# Revisiting COCO with automated benchmarking in mind
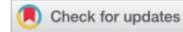
Nikolaus Hansen
Inria
CMAP, CNRS, Ecole Polytechnique, Institut Polytechnique de Paris, France

Presented as invited talk at *BENCH@GECCO25 - Workshop on Good Benchmarking Practices for Evolutionary Computation*, July 2025

feel free…

# …feel free to ask questions…

# COCO: a platform for comparing continuous optimizers in a black-box setting

Nikolaus Hansen[a,b], Anne Auger[a,b], Raymond Ros[c], Olaf Mersmann[d], Tea Tušar[e] and Dimo Brockhoff[a,b]

[a]Randopt Team, Inria, Palaiseau, France; [b]CMAP, CNRS, Ecole Polytechnique, Institut Polytechnique de Paris, Palaiseau, France; [c]LRI, Université Paris-Sud, Orsay, France; [d]Computational Statistics, TU Dortmund University, Dortmund, Germany; [e]Jožef Stefan Institute, Ljubljana, Slovenia

## ABSTRACT

We introduce COCO, an open-source platform for Comparing Continuous Optimizers in a black-box setting. COCO aims at automating the tedious and repetitive task of benchmarking numerical optimization algorithms to the greatest possible extent. The platform and the underlying methodology allow to benchmark in the same framework deterministic and stochastic solvers for both single and multiobjective optimization. We present the rationals behind the (decade-long) development of the platform as a general proposition for guidelines towards better benchmarking. We detail underlying fundamental concepts of COCO such as the definition of a problem as a function instance, the underlying idea of instances, the use of target values, and runtime defined by the number of function calls as the central performance measure. Finally, we give a quick overview of the basic code structure and the currently available test suites.

## 1. Introduction

We consider the continuous black-box optimization or search problem to minimize

$$f : X \subset \mathbb{R}^n \to \mathbb{R}^m \quad n, m \geq 1, \tag{1}$$

where the search domain $X$ is typically a bounded hypercube or the entire continuous space.[1] More specifically, we aim to find, as quickly as possible, one or several solutions $\boldsymbol{x}$ in the search space $X$ with *small* value(s) of $f(\boldsymbol{x}) \in \mathbb{R}^m$.

A continuous optimization algorithm, denoted as *solver*, addresses the above problem. In this paper, we only consider zero-order black-box optimization [19,57,58]: while the search domain $X \subset \mathbb{R}^n$ and its boundaries are accessible, no other prior knowledge about $f$ is available to the solver.[2] That is, $f$ is considered as a black-box, also known as an oracle, and the *only* way the solver can acquire information on $f$ is by querying the value $f(\boldsymbol{x})$ of a solution $\boldsymbol{x} \in X$. Zero-order black-box optimization is thus a derivative-free optimization setting.[3] We generally consider 'time' to be the number of calls to the function $f$ and will define 'runtime' correspondingly.

---

# Anytime Performance Assessment in Blackbox Optimization Benchmarking

Nikolaus Hansen, Anne Auger, Dimo Brockhoff, and Tea Tušar

*Abstract*—We present concepts and recipes for the anytime performance assessment when benchmarking optimization algorithms in a blackbox scenario. We consider runtime—oftentimes measured in the number of blackbox evaluations needed to reach a target quality—to be a universally measurable cost for solving a problem. Starting from the graph that depicts the solution quality versus runtime, we argue that runtime is the *only* performance measure with a generic, meaningful, and quantitative interpretation. Hence, our assessment is solely based on runtime measurements. We discuss proper choices for solution quality indicators in single- and multi-objective optimization, as well as in the presence of noise and constraints. We also discuss the choice of the target values, budget-based targets, and the aggregation of runtimes by using simulated restarts, averages, and empirical cumulative distributions which generalize convergence graphs of single runs. The presented performance assessment is to a large extent implemented in the comparing continuous optimizers (COCO) platform freely available at https://github.com/numbbo/coco.

*Index Terms*—Anytime optimization, benchmarking, blackbox optimization, performance assessment, quality indicator.

## I. INTRODUCTION

WE PRESENT practical concepts and ideas for the performance assessment of optimization algorithms when benchmarked in a blackbox and anytime scenario. Going beyond a simple ranking of algorithms, we aim to provide a *quantitative* and *meaningful* performance assessment, which allows for conclusions like *algorithm A is seven times faster than algorithm B* in solving a given problem or in solving problems with certain characteristics. To achieve this end in a comparative and timeless manner, we argue that we should measure the *number of blackbox evaluations* to reach a predefined *quality indicator* value (a target). More generally, we argue to measure a cost that is defined on a ratio scale and is comparable across publications. We call this measure

the *runtime* of the algorithm to reach a given target. Yet, our assessment methodology does not depend on any specific cost measure, as long as the costs are quantitative and comparable.[1]

In this article, we formalize the optimization goal by a so-called quality indicator. Its definition may heavily depend on the optimization scenario, e.g., the number of objectives or constraints. Broadly speaking, a quality indicator is based on the sequence of all so-far visited solutions. In the simplest case, it is the objective function value of the last visited solution.

Runtimes represent the cost of optimization. Compared to the quality indicator, the definition of costs depends to a lesser extent on the specific optimization scenario. To sustain reproducibility and comparability across publications, we recommend against CPU or wall-clock time as cost measure[2] (see also Hooker [22] for a further discussion on the unwanted consequences of benchmarking based on CPU time).

Benchmarking is usually computationally expensive and benchmarking for a *single* budget seems vastly inefficient by 1) addressing only one of many possible budget scenarios (scenarios heavily depend on the software and hardware environment) and 2) throwing away most of the data generated during the experiment. An anytime approach to benchmarking prevents these drawbacks. To allow for a budget-free performance assessment even for non-anytime algorithms that have a maximum or timeout budget as decisive or mandatory *input* parameter (decided by the user), we collect data with an any-budget *experimental procedure* that runs repeated experiments with increasing input budget [31].[3] Non-anytime algorithms that do not take a maximum budget as an input parameter can be accurately assessed only by the time of their final solution proposal.

In this article, we advocate to routinely use (anytime) *empirical runtime distributions* to assess the performance of optimization algorithms. We demonstrate how to directly compare the runtime distributions of algorithms that have

---

[1]We are grateful to the anonymous reviewer pointing this out to us.

[2]An exploratory CPU timing experiment to get an estimate of the internal time complexity of the algorithm is still advisable, like it is prescribed in the comparing continuous optimizer (COCO) platform [20].

[3]We can stop the procedure when the last budget was not fully exhausted. Increasing the budget each time by a factor of $r > 1$ adds to the overall computational costs for the experimentation less than $\sum_{i=0}^{\infty} 1/r^i = r/(r-1)$ times the last consumed budget. For the performance assessment, always the data from the smallest eligible budget is used. The performance assessment will be too optimistic by tacitly assuming that the budget can be set properly without additional costs. On the other hand, runtimes may be overestimated (by less than a factor of $r$). A code example is provided in

# COCO/BBOB Flow Chart



Figure by Tea Tušar, in Hansen et al (2020), COCO: A platform for comparing continuous optimizers in a black-box setting. *Optimization Methods and Software*, published online, Aug 2020.

# Benchmarking: Possible Goals

- Understanding algorithms

- Regression testing

- Measuring performance in a systematic way (a performance "profile")

- Running a competition → ranking (or selection)

- Algorithm selection (or ranking), could be automated

# Automated Benchmarking: Possible Goals

- ~~Understanding algorithms~~

- ~~Regression testing~~

- Measuring performance in a systematic way (a performance "profile")

- Running a competition → ranking (or selection)

- Algorithm selection (or ranking), could be automated

how…

# Benchmarking: The Global Picture

Two *surprisingly (but not completely) independent issues* to address

1. **What to benchmark**? For example, which collection of test problems?

2. **How to assess performance**?

   - experimental setup

   - data collection

   - measures used

suites…

# What to benchmark: COCO/BBOB Test Suite(s)

- Functions are

  - Based on known (analytical) functions, modelling a "known" difficulty

  - Comprehensible

    aiding human interpretation, however, in science, comprehensible is invariably superior to incomprehensible

  - Scalable

  - Difficult (also: non-separable)

    compared to the "typical standard" (at that time)

  - Quasi-randomized as instances

    with arbitrary shifts and smallish irregularities to avoid artificial exploits and mitigate overfitting, emulates repetition of experiments

  - Come sometimes in pairs

    to observe the effect of a single property

- Require to define target values (function instance + target = problem)

  natural targets in the discrete search domain are known fitness levels and the global optimum, we may need experiments to define useful targets

revisiting...

# What to automated benchmark

- Functions should be

  ✓ Difficult (also: non-separable)
  
  compared to the "typical standard" (at that time)

  ✓ Quasi-randomized as instances
  
  with arbitrary shifts and smallish irregularities to avoid artificial exploits and mitigate overfitting, emulates repetition of experiments

- Scalable

- ~~Based on known (analytical) functions, modelling a "known" difficulty~~

- ~~Comprehensible~~

  aiding human interpretation, however, in science, comprehensible is invariably superior to incomprehensible

- ~~Come sometimes in pairs~~

  to observe the effect of a single property

- Require to define target values (function instance + target = problem)
  
  natural targets in the discrete search domain are known fitness levels and the global optimum, we may need experiments to define useful targets

representative …

✓ Difficult (also: non-separable)

compared to the "typical standard" (at that time)

✓ Quasi-randomized as instances

with arbitrary shifts and smallish irregularities to avoid artificial exploits and mitigate overfitting, emulates repetition of experiments

- Scalable

- ~~Based on known (analytical) functions, modelling a "known" difficulty~~

- ~~Comprehensible~~

aiding human interpretation, however, in science, comprehensible is invariably superior to incomprehensible

- ~~Come sometimes in pairs~~

to observe the effect of a single property

- Require to define target values (function instance + target = problem)

natural targets in the discrete search domain are known fitness levels and the global optimum, we may need experiments to define useful targets

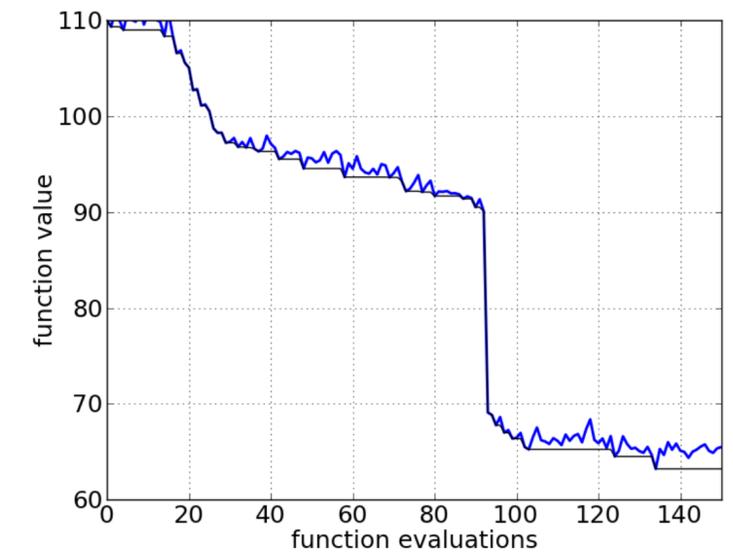- **Benchmark suites "*ultimately define the anticipated purpose*" of the benchmarked algorithms and hence should (must?) be representative**

assessment…

# "quality indicator" versus "time" convergence graphs

is all we have



time means...

# Specifically

- **time**: we use number of function evaluations

  the number is comparable across publications and <span style="color:red">invariant</span> under changes of computer hardware, programming language, compiler, OS,…

- **quality indicator**:

  - single-objective: affine transformation of the function value (to be minimized)

    different for each instance
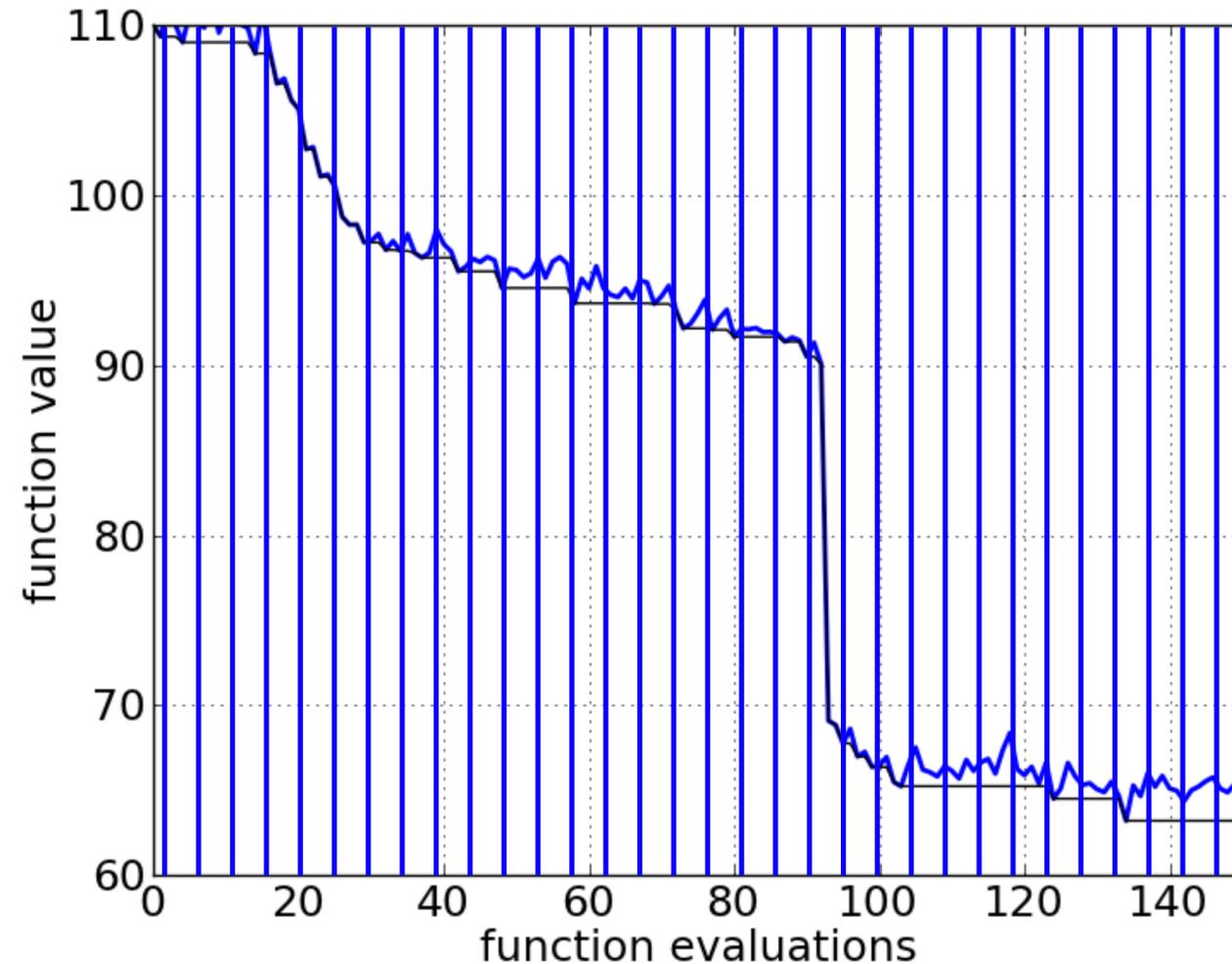
  - multi-objective: negative hypervolume value after objective-wise affine transformation (to be minimized)

  Affine transformations are considered as part of the function definition (benchmark suite definition)

  they also affect the target values that define a problem: target precisions are defined identical for all functions in a suite

  two possibilities to collect comparable…

# Discretization: Two Possibilities

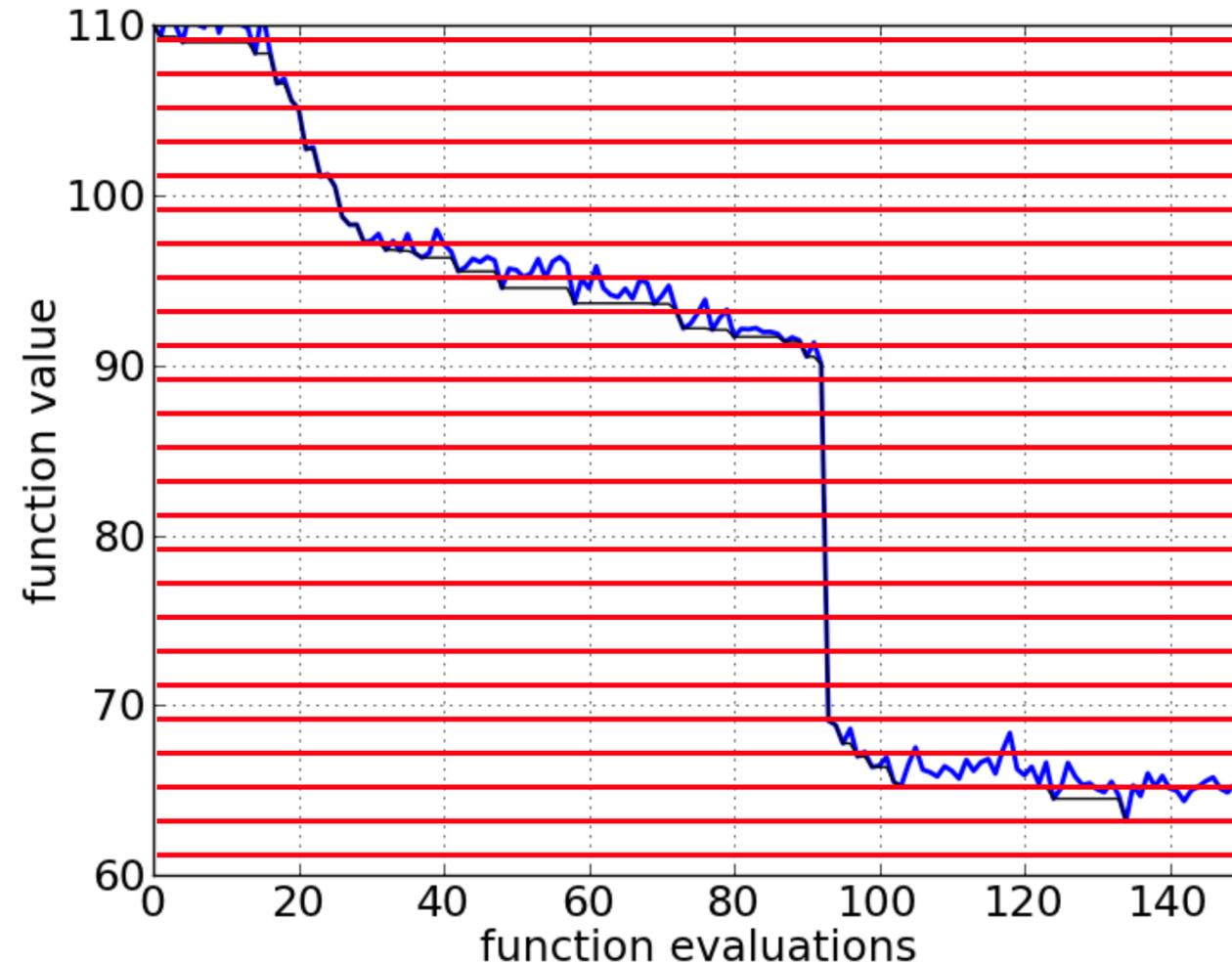- a convergence graph

- lower envelope (a monotonous graph)



- **vertical** (by evaluation) is a natural discretization

  for wall clock or CPU time we would need to determine discretization intervals

- $\implies$ evaluations are the independent variable

  function value (or its attainment fraction) is the dependent variable, the measurement

# Discretization: Two Possibilities



- a convergence graph
- lower envelope (a monotonous graph)

- **horizontal** appears to be an "unnatural" discretization

we need to determine discretization values

- function "target" values (or attainment fractions) are the independent variable

time is the dependent variable, the measurement

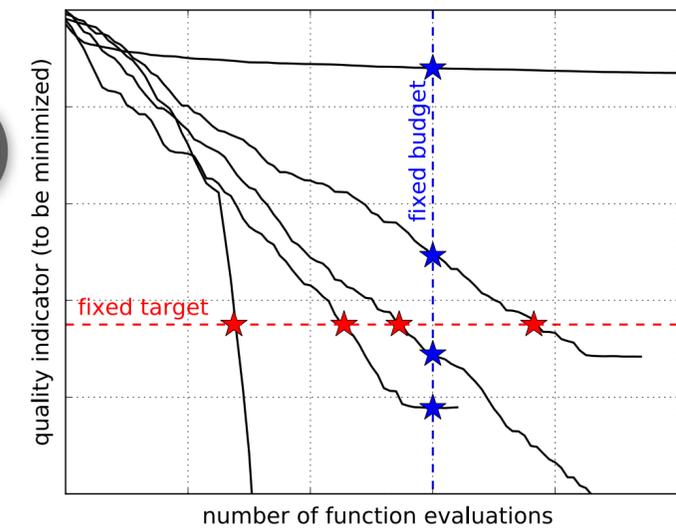- both discretizations recover the original data (in the limit)

a time measurement for each discretization function value, these measurements can be plotted as ECDF

resulting measurements...

# COCO/BBOB: Fixed Target(s) versus Fixed Budget(s)



The resulting measurement

- Vertical discretization = fixed budget design **measures *quality indicator values** and/or their attainment fractions*

- Horizontal discretization = fixed target design **measures *evaluations**, aka *runtime*

Runtime, when measured in number of evaluations,

- is independent of the computational platform, language, compiler, coding style, and other specific experimental conditions,

- is independent, as a measurement, of the specific function on which it has been obtained, that is, *'taking 42 evaluations'* has the same meaning on any function, while *'reaching a function value of 42'* has not,

- is relevant, meaningful and easily interpretable without expert domain knowledge,

- is quantitative on the ratio scale    [66],

- assumes a wide range of values, and

- aggregates over a collection of values in a meaningful way.

Source:  Hansen et al. 2021

The resulting measurement

- **Vertical discretization = fixed budget design measures *quality indicator values and/or their attainment fractions***

- **Horizontal discretization = fixed target design measures *evaluations*, aka *runtime***

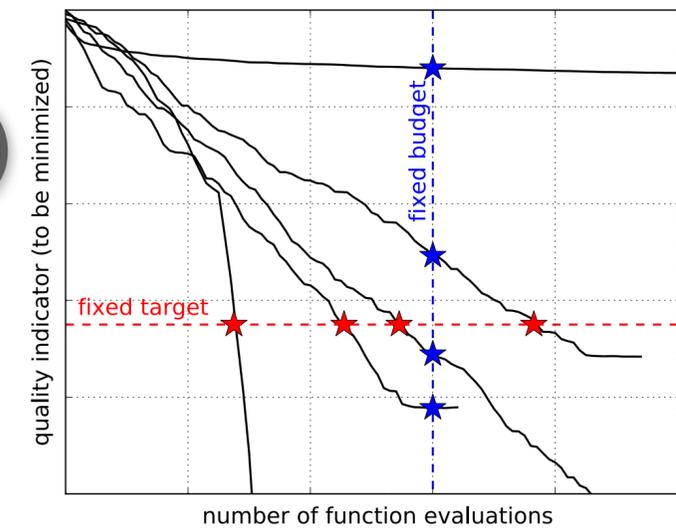Runtime, when measured in number of evaluations,

- is independent of the computational platform, language, compiler, coding style, and other specific experimental conditions,
- is independent, as a measurement, of the specific function on which it has been obtained, that is, *'taking 42 evaluations'* has the same meaning on any function, while *'reaching a function value of 42'* has not,[1]
- is relevant, meaningful and easily interpretable without expert domain knowledge,
- is quantitative on the ratio scale    [66],
- assumes a wide range of values, and
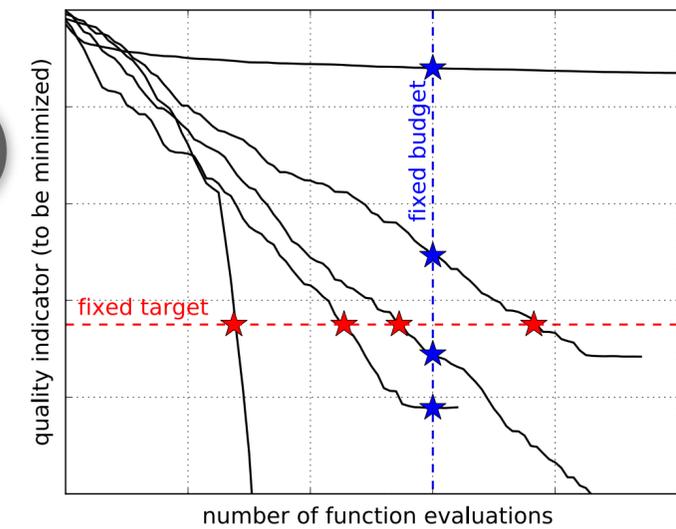- aggregates over a collection of values in a meaningful way.

Source:  Hansen et al. 2021

[1]and the value itself can even be a somewhat arbitrary level set annotation

automated…

The resulting measurement

- Vertical discretization = fixed budget design **measures *quality indicator values* and/or their attainment fractions**

- Horizontal discretization = fixed target design **measures *evaluations*, aka *runtime***

Runtime, when measured in number of evaluations,

- is independent of the computational platform, language, compiler, coding style, and other specific experimental conditions,
- is independent, as a measurement, of the specific function on which it has been obtained, that is, *'taking 42 evaluations'* has the same meaning on any function, while *'reaching a function value of 42'* has not,[1]
- is relevant, meaningful ~~and easily interpretable without expert domain knowledge~~,
- ~~is quantitative on the ratio scale    [66]~~,
- ~~assumes a wide range of values, and~~
- aggregates over a collection of values in a meaningful way.

Source:  Hansen et al. 2021

[1]and the value itself can even be a somewhat arbitrary level set annotation

we only use…

# COCO/BBOB



uses (almost) always

# horizontal discretization

example…

# From a Convergence Graph to the Empirical Runtime Distribution

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

discretization…

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution



equidistance "target" values

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

stars…

20

# From a Convergence Graph to the Empirical Runtime Distribution



Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

bars...

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution



Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

22

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution



Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

23

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution



for the remaining construction, we could use any runtimes, for example, from different runs or different functions

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

empirical distribution…

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution



Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

two curves…

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

flip…

# From a Convergence Graph to the Empirical Runtime Distribution



Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

28

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution



Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution



Hansen et al. 2022. Anytime performance assessment in blackbox
optimization benchmarking. *IEEE Trans. on EC*, 26(6).

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution



when we maximize (instead of minimize), the graph can be considered as an empirical runtime distribution *as is*

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

horizontal bars…

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution



the runtime distribution represents the ordered runtimes depicted as horizontal bars

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

two readings…

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# From a Convergence Graph to the Empirical Runtime Distribution

**Two readings of ECDFs (CDFs)**
Conventional:
 attainment fractions (probabilities)
My (strong) preference:
 runtime values

the runtime distribution
represents the ordered
runtimes depicted as
horizontal bars



López-Ibáñez, et al. 2024. Using the empirical attainment function
for analyzing single-objective black-box optimization
algorithms. *IEEE Trans. on EC, 29(2)*.

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# The Area Above the Empirical Runtime Distribution



the area above the runtime distribution represents a (truncated) average runtime

Specifically, the logarithm of the truncated geometric average equals the bright orange area divided by the truncation fraction

$$\int_0^{\frac{\#\text{solved}}{\#\text{all}}} \log(\#\textbf{evals}(\Delta f_{i(r)}))\,\mathrm{d}r$$

evaluations

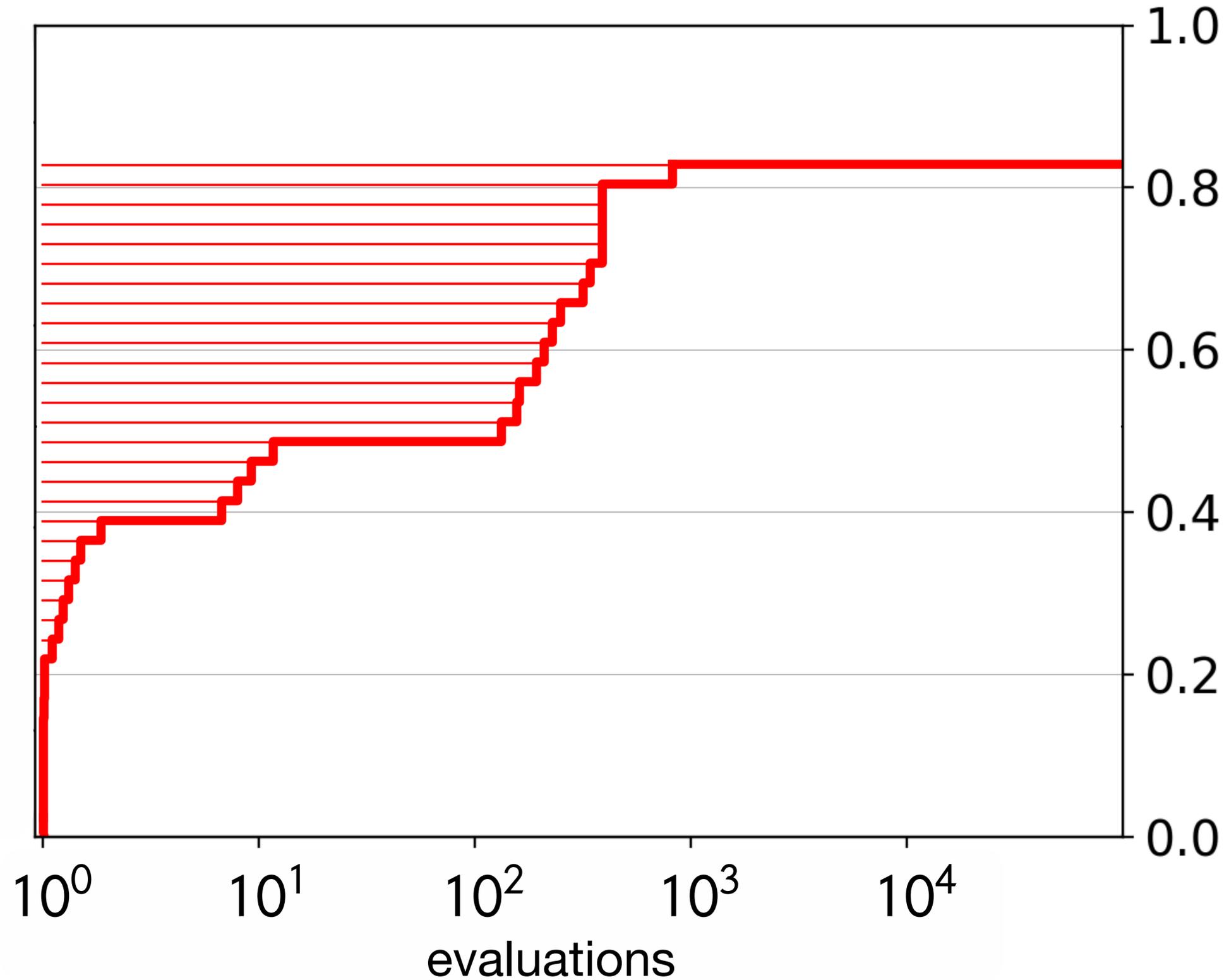Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

formula…

34

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# The Area Above the Empirical Runtime Distribution

the area above the runtime distribution represents a (truncated) average runtime

Specifically, the logarithm of the truncated geometric average equals the bright orange area divided by the truncation fraction

$$\int_0^{\frac{\#solved}{\#all}} \log(\textbf{\#evals}(\Delta f_{i(r)}))\, dr$$

evaluations

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

formula…

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

# The Area Above the Empirical Runtime Distribution

the area above the runtime distribution represents a (truncated) average runtime

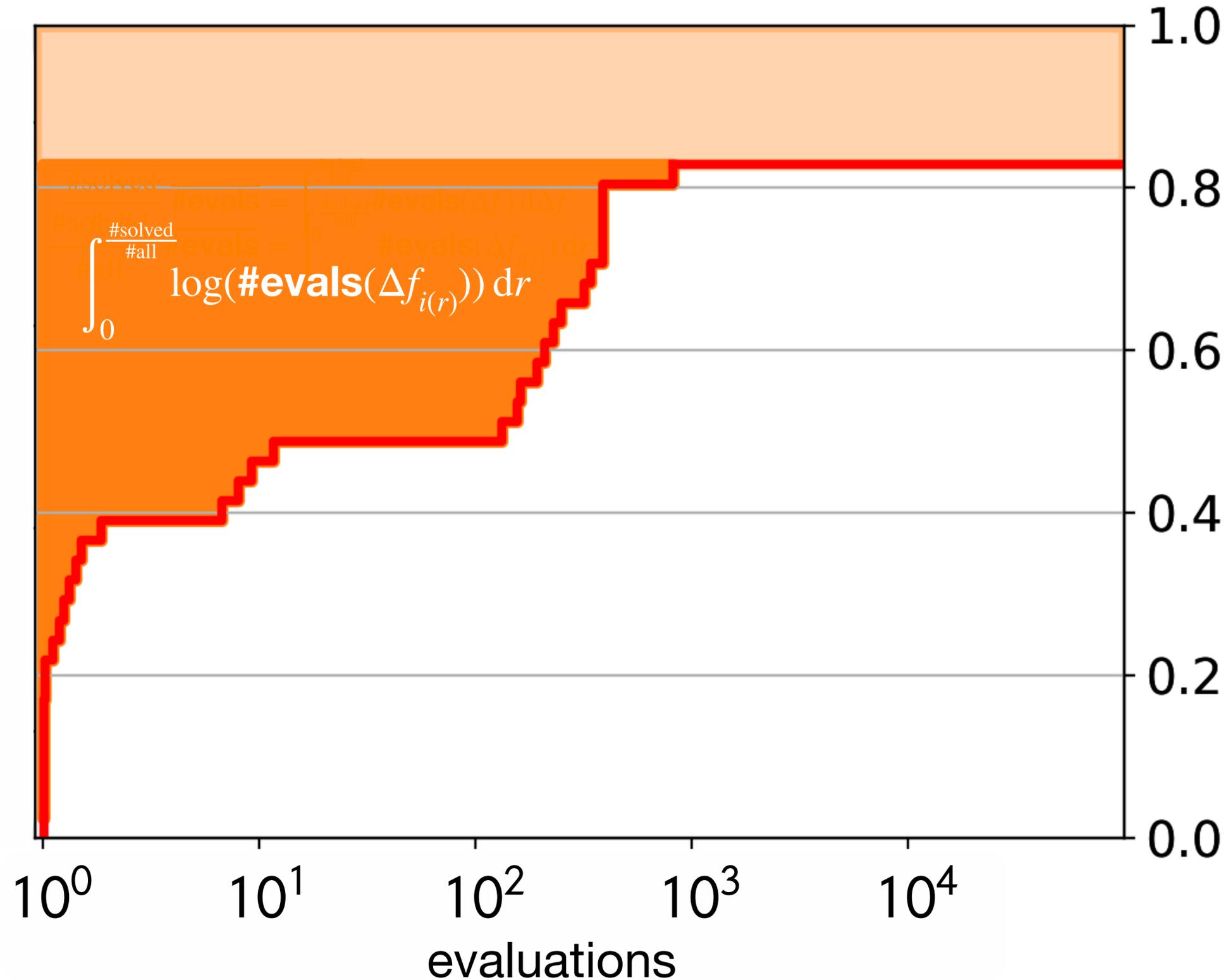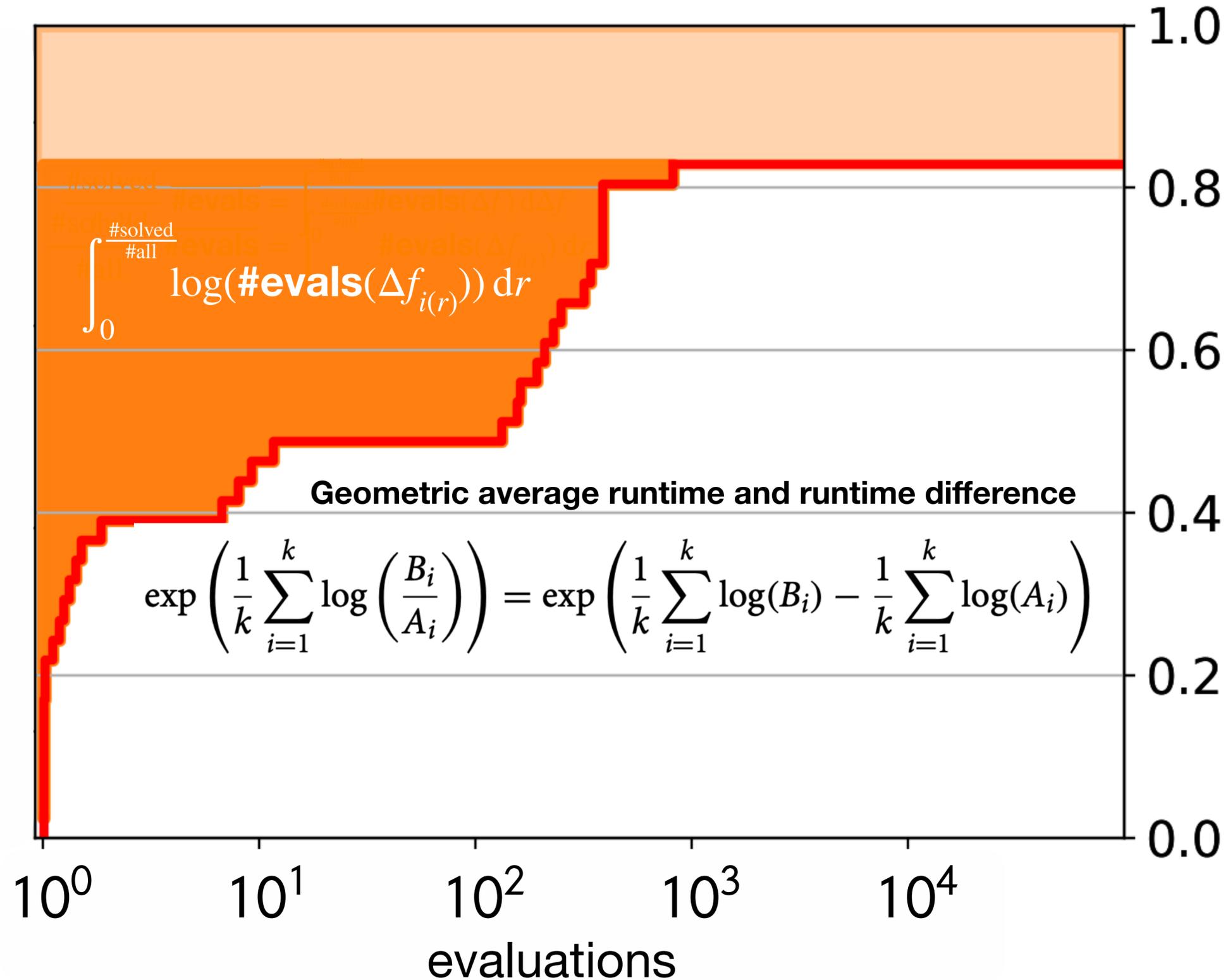Specifically, the logarithm of the truncated geometric average equals the bright orange area divided by the truncation fraction

$$\int_0^{\frac{\#solved}{\#all}} \log(\mathbf{\#evals}(\Delta f_{i(r)}))\, \mathrm{d}r$$

**Geometric average runtime and runtime difference**

$$\exp\left(\frac{1}{k}\sum_{i=1}^{k}\log\left(\frac{B_i}{A_i}\right)\right) = \exp\left(\frac{1}{k}\sum_{i=1}^{k}\log(B_i) - \frac{1}{k}\sum_{i=1}^{k}\log(A_i)\right)$$

evaluations

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

truncation…

36

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

the area above the runtime distribution represents a (truncated) average runtime

Specifically, the logarithm of the truncated geometric average equals the bright orange area divided by the truncation fraction
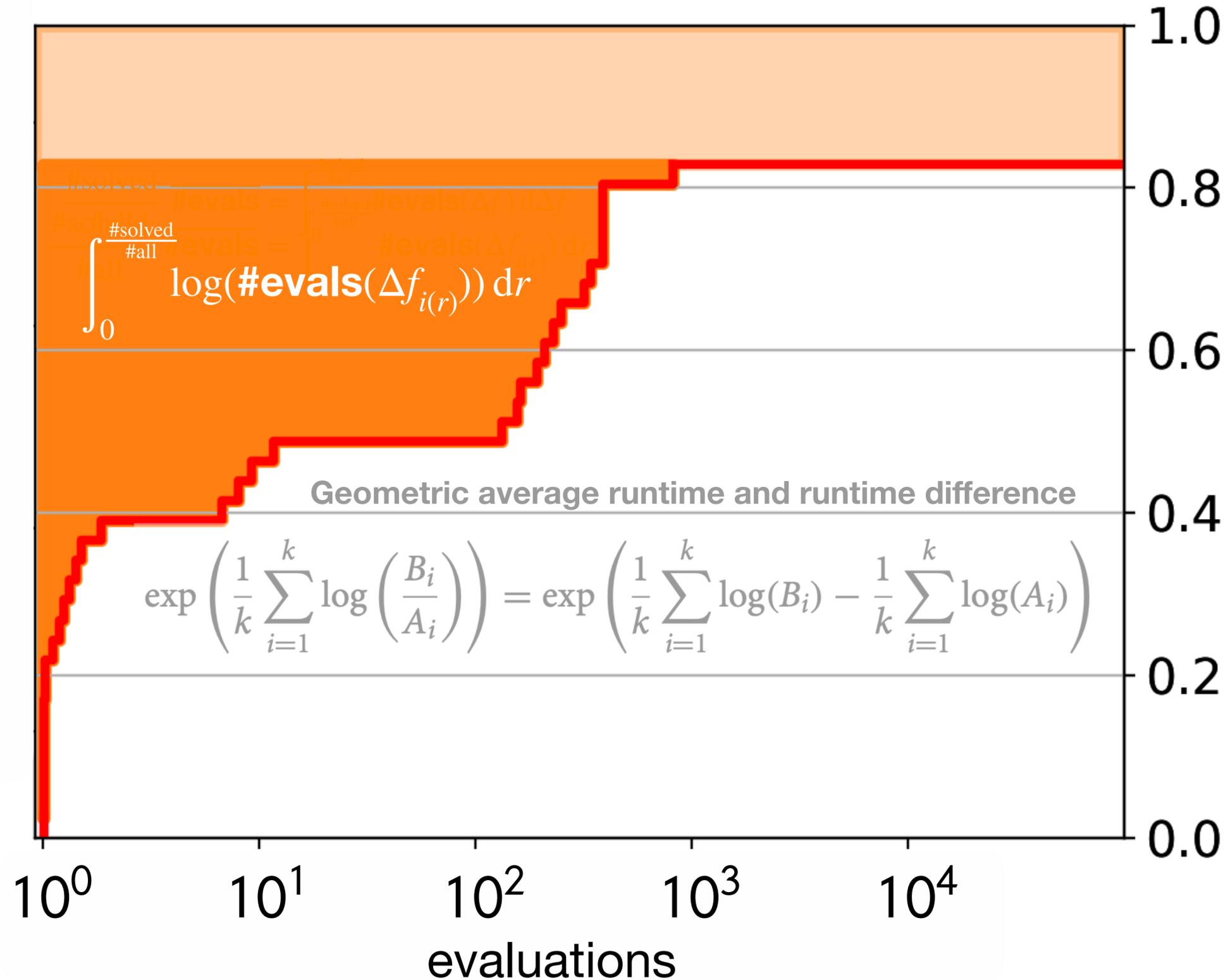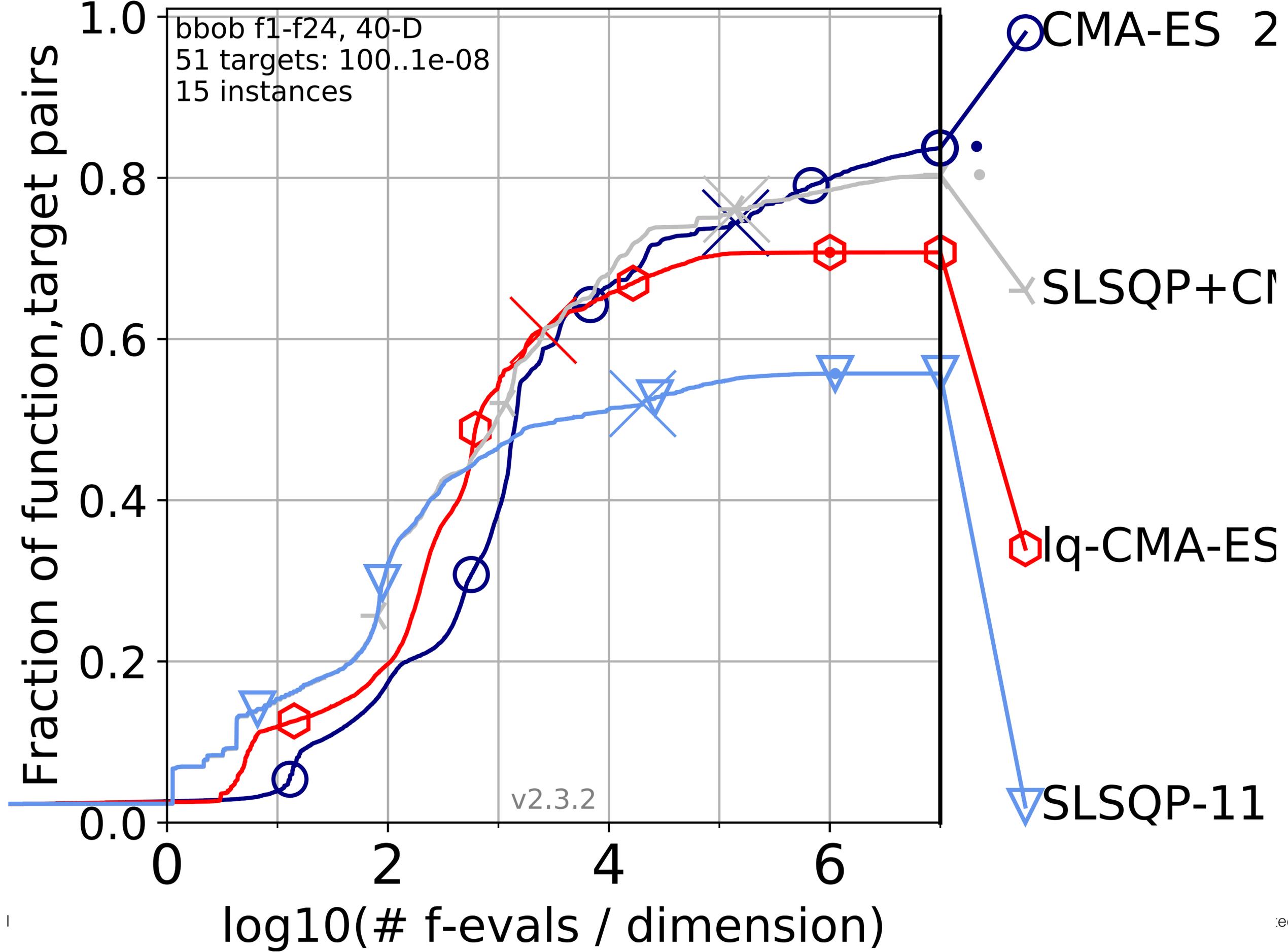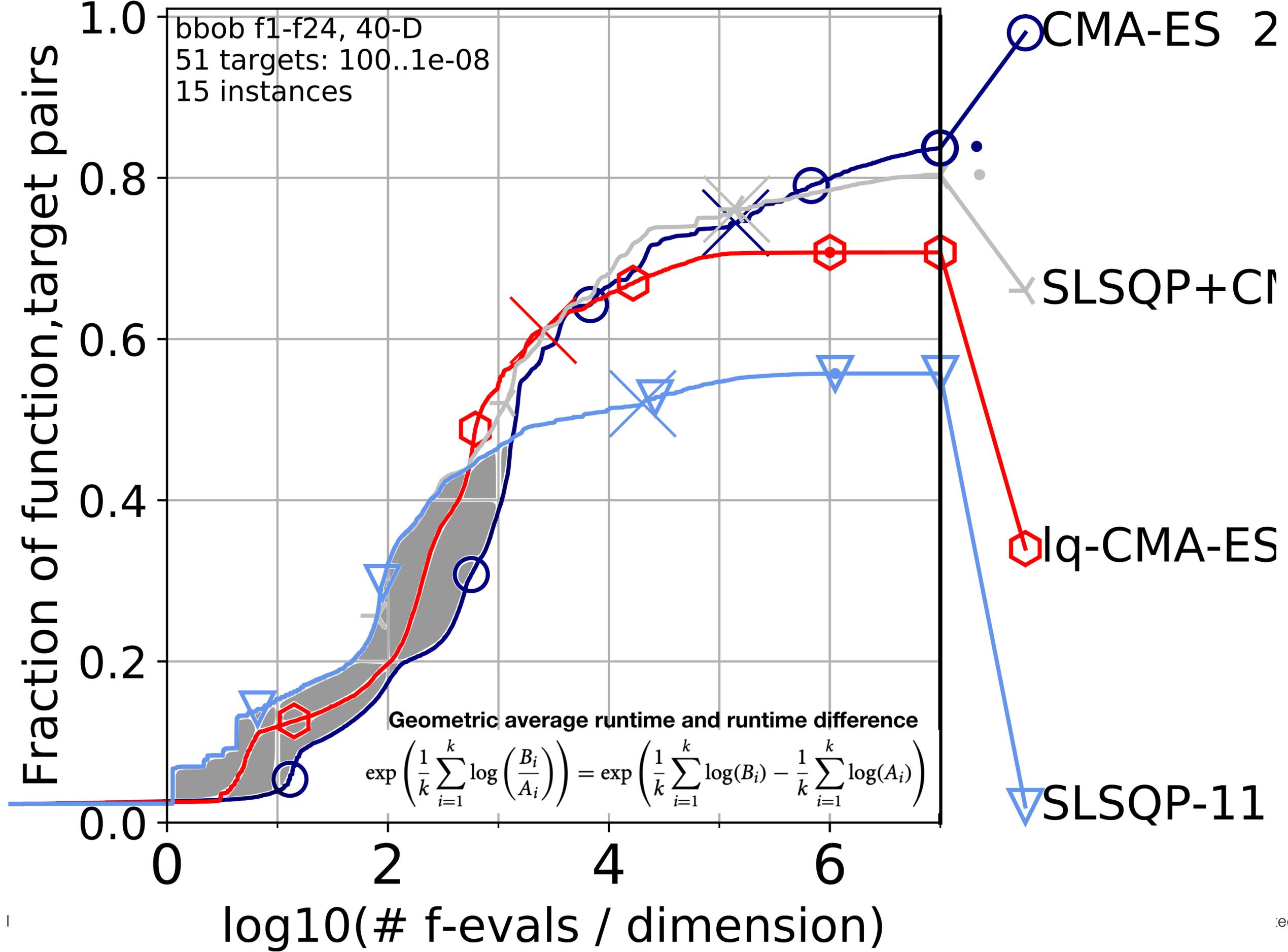
We may choose a dynamic truncation budget for automated benchmarking



$$\int_0^{\frac{\#solved}{\#all}} \log(\textbf{#evals}(\Delta f_{i(r)}))\, dr$$

Geometric average runtime and runtime difference

$$\exp\left(\frac{1}{k}\sum_{i=1}^{k}\log\left(\frac{B_i}{A_i}\right)\right) = \exp\left(\frac{1}{k}\sum_{i=1}^{k}\log(B_i) - \frac{1}{k}\sum_{i=1}^{k}\log(A_i)\right)$$

evaluations

Hansen et al. 2022. Anytime performance assessment in blackbox optimization benchmarking. *IEEE Trans. on EC*, 26(6).

comparing graphs…

Nikolaus Hansen, Inria & Institute Polytechnique de Paris

bbob f1-f24, 40-D
51 targets: 100..1e-08
15 instances

1.0

0.8

0.6

Fraction of function,target pairs

0.4

0.2

0.0

CMA-ES  2

SLSQP+CM

lq-CMA-ES

SLSQP-11

v2.3.2

0        2        4        6

log10(# f-evals / dimension)

difference area…

Nikolaus Hansen, I                                                                    ed benchmarking in mind

bbob f1-f24, 40-D
51 targets: 100..1e-08
15 instances

Geometric average runtime and runtime difference

$$\exp\left(\frac{1}{k}\sum_{i=1}^{k}\log\left(\frac{B_i}{A_i}\right)\right) = \exp\left(\frac{1}{k}\sum_{i=1}^{k}\log(B_i) - \frac{1}{k}\sum_{i=1}^{k}\log(A_i)\right)$$

CMA-ES  2

SLSQP+CM

lq-CMA-ES

SLSQP-11

Fraction of function,target pairs

log10(# f-evals / dimension)

bbob f1-f24, 40-D
51 targets: 100..1e-08
15 instances

CMA-ES  2

SLSQP+CM

lq-CMA-ES

SLSQP-11

Geometric average runtime and runtime difference

$$\exp\left(\frac{1}{k}\sum_{i=1}^{k}\log\left(\frac{B_i}{A_i}\right)\right) = \exp\left(\frac{1}{k}\sum_{i=1}^{k}\log(B_i) - \frac{1}{k}\sum_{i=1}^{k}\log(A_i)\right)$$

Fraction of function,target pairs

log10(# f-evals / dimension)

bbob f1-f24, 40-D
51 targets: 100..1e-08
15 instances

CMA-ES 2

SLSQP+CM

lq-CMA-ES

SLSQP-11

**Geometric average runtime and runtime difference**

$$\exp\left(\frac{1}{k}\sum_{i=1}^{k}\log\left(\frac{B_i}{A_i}\right)\right) = \exp\left(\frac{1}{k}\sum_{i=1}^{k}\log(B_i) - \frac{1}{k}\sum_{i=1}^{k}\log(A_i)\right)$$

Fraction of function,target pairs

log10(# f-evals / dimension)

fixed target vs budget properties...

# COCO/BBOB: Fixed Target(s) versus Fixed Budget



- The fixed budget (vertical) design is (much) easier to set up

  choosing a budget is simpler than choosing a target and we need to chose a maximal "timeout" budget either way

- For the (very) same reason, results from the fixed target (horizontal) design are (much) simpler to interpret and more conclusive

  without specific insight, a function value is impossible to interpret beyond ordering

- Fixed target results

  - can't miss a big fitness change (improvement)

    whereas with budget (vertical) discretization catching a fitness jump timely entirely relies on a fine discretization

  - are "budget-free" even for a single target

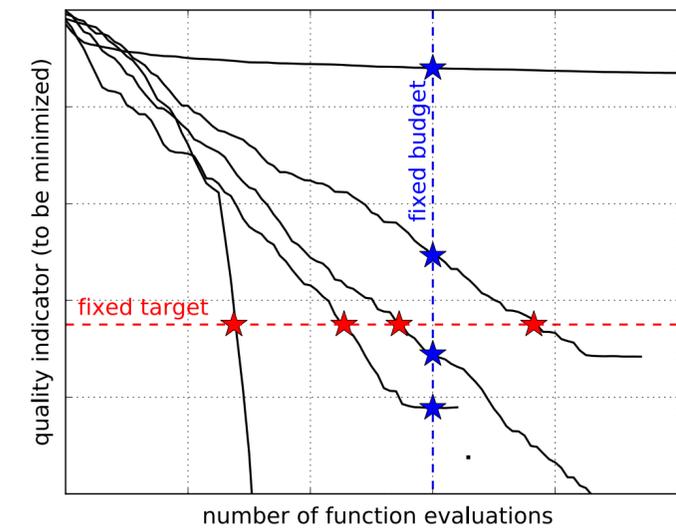    we can compare results from runs with different maximal "timeout" budgets

  - can be meaningfully aggregated in ECDFs and geometric averages

    whereas function values from different functions are not commensurable

  - can give different weights to *different problems* to make the geometric average more representative, as in

$$\sum_i w_i \log(\mathrm{RT}_i)$$

# COCO/BBOB: Fixed Target(s) versus Fixed Budget



- The fixed budget (vertical) design is (much) easier to set up

  choosing a budget is simpler than choosing a target and we need to chose a maximal "timeout" budget either way

- For the (very) same reason, results from the fixed target (horizontal) design are (much) simpler to interpret and more conclusive

  without specific insight, a function value is impossible to interpret beyond ordering

- Fixed target results

  - can't miss a big fitness change (improvement)

    whereas with budget (vertical) discretization catching a fitness jump timely entirely relies on a fine discretization

  - are "budget-free" even for a single target

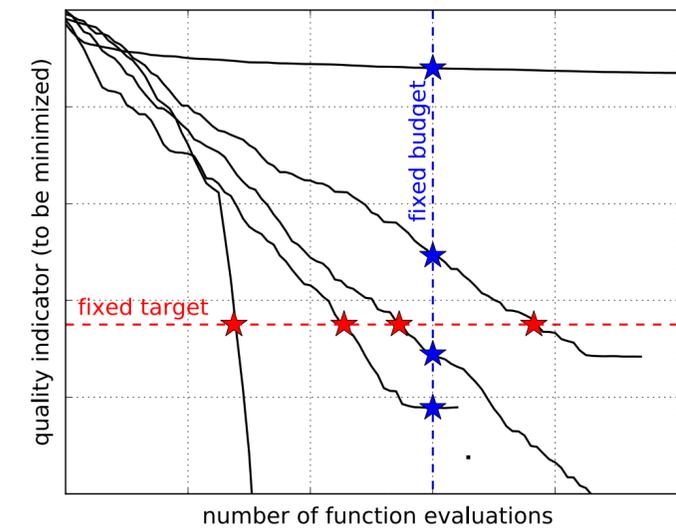    we can compare results from runs with different maximal "timeout" budgets

  - can be meaningfully aggregated in ECDFs and geometric averages

    whereas function values from different functions are not commensurable

  - can give different weights to different problems to make the geometric average more representative, as in

    $$\sum_i w_i \log(\mathrm{RT}_i)$$

# COCO/BBOB: Fixed Target(s) versus Fixed Budget



- The fixed budget (vertical) design is (much) easier to set up

  choosing a budget is simpler than choosing a target and we need to chose a maximal "timeout" budget either way

- For the (very) same reason, results from the fixed target (horizontal) design are (much) simpler to interpret and more conclusive

  without specific insight, a function value is impossible to interpret beyond ordering

- Fixed target results

  - can't miss a big fitness change (improvement)

    whereas with budget (vertical) discretization catching a fitness jump timely entirely relies on a fine discretization

  - are "budget-free" even for a single target

    we can compare results from runs with different maximal "timeout" budgets

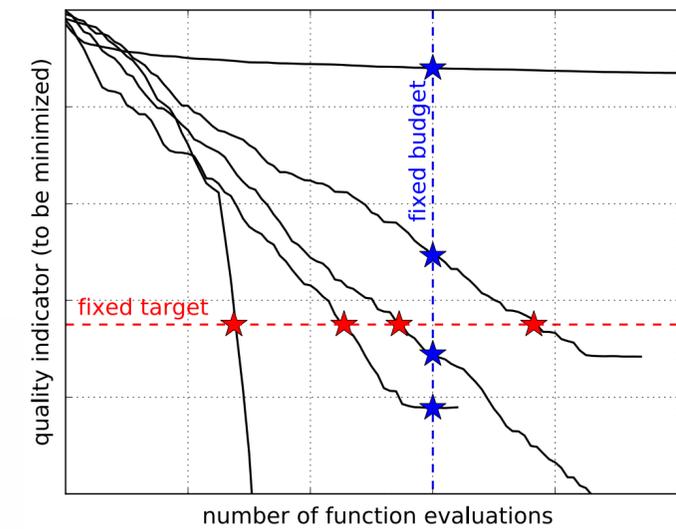  - can be meaningfully aggregated in ECDFs and geometric averages

    whereas function values from different functions are not commensurable

  - can give different weights to different problems to make the geometric average more representative, as in
    $$\sum_i w_i \log(\text{RT}_i)$$

    automated…

# Automated BBOB: Fixed Target(s) versus Fixed Budget



- The fixed budget (vertical) design is (much) easier to set up

  choosing a budget is simpler than choosing a target and we need to chose a maximal "timeout" budget either way

- For the (very) same reason, results from the fixed target (horizontal) design are (much) simpler to interpret and more conclusive

  without specific insight, a function value is impossible to interpret beyond ordering

- Fixed target results

  - can't miss a big fitness change (improvement)

    whereas with budget (vertical) discretization catching a fitness jump timely entirely relies on a fine discretization

  - are "budget-free" even for a single target

    we can compare results from runs with different maximal "timeout" budgets

- can be meaningfully aggregated in ECDFs and geometric averages

  whereas function values from different functions are not commensurable

- can give different weights to *different problems* to make the geometric average more representative, as in
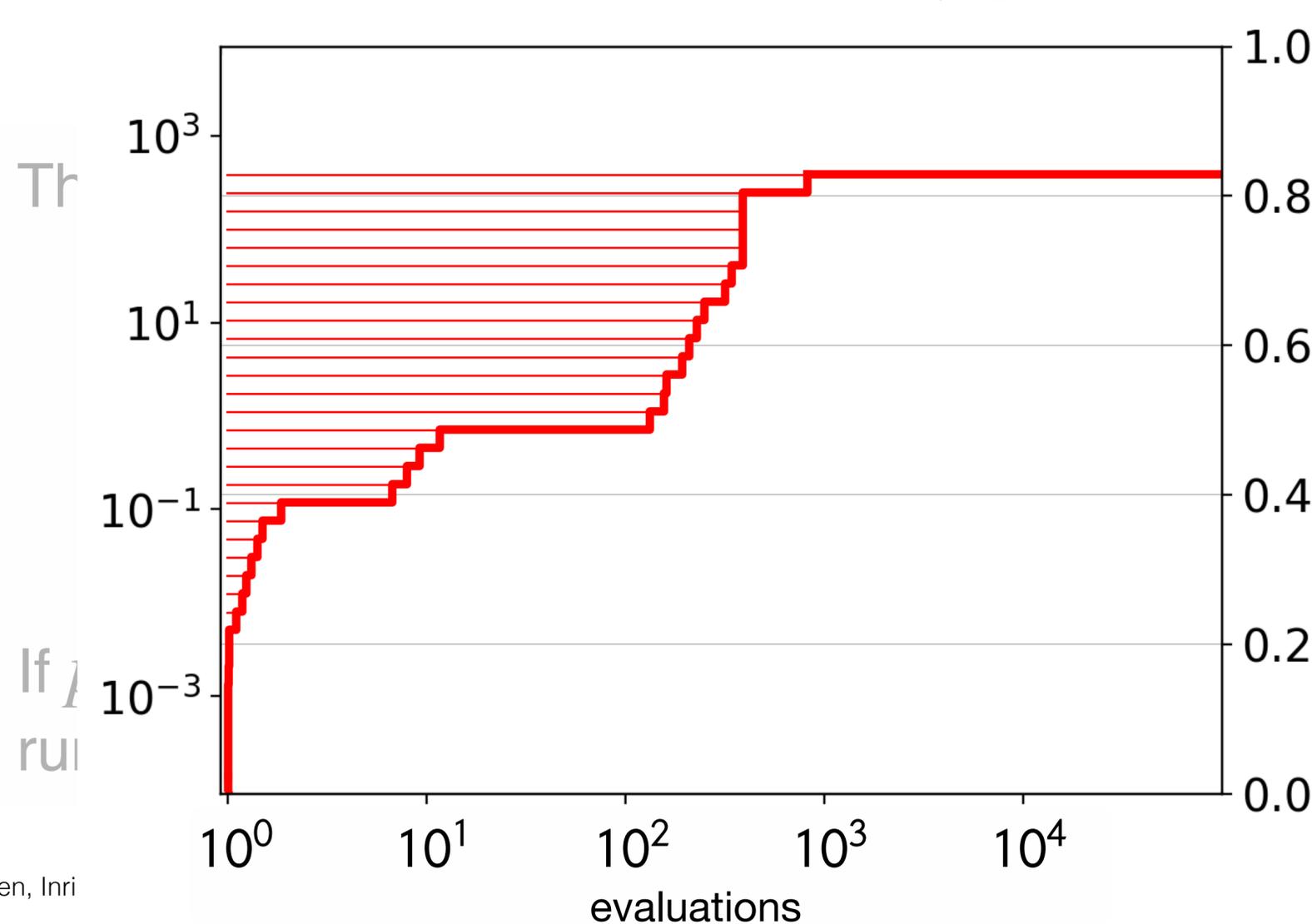
  $$\sum_i w_i \log(\text{RT}_i)$$

  simulated runtimes…

# Converting Positive Success Probabilities into Runtimes

## Solving the fast-versus-successful comparison dilemma

We construct runtimes <span style="color:magenta">with simulated restarts</span>

$$\text{RT}^{\text{sim}} = \text{RT}^{\text{succ}}_j + \sum_{i=1}^{N} \text{RT}^{\text{fail}}_i \approx (1 + N) \times \text{RT}^{\text{average}}$$



runtimes are drawn uniformly at random until the first success, $N \geq 0$ is hence a random variable

Th

If $\lambda$
rui

$$\frac{\text{all runs}}{} = \frac{\sum_i \text{RT}_i}{\# \text{ successes}}$$

runs gives an estimate of a

ERT…

# Converting Positive Success Probabilities into Runtimes

## Solving the fast-versus-successful comparison dilemma

We construct runtimes with simulated restarts

$$\text{RT}^{\text{sim}} = \text{RT}_j^{\text{succ}} + \sum_{i=1}^{N} \text{RT}_i^{\text{fail}} \approx (1 + N) \times \text{RT}^{\text{average}}$$

runtimes are drawn uniformly at random until the first success,
$N \geq 0$ is hence a random variable

Their expected runtime is

$$\text{ERT} = \frac{\text{average evaluations to target of all runs}}{p_{\text{success}}} = \frac{\sum_i \text{RT}_i}{\# \text{ successes}}$$
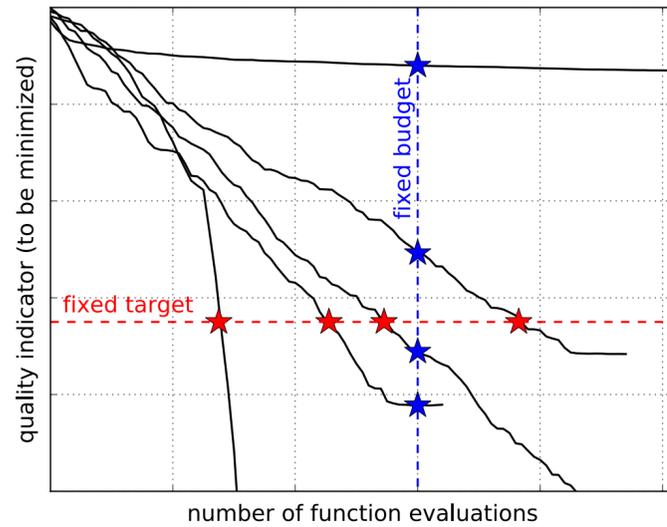
If $p_{\text{success}} = 0$, the inequality $p_{\text{success}} \leq 1/N_{\text{runs}}$ gives an estimate of a runtime lower bound
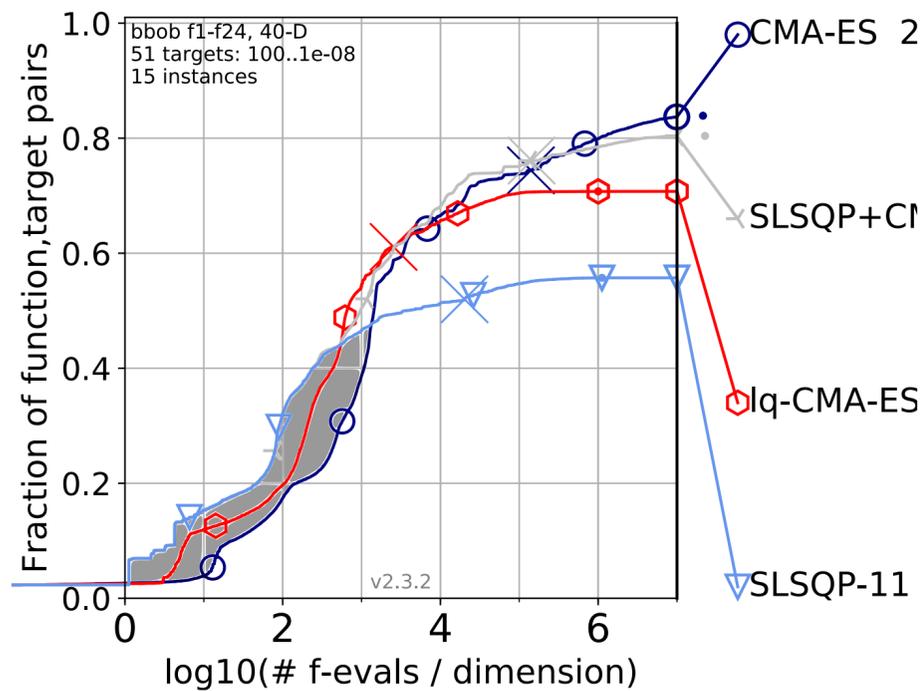
# Thank you

We construct runtimes with simulated restarts

$$RT^{sim} = RT_j^{succ} + \sum_{i=1}^{N} RT_i^{fail} \approx (1 + N) \times RT^{average}$$

runtimes are drawn uniformly at random until the first success,
$N \geq 0$ is hence a random variable

# Time for discussion



HANSEN et al.: ANYTIME PERFORMANCE ASSESSMENT IN BLACKBOX OPTIMIZATION BENCHMARKING                                1297
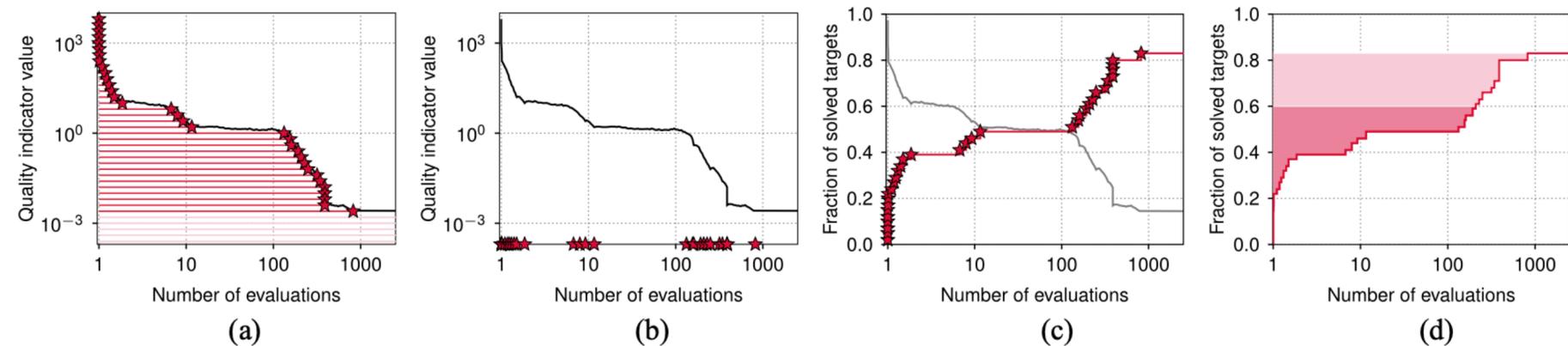


Fig. 2.   Construction of the (anytime) empirical runtime distribution (RTD) from the quality indicator convergence graph and equally spaced target values. The lengths of the horizontal lines in (a) represent the runtime measurements for different target values taken from the quality indicator convergence graph, as do the stars in (a)–(c). When the target values are evenly spaced on the given $y$-axis, the RTD reconstructs, up to discretization, the lower envelope of the convergence graph flipped upside down (c). The areas above the RTD in (d) represent average runtimes, where the logarithm of the geometric average of a fraction of the smallest runtimes equals the area size divided by the trim fraction. The dark area corresponds to the shortest 60%, and dark plus light area correspond to the fraction of all *solved* problems.