# Ten+ Years of Benchmarking with COCO/BBOB

Nikolaus Hansen
Inria
CMAP, CNRS, Ecole Polytechnique, Institut Polytechnique de Paris, France

COCO — Comparing Continuous Optimisers

- is a (software) platform for comparing continuous optimisers in a black-box scenario

  https://github.com/numbbo/coco

- *automatises* the tedious and repetitive task of *benchmarking numerical optimisation algorithms in a black-box setting*

- advantage: saves time and prevents common (and not so common) pitfalls

COCO provides

- experimental and measurement *methodology*

  main decision: what is the end point of measurement

- suites of benchmark functions

  single objective, bi-objective, noisy, mixed-integer, more to come…

- data of already benchmarked algorithms to compare with

# Benchmarking: Related Goals

1. Understanding algorithms

2. Measuring performance in a systematic way (a performance "profile")

3. Running a competition

# Benchmarking: The Global Picture

Two *surprisingly (but not completely) independent* puzzles to solve

- **What to benchmark**: for example, which collection of test problems?

- **How to assess performance**?

  - experimental setup

  - data collection

  - measures used and presented
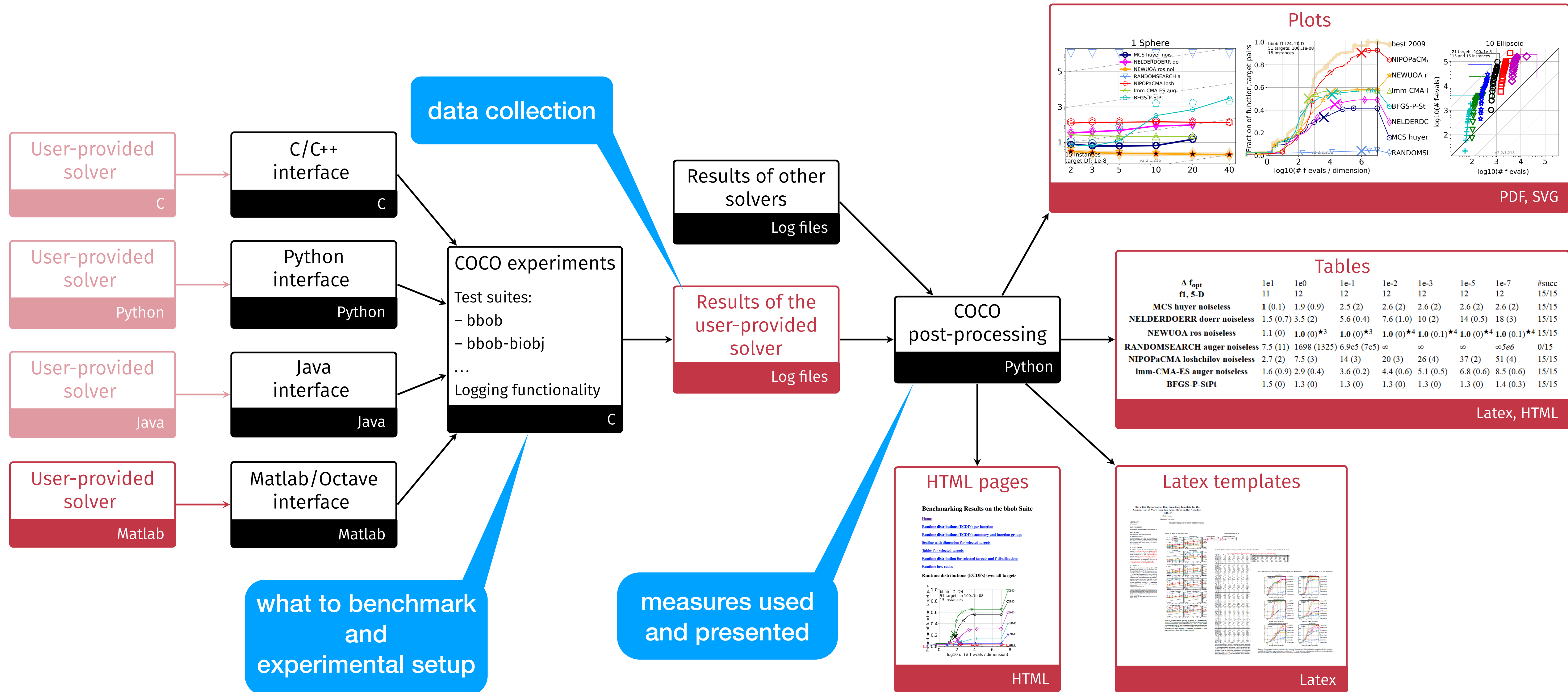
# COCO/BBOB: The Global Picture

**data collection**

**what to benchmark and experimental setup**

**measures used and presented**

User-provided solver — C
C/C++ interface — C

User-provided solver — Python
Python interface — Python

User-provided solver — Java
Java interface — Java

User-provided solver — Matlab
Matlab/Octave interface — Matlab

COCO experiments

Test suites:
– bbob
– bbob-biobj
…
Logging functionality — C

Results of other solvers — Log files

Results of the user-provided solver — Log files

COCO post-processing — Python

### Plots

1 Sphere

- MCS huyer nois
- NELDERDOERR do
- NEWUOA ros noi
- RANDOMSEARCH a
- NIPOPaCMA losh
- lmm-CMA-ES aug
- BFGS-P-StPt

15 instances
target Df: 1e-8

best 2009
- NIPOPaCM,
- NEWUOA r
- lmm-CMA-I
- BFGS-P-St
- NELDERDC
- MCS huyer
- RANDOMSI

10 Ellipsoid

PDF, SVG

### Tables

| Δ $f_{opt}$ | 1e1 | 1e0 | 1e-1 | 1e-2 | 1e-3 | 1e-5 | 1e-7 | #succ |
|---|---|---|---|---|---|---|---|---|
| f1, 5-D | 11 | 12 | 12 | 12 | 12 | 12 | 12 | 15/15 |
| MCS huyer noiseless | **1** (0.1) | 1.9 (0.9) | 2.5 (2) | 2.6 (2) | 2.6 (2) | 2.6 (2) | 2.6 (2) | 15/15 |
| NELDERDOERR doerr noiseless | 1.5 (0.7) | 3.5 (2) | 5.6 (0.4) | 7.6 (1.0) | 10 (2) | 14 (0.5) | 18 (3) | 15/15 |
| NEWUOA ros noiseless | 1.1 (0) | **1.0** (0)★3 | **1.0** (0)★3 | **1.0** (0)★4 | **1.0** (0.1)★4 | **1.0** (0)★4 | **1.0** (0.1)★4 | 15/15 |
| RANDOMSEARCH auger noiseless | 7.5 (11) | 1698 (1325) | 6.9e5 (7e5) | ∞ | ∞ | ∞ | ∞5e6 | 0/15 |
| NIPOPaCMA loshchilov noiseless | 2.7 (2) | 7.5 (3) | 14 (3) | 20 (3) | 26 (4) | 37 (2) | 51 (4) | 15/15 |
| lmm-CMA-ES auger noiseless | 1.6 (0.9) | 2.9 (0.4) | 3.6 (0.2) | 4.4 (0.6) | 5.1 (0.5) | 6.8 (0.6) | 8.5 (0.6) | 15/15 |
| BFGS-P-StPt | 1.5 (0) | 1.3 (0) | 1.3 (0) | 1.3 (0) | 1.3 (0) | 1.3 (0) | 1.4 (0.3) | 15/15 |

Latex, HTML

### HTML pages

Benchmarking Results on the bbob Suite

Home
Runtime distributions (ECDFs) per function
Runtime distributions (ECDFs) summary and function groups
Scaling with dimension for selected targets
Tables for selected targets
Runtime distribution for selected targets and f-distributions
Runtime loss ratios

Runtime distributions (ECDFs) over all targets

HTML

### Latex templates

Latex

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.

…feel free to ask questions…

# COCO/BBOB: Test Suite(s)

- Functions are

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

  - Comprehensible

  - 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

- The bad

  - Rastrigin function type is somewhat overrepresented

    partly due to function pairing

  - 10% of the default targets for F23 Katsuuras are trivial to hit

    evaluating the domain middle at first is a good "algorithm"

- Require to define target values (function + 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

# Data Format

with hindsight 20/20

- The good:

  - scattered experiments can be "merged" (and "unmerged") with a single "drag-and-drop"

  - separation between .info (meta- and summary-data) and .dat files is helpful

  - 10+ years old data are still smoothly usable

  - backwards compatible adjustments are/were possible

- The bad:

  - slightly too few targets (too coarse discretization, not a *format* issue though with backward compatible fix)

  - "handling" of restarts is suboptimal

  - meta-data are not json-style (key-value Python-dict-style) formatted

  - COCO maintains/writes *two* somewhat incompatible formats

# …feel free to ask questions…

# "quality indicator" versus "time" convergence graphs

is all we have (and all we use)

# Specifically

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

  is invariant under changes of computer hardware, OS, programming language, compiler, …

- **quality indicator**:

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

    different for each instance

  - MO: 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

# Convergence Graphs is All We Have



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

we only use the lower envelope

# 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

- evaluations are the independent variable

  function value is the dependent variable, the measurement

# Discretization: Two Possibilities



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

- **horizontal**: not a "natural" discretization

  we need to determine discretization intervals

- function "target" values are the independent variable

  time is the dependent variable, the measurement

- still recovers the original data

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

# Runtime distribution from a single graph

the ECDF recovers the monotonous graph

AKA runtime distribution

the ECDF recovers
the monotonous
graph, discretised
and flipped

AKA runtime distribution

the ECDF recovers
the monotonous
graph, discretised
and flipped

- recovering the convergence graph from discretized data

- collecting runtimes from a single experiments as ECDF

are two interpretations of the same thing
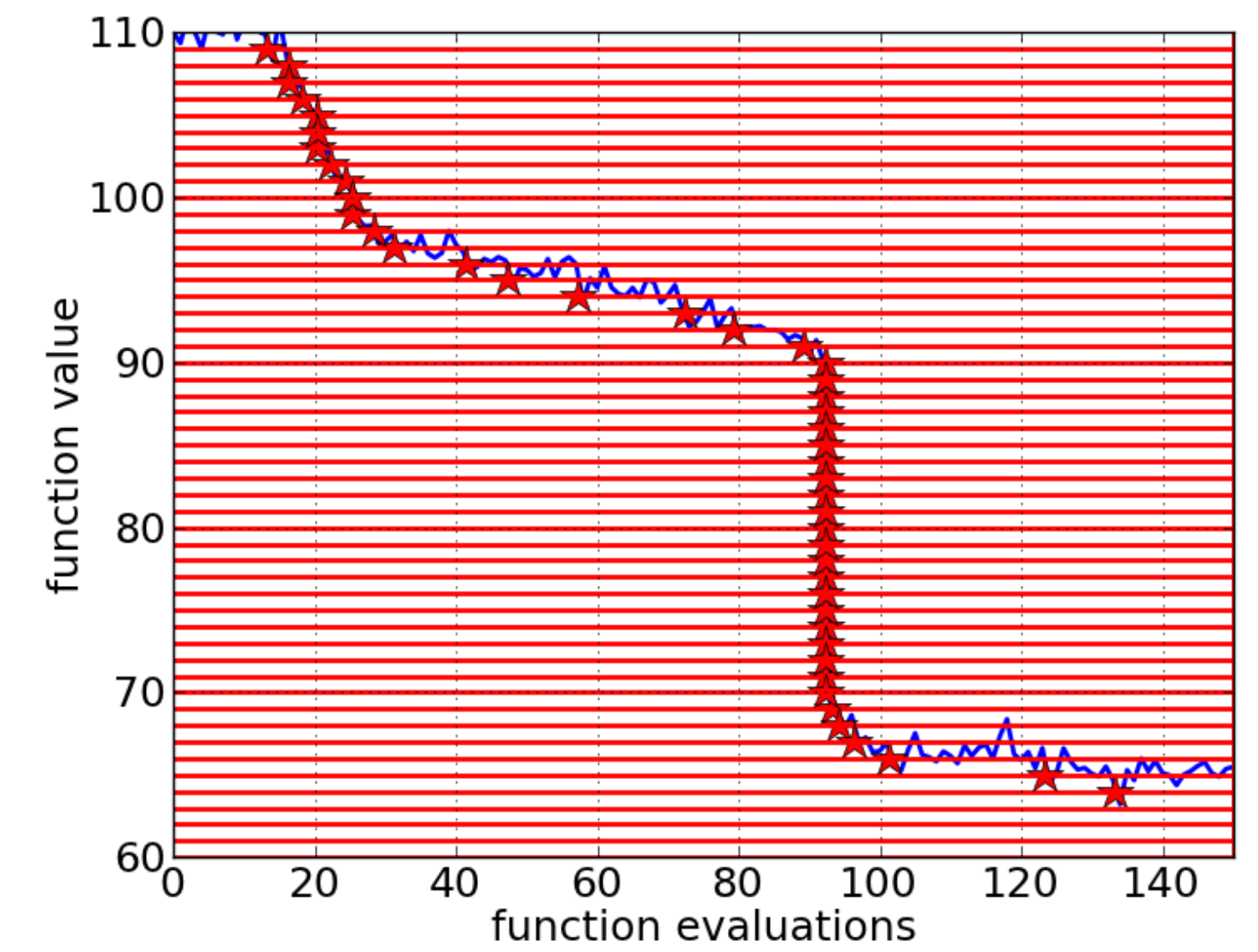
# Runtime distribution from a single graph



the ECDF recovers the monotonous graph, discretised and flipped

the area over the ECDF curve is the average runtime (the geometric average if the x-axis is in log scale)

# COCO/BBOB



uses only
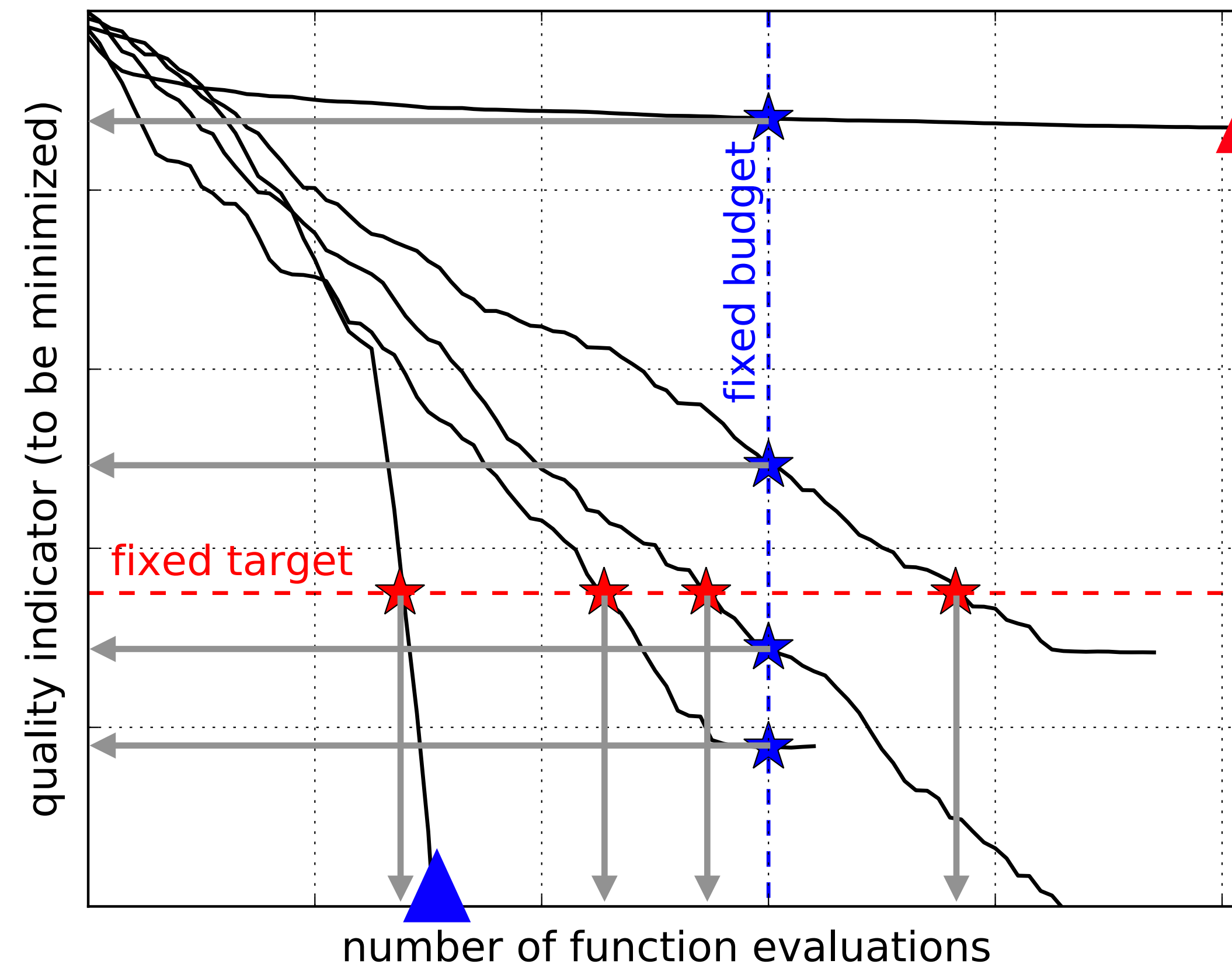
# horizontal discretization

# COCO/BBOB

this is

# not

just

# a technical subtlety

because it crucially determines what measurement we are looking at in the end

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



- five convergence graphs "quality indicator" versus "time"

- **Leads to _different imprecise data_ in both cases**

  - <span style="color:red">"too" bad</span> performance
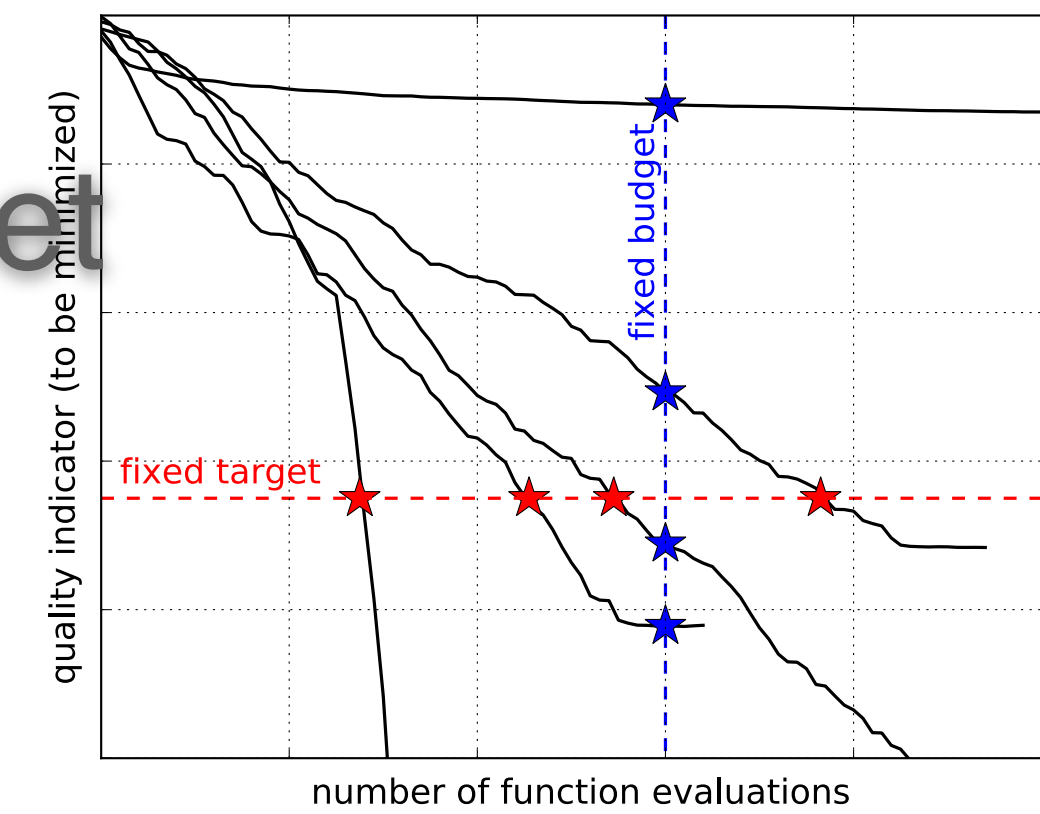
    then the data only provide a lower bound estimate for the runtime (and a fixed budget measure at maximum budget)

  - <span style="color:blue">"too" good</span> performance

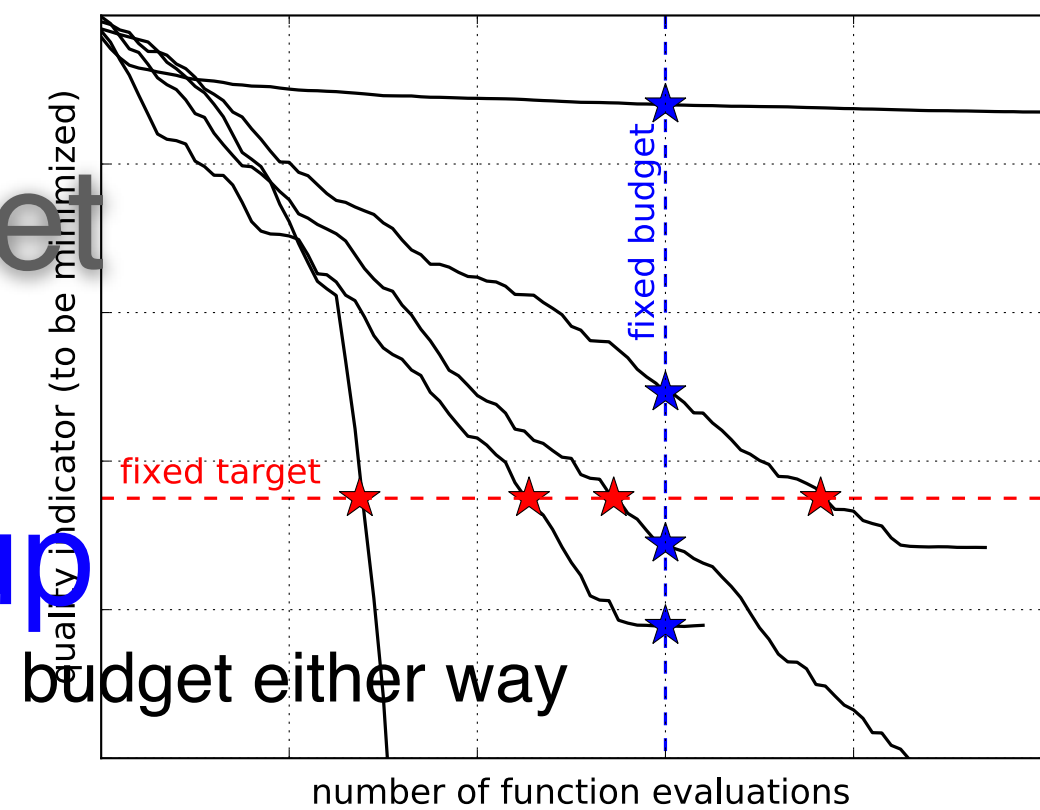    (reached global optimum up to the relevant or numerical precision before the given budget)

## The resulting measurement

- Fixed budget (vertical) design: function values

- Fixed target design: evaluations

- 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 results are (much) simpler to interpret and more conclusive

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

- Fixed target results are "budget-free"

  we can compare results run with different maximal "timeout" budgets

- Fixed target results can be meaningfully aggregated in ECDFs and geometric averages

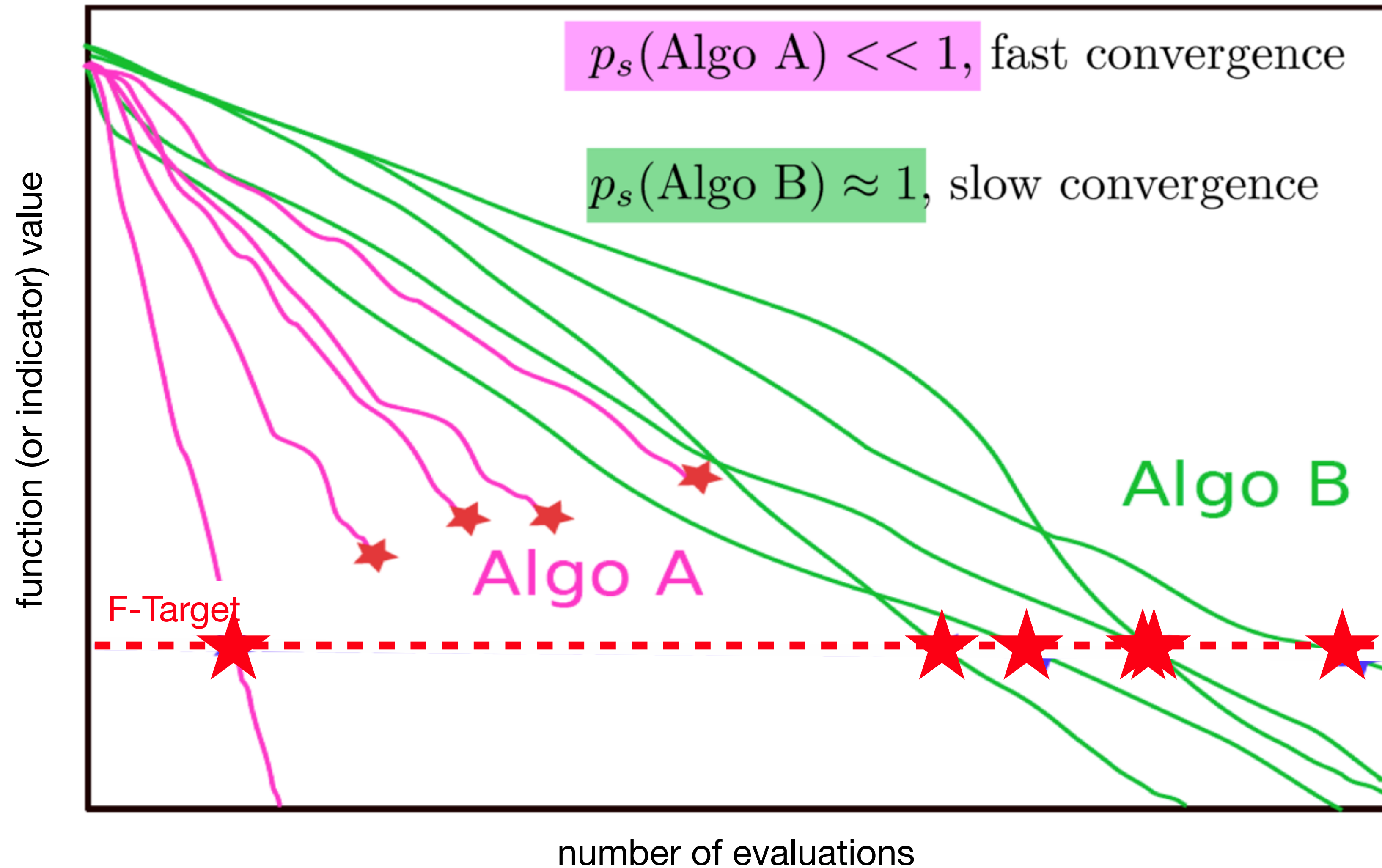  whereas function values from different functions are not commensurable

# Scales of Measurement ("Quality" of Data)

- Nominal - categorial, define a classification

- Ordinal - define an order, ranks, function values (fixed budget)

- Interval - differences are meaningful

- Rational - ratios are meaningful, we can take the logarithm, time (function evaluations, fixed target)

CAVEAT: mathematical and semantic treatment of data is not the same. From a classification with values {1, 2} we can *mathematically* take differences and ratios of the values, but they have no meaningful *semantic interpretation*.

# Treating Success Probabilities

## Solving the fast-versus-successful comparison dilemma



$p_s(\text{Algo A}) << 1$, fast convergence

$p_s(\text{Algo B}) \approx 1$, slow convergence

Algo B

Algo A

F-Target

function (or indicator) value

number of evaluations

# Treating Success Probabilities

## Solving the fast-versus-successful comparison dilemma

Short answer: consider as runtime

$$\frac{\text{something}}{p\text{success}}$$

that is, roughly,

$$\text{runtime} \propto \frac{1}{p\text{success}}$$

# Treating Success Probabilities

## Solving the fast-versus-successful comparison dilemma

We can **simulate a runtime distribution** by simulated (artificial) restarts using the given independent runs



Algo Restart A:

$$p_s(\text{Algo Restart A}) = 1$$

Algo Restart B:

$$p_s(\text{Algo Restart B}) = 1$$

Caveat: the performance of algorithm A critically depends on termination methods (before to hit the target)

**which reflects the situation on a practical problem** unless many runs can be done in parallel

# Treating Success Probabilities
## Solving the fast-versus-successful comparison dilemma

Replacing the success probability with the expected runtime (ERT, aka Enes, SP2, aRT) to hit a target value in #evaluations is computed (estimated) as:

$$\mathrm{ERT} = \frac{\#\text{evaluations(until to hit the target)}}{\#\text{successes}}$$

unsuccessful runs count (only) in the nominator

$$= \mathrm{avg}(\mathrm{evals_{succ}}) + \overbrace{\frac{N_\mathrm{unsucc}}{N_\mathrm{succ}}}^{\text{odds ratio}} \times \mathrm{avg}(\mathrm{evals_{unsucc}})$$

$$\approx \mathrm{avg}(\mathrm{evals_{succ}}) + \frac{N_\mathrm{unsucc}}{N_\mathrm{succ}} \times \mathrm{avg}(\mathrm{evals_{succ}})$$

$$= \frac{N_\mathrm{succ} + N_\mathrm{unsucc}}{N_\mathrm{succ}} \times \mathrm{avg}(\mathrm{evals_{succ}})$$

$$= \frac{1}{\text{success rate}} \times \mathrm{avg}(\mathrm{evals_{succ}})$$

defined (only) for #successes > 0. The last three lines are AKA Q-measure or SP1 (success performance).

See [Price 1997] and [Auger&Hansen 2005]

# Data Sets and Usage Statistics

**Table 1. Visibility of COCO.** All citations as of November 19, 2019, in Google Scholar.

| | | |
|---|---|---|
| Data sets online | `bbob` suite | 227 |
| | `bbob-noisy` suite | 45 |
| | `bbob-biobj` suite | 32 |
| | `bbob-largescale` suite | 11 |
| | `bbob-mixint` suite | 4 |
| BBOB workshop papers using COCO | | 143 |
| Unique authors on the workshop papers | | 109 from 28 countries |
| Papers in Google Scholar found with the search phrase *"comparing continuous optimizers"* OR *"black-box optimization benchmarking (BBOB)"* | | 559 |
| Citations to the COCO documentation including | | 1,455 |

Any `cocopp.archiving.create(folder)`-ed data sets provided under an URL can be loaded with `av = cocopp.archiving.get(URL)` and used in the data processing. See [Hansen et al 2020].

30