAN, Karthik; HONAVAR, Vasant G.; RUDOLPH, Günter; WEGENER, Joachim; BULL, Larry; POTTER, Mitchell A.; SCHULTZ, Alan C.; MILLER, Julian F.; BURKE, Edmund K.; JONOSKA, Natasa (eds.). *GECCO 2002: Proceedings of the Genetic and Evolutionary Computation Conference, New York, USA, 9-13 July 2002*. Morgan Kaufmann, 2002, pp. 1105–1112.
  14. HANSEN, Nikolaus; AUGER, Anne; ROS, Raymond; MERSMANN, Olaf; TUŠAR, Tea; BROCKHOFF, Dimo. COCO: a platform for comparing continuous optimizers in a black-box setting. *Optimization Methods and Software*. 2020, vol. 36, no. 1, pp. 114–144. ISSN 1029-4937. Available from DOI: 10.1080/10556788.2020.1808977.
  15. HANSEN, Nikolaus; FINCK, Steffen; ROS, Raymond; AUGER, Anne. *Real-Parameter Black-Box Optimization Benchmarking 2009: Noiseless Functions Definitions*. 2009. Research Report, RR-6829. INRIA. Available also from: <https://inria.hal.science/inria-00362633>.

16. AUGER, A.; HANSEN, N.; PEREZ ZERPA, J. M.; ROS, R.; SCHOENAUER, M. Experimental Comparisons of Derivative Free Optimization Algorithms. In: VAHRENHOLD, Jan (ed.). *Experimental Algorithms*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 3–15. ISBN 978-3-642-02011-7.
17. *bbob* — *numbbo.github.io* [<https://numbbo.github.io/coco/testsuites/bbob>]. [N.d.]. [Accessed 05-05-2024].
18. FINCK, Steffen; HANSEN, Nikolaus; ROS, Raymond; AUGER, Anne. *Real-Parameter Black-Box Optimization Benchmarking 2010: Presentation of the Noiseless Functions*. 2009. Tech. rep., 2009/20. Research Center PPE. Available also from: <https://numbbo.github.io/gforge/downloads/download16.00/bbobdocfunctions.pdf>. Accessed: May 2024.
19. HANSEN, Nikolaus; ROS, Raymond; MAUNY, Nikolas; SCHOENAUER, Marc; AUGER, Anne. Impacts of Invariance in Search: When CMA-ES and PSO Face Ill-Conditioned and Non-Separable Problems. *Applied Soft Computing*. 2011, vol. 11, pp. 5755–5769. Available from DOI: 10.1016/j.asoc.2011.03.001.
20. HANSEN, Nikolaus; AUGER, Anne; BROCKHOFF, Dimo; TUSAR, Dejan; TUŠAR, Tea. *COCO: Performance Assessment*. 2016. Available also from: <https://inria.hal.science/hal-01315318>. ArXiv e-prints, arXiv:1605.03560.
21. KHALDI, Mohammed Imed Eddine; DRAA, Amer. Surrogate-assisted evolutionary optimisation: a novel blueprint and a state of the art survey. *Evolutionary Intelligence*. 2023. ISSN 1864-5917. Available from DOI: 10.1007/s12065-023-00882-8.
22. JIN, Yaochu; OLHOFER, M.; SENDHOFF, B. A framework for evolutionary optimization with approximate fitness functions. *IEEE Transactions on Evolutionary Computation*. 2002, vol. 6, no. 5, pp. 481–494. Available from DOI: 10.1109/TEVC.2002.800884.
23. BAJER, Lukáš; PITRA, Zbyněk; REPICKÝ, Jakub; HOLEŇA, Martin. Gaussian Process Surrogate Models for the CMA Evolution Strategy. *Evolutionary Computation*. 2019, vol. 27, no. 4, pp. 665–697. ISSN 1063-6560. Available from DOI: 10.1162/evco\_a\_00244.
24. HANSEN, Nikolaus. A Global Surrogate Assisted CMA-ES. In: *GECCO 2019 - The Genetic and Evolutionary Computation Conference*. Prague, Czech Republic: ACM, 2019, pp. 664–672. Available from DOI: 10.1145/3321707.3321842.
25. KIM, Samuel; LU, Peter Y.; LOH, Charlotte; SMITH, Jamie; SNOEK, Jasper; SOLJAČIĆ, Marin. *Deep Learning for Bayesian Optimization of Scientific Problems with High-Dimensional Structure*. 2022. Available from arXiv: 2104.11667 [cs.LG].

26. KOZA, Jan; TUMPACH, Jiri; PITRA, Zbynek; HOLENA, Martin. Combining Gaussian Processes and Neural Networks in Surrogate Modelling for Covariance Matrix Adaptation Evolution Strategy. In: *ITAT*. 2021, pp. 29–38.
27. RUZICKA, Jiri; KOZA, Jan; TUMPACH, Jiri; PITRA, Zbynek; HOLENA, Martin. Combining Gaussian Processes with Neural Networks for Active Learning in Optimization. In: *IAL@ PKDD/ECML*. 2021, pp. 105–120.
28. LU, Minfang; NING, Shuai; LIU, Shuangrong; SUN, Fengyang; ZHANG, Bo; YANG, Bo; WANG, Lin. OPT-GAN: A Broad-Spectrum Global Optimizer for Black-Box Problems by Learning Distribution. *Proceedings of the AAAI Conference on Artificial Intelligence*. 2023, vol. 37, no. 10, pp. 12462–12472. ISSN 2159-5399. Available from DOI: 10.1609/aaai.v37i10.26468.
29. MÜLLER, Samuel; FEURER, Matthias; HOLLMANN, Noah; HUTTER, Frank. *PFNs4BO: In-Context Learning for Bayesian Optimization*. 2023. Available from arXiv: 2305.17535 [cs.LG].
30. HE, Bobby; LAKSHMINARAYANAN, Balaji; TEH, Yee Whye. *Bayesian Deep Ensembles via the Neural Tangent Kernel*. 2020. Available from arXiv: 2007.05864 [stat.ML].
31. NOVAK, Roman; XIAO, Lechao; HRON, Jiri; LEE, Jaehoon; ALEMI, Alexander A.; SOHL-DICKSTEIN, Jascha; SCHOENHOLZ, Samuel S. *Neural Tangents: Fast and Easy Infinite Neural Networks in Python*. 2019. Available from arXiv: 1912.02803 [stat.ML].
32. MAGUOLO, Gianluca; NANNI, Loris; GHIDONI, Stefano. *Ensemble of Convolutional Neural Networks Trained with Different Activation Functions*. 2020. Available from arXiv: 1905.02473 [cs.CV].
33. LAKSHMINARAYANAN, Balaji; PRITZEL, Alexander; BLUNDELL, Charles. *Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles*. 2017. Available from arXiv: 1612.01474 [stat.ML].
34. STOYANOVA, Yana; GHANDI, Soroush; TAVAKOL, Maryam. Toward robust uncertainty estimation with random activation functions. In: *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence*. AAAI Press, 2023. AAAI’23/IAAI’23/EAAI’23. ISBN 978-1-57735-880-0. Available from DOI: 10.1609/aaai.v37i12.26768.
35. PEARCE, Tim; LEIBFRIED, Felix; BRINTRUP, Alexandra. Uncertainty in Neural Networks: Approximately Bayesian Ensembling. In: CHIAPPA, Silvia; CALANDRA, Roberto (eds.). *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, vol. 108, pp. 234–244. Proceedings of Machine Learning Research.



Available also from: <https://proceedings.mlr.press/v108/pearce20a.html>.

36. ABDAR, Moloud; POURPANAH, Farhad; HUSSAIN, Sadiq; REZAZADEGAN, Dana; LIU, Li; GHAVAMZADEH, Mohammad; FIEGUTH, Paul; CAO, Xiaochun; KHOSRAVI, Abbas; ACHARYA, U. Rajendra; MAKARENKOV, Vladimir; NAHAVANDI, Saeid. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*. 2021, vol. 76, pp. 243–297. ISSN 1566-2535. Available from DOI: 10.1016/j.inffus.2021.05.008.
37. HANSEN, Nikolaus; AKIMOTO, Youhei; BAUDIS, Petr. *CMA-ES/pycma on Github* [Zenodo, DOI:10.5281/zenodo.2559634]. 2019. Available from DOI: 10.5281/zenodo.2559634.
38. NOBEL, Jacob de; VERMETTEN, Diederick; WANG, Hao; DOERR, Carola; BÄCK, Thomas. Tuning as a Means of Assessing the Benefits of New Ideas in Interplay with Existing Algorithmic Modules. In: *Proceedings of the Genetic and Evolutionary Computation Conference Companion*. Lille, France: Association for Computing Machinery, 2021, pp. 1375–1384. GECCO '21. ISBN 9781450383516. Available from DOI: 10.1145/3449726.3463167.
39. DOERR, Carola; WANG, Hao; YE, Furong; RIJN, Sander van; BÄCK, Thomas. *IOHprofiler: A Benchmarking and Profiling Tool for Iterative Optimization Heuristics*. arXiv, 2018. Available from DOI: 10.48550/ARXIV.1810.05281.
40. HANSEN, Nikolaus; AUGER, Anne; ROS, Raymond; FINCK, Steffen; POŠÍK, Petr. Comparing results of 31 algorithms from the black-box optimization benchmarking BBOB-2009. In: *Proceedings of the 12th Annual Conference Companion on Genetic and Evolutionary Computation*. Portland, Oregon, USA: Association for Computing Machinery, 2010, pp. 1689–1696. GECCO '10. ISBN 9781450300735. Available from DOI: 10.1145/1830761.1830790.

# List of abbreviations

AF	Activation Function
BBOB	Black-box Optimization Benchmarking test suite
GAN	Generative Adversarial Network
GP	Gaussian Process
NN	Neural Network
NTK	Neural Tangent Kernel
PFN	Prior-data Fitted Network
RAF	Random Activation Function
VAE	Variational autoencoder

# Attachment Contents

```
| readme.adoc.....information about other files in this directory
| nncmaes.....implementation source code
| tex.....LATEX source code
| thesis.pdf.....thesis PDF
```