# BBOB Black-Box Optimization Benchmarking with CoCO (Comparing Continuous Optimizers)

## The Turbo-Intro

# Black-Box Optimization (Search)

Minimize (or maximize) a continuous domain objective (cost, loss, error, fitness) function

$$f : \mathbb{R}^d \to \mathbb{R}$$

in a black-box scenario (direct search)

$$x \longrightarrow \blacksquare \longrightarrow f(x)$$

where

- gradients are not available or useful

- problem specific knowledge is used only *within* the black box, e.g. with an appropriate encoding
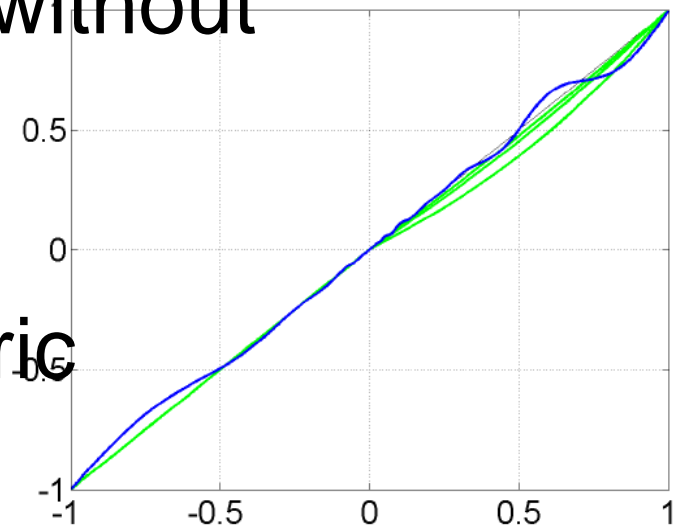
The search costs are the number of function evaluations

# CoCO: the noiseless functions

24 functions within five sub-groups

- Separable functions

- Essential unimodal functions

- Ill-conditioned unimodal functions

- Multimodal structured functions

- Multimodal functions with weak or without structure

functions are not perfectly symmetric
   and are locally deformed

# CoCO: the noisy functions

three noise-"models", so-called:

- Gauss, Uniform (severe), Cauchy (outliers)
- Utility-free noise

$$E(f(x)) \leq E(f(y)) \Rightarrow U(f(x)) \leq U(f(y)) \; \forall x, y, U$$

30 functions with three sub-groups
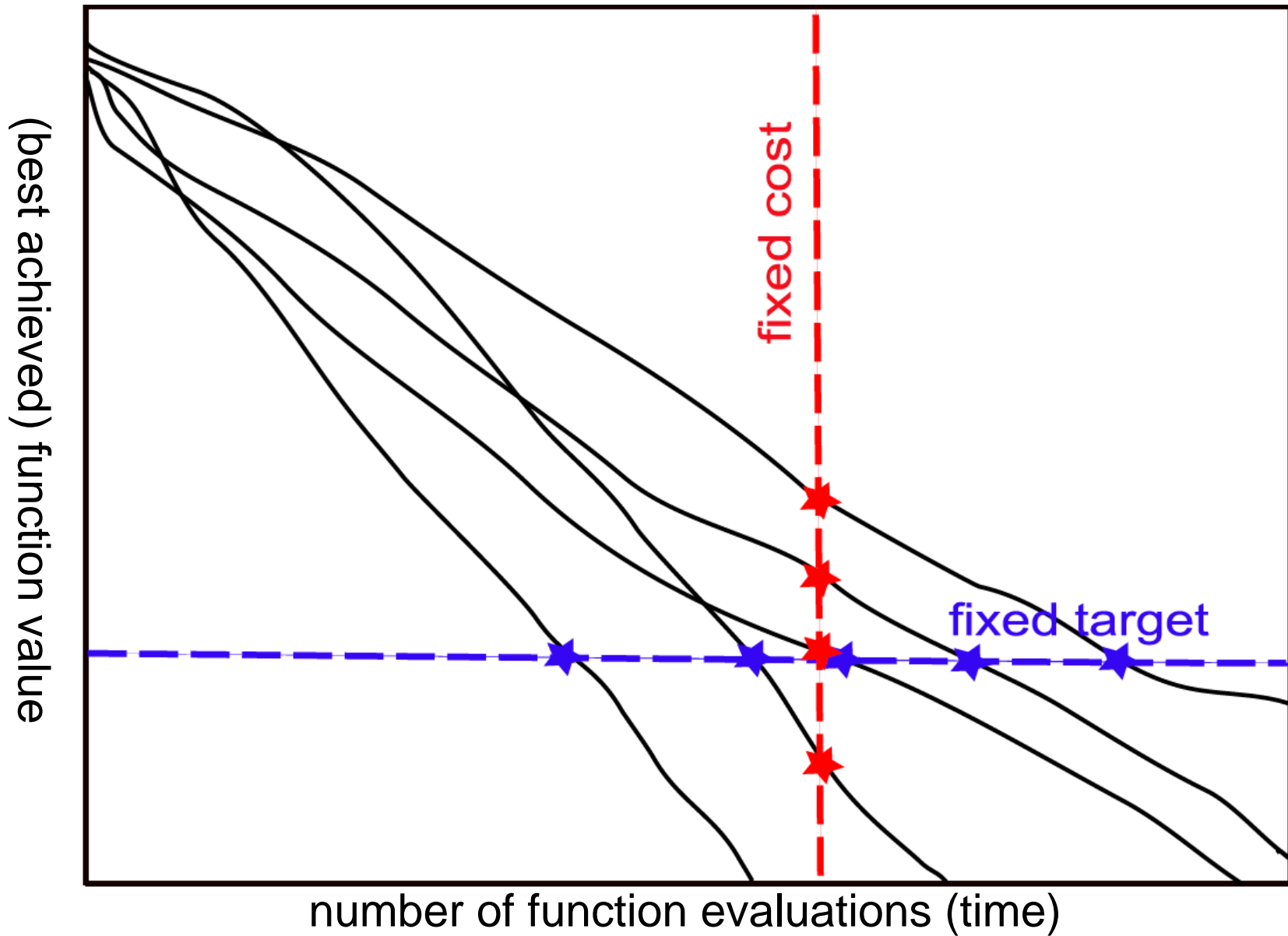
- 2x3 functions with weak noise
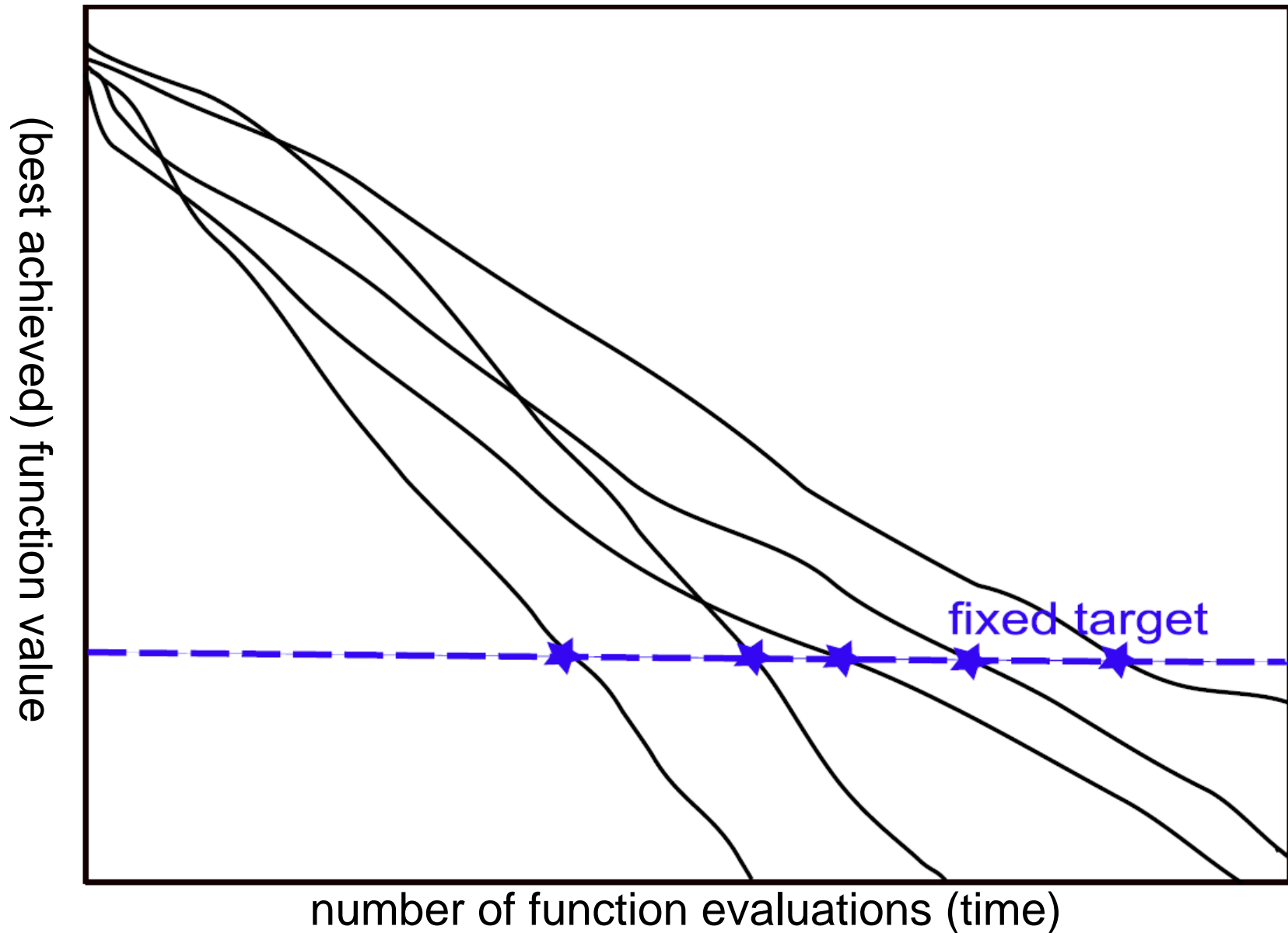- 5x3 unimodal functions
- 3x3 multimodal functions

convergence graphs is
all we have to start with
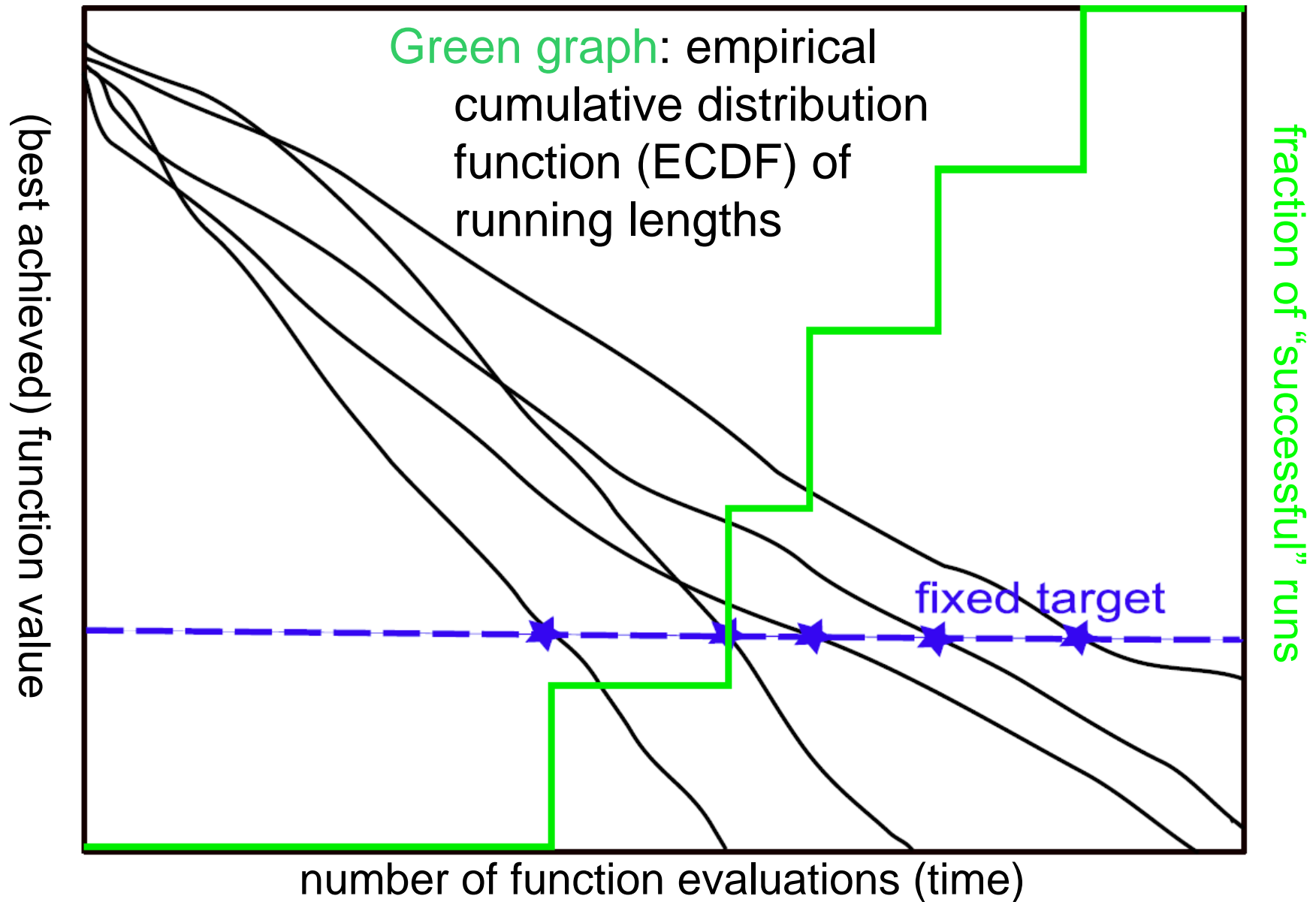
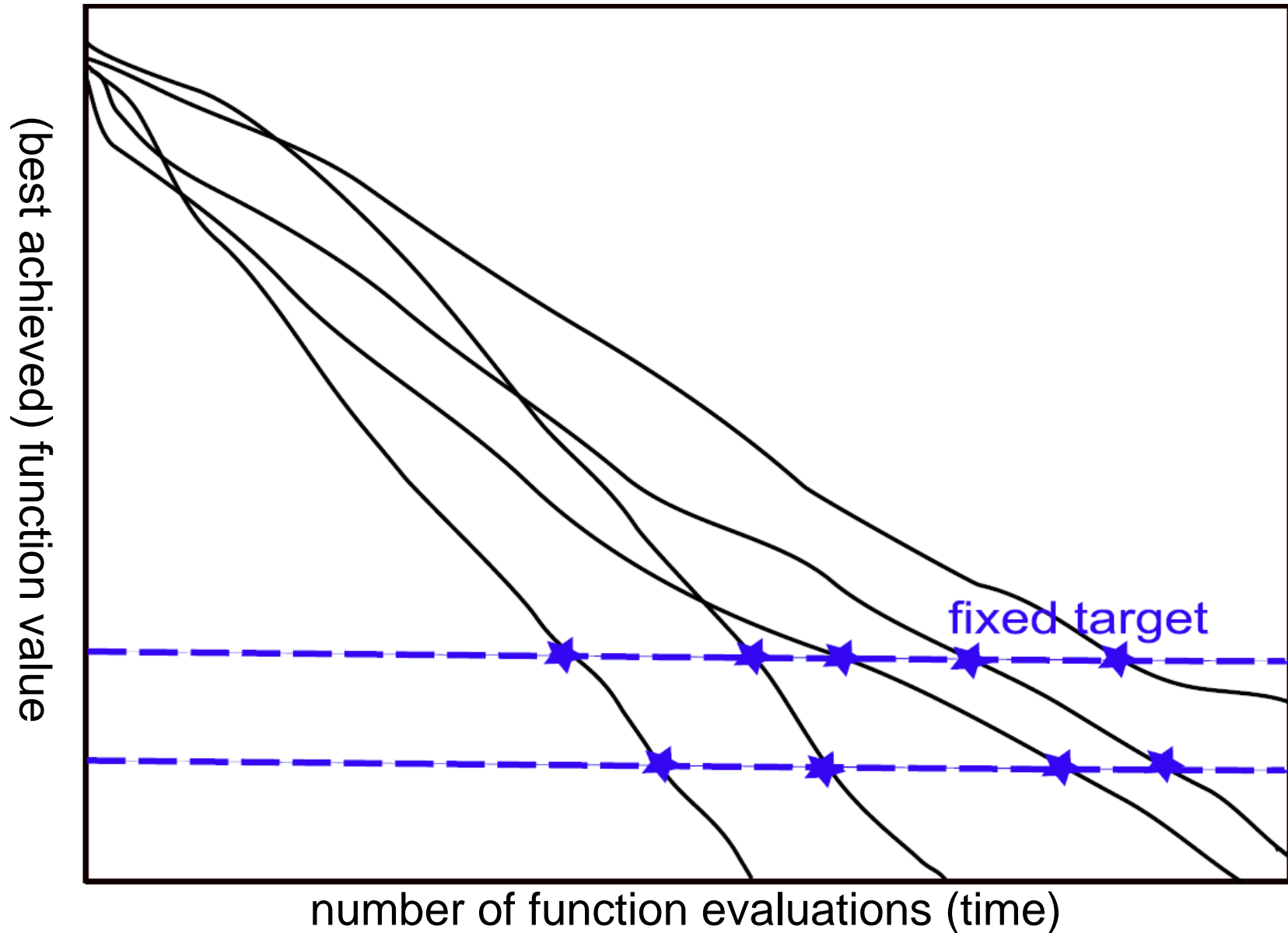# Measuring Performance from Convergence Graphs
# fixed-cost versus fixed-target

# Empirical Cumulative Distribution
## with a given target value



(best achieved) function value

number of function evaluations (time)

fixed target

# Empirical Cumulative Distribution with a <span style="color:blue">given target value</span>



**Green graph**: empirical cumulative distribution function (ECDF) of running lengths

fixed target

(best achieved) function value

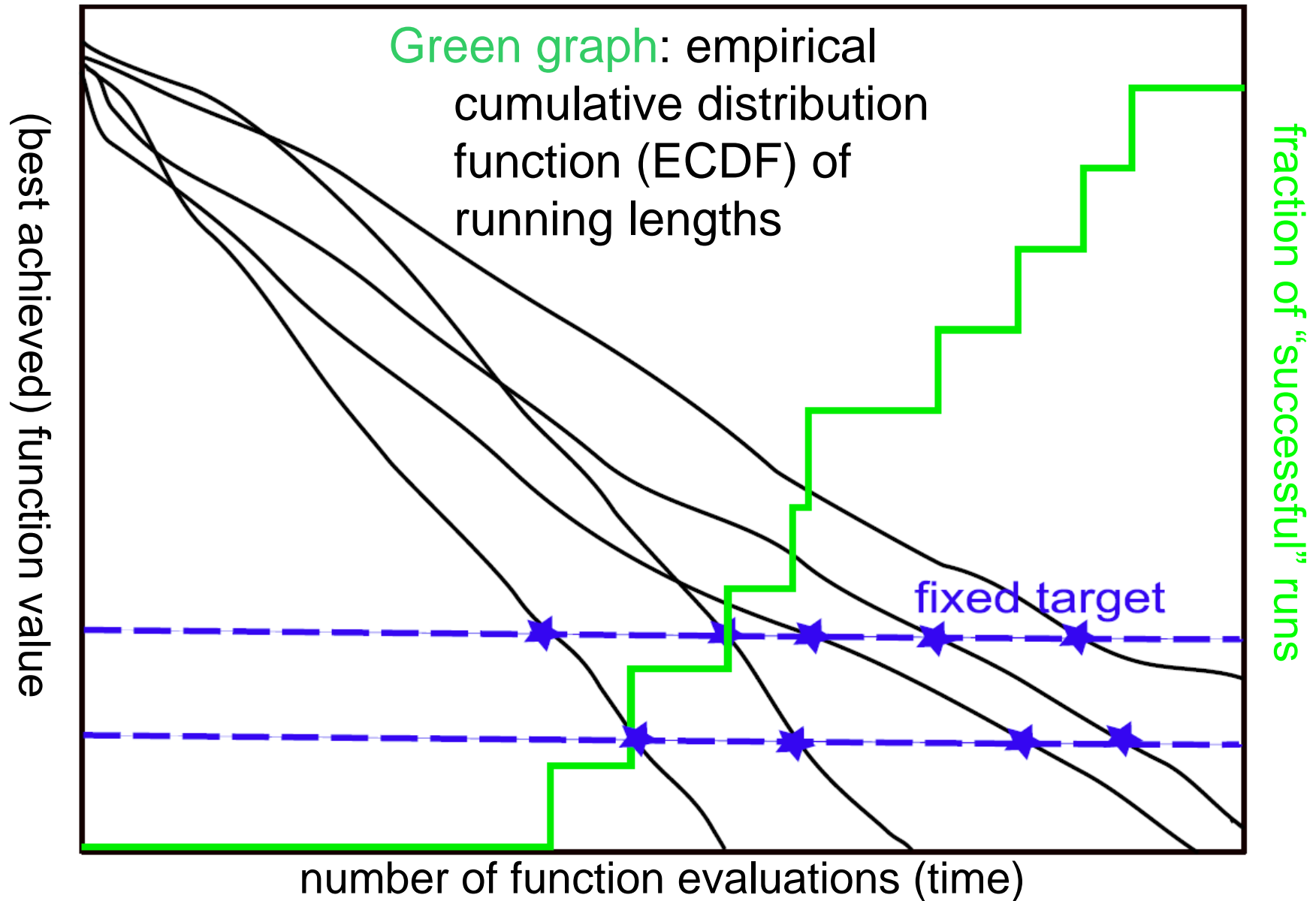number of function evaluations (time)

fraction of "successful" runs

# Empirical Cumulative Distribution with two given target values

# Empirical Cumulative Distribution with two given target values



(best achieved) function value

Green graph: empirical cumulative distribution function (ECDF) of running lengths

fixed target

fraction of "successful" runs

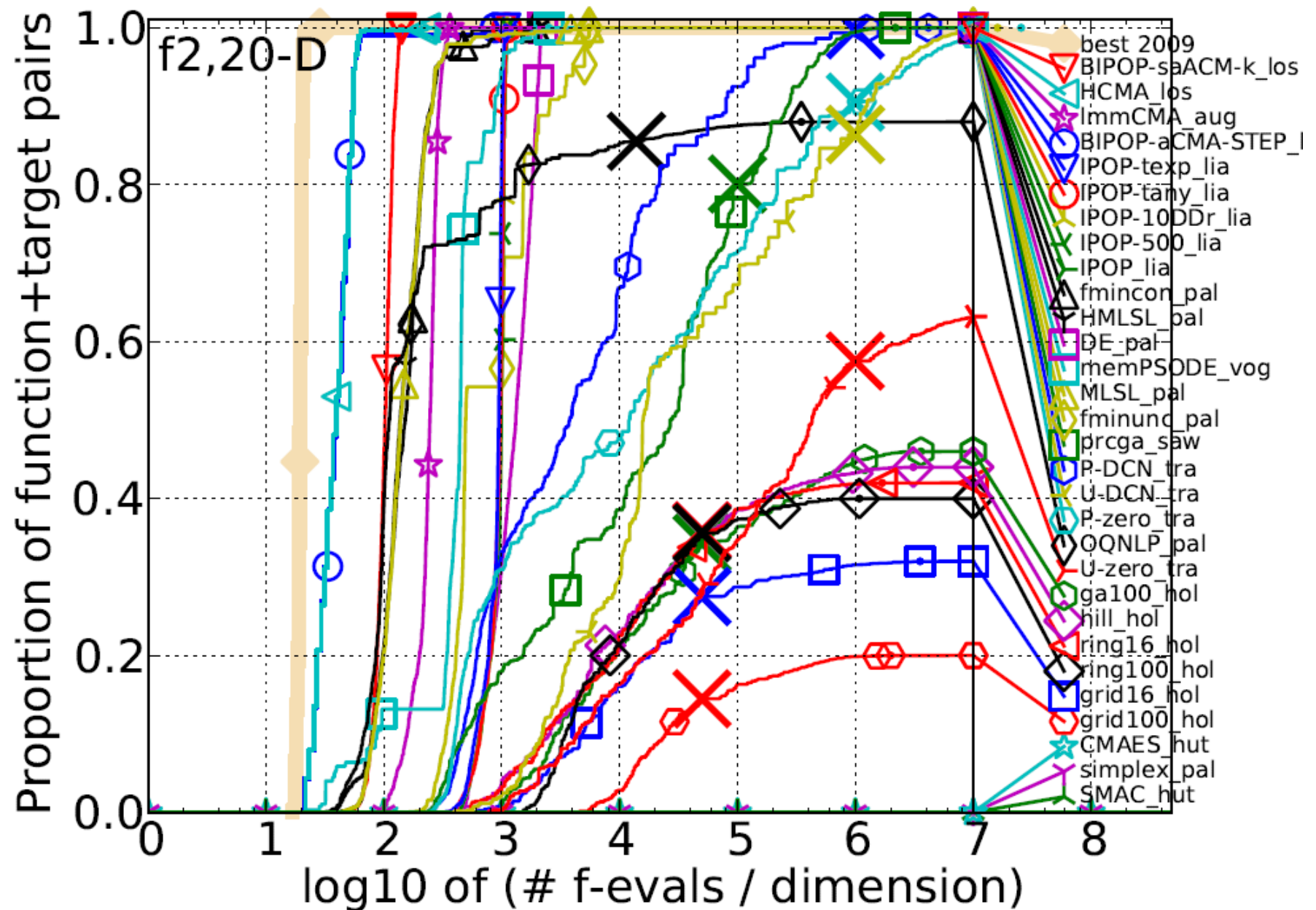number of function evaluations (time)

# Cumulative Distribution of Runtimes

Runtime ECDFs (empirical cumulative distribution function) display a set of runlengths

- they can aggregate over any set of functions and target values

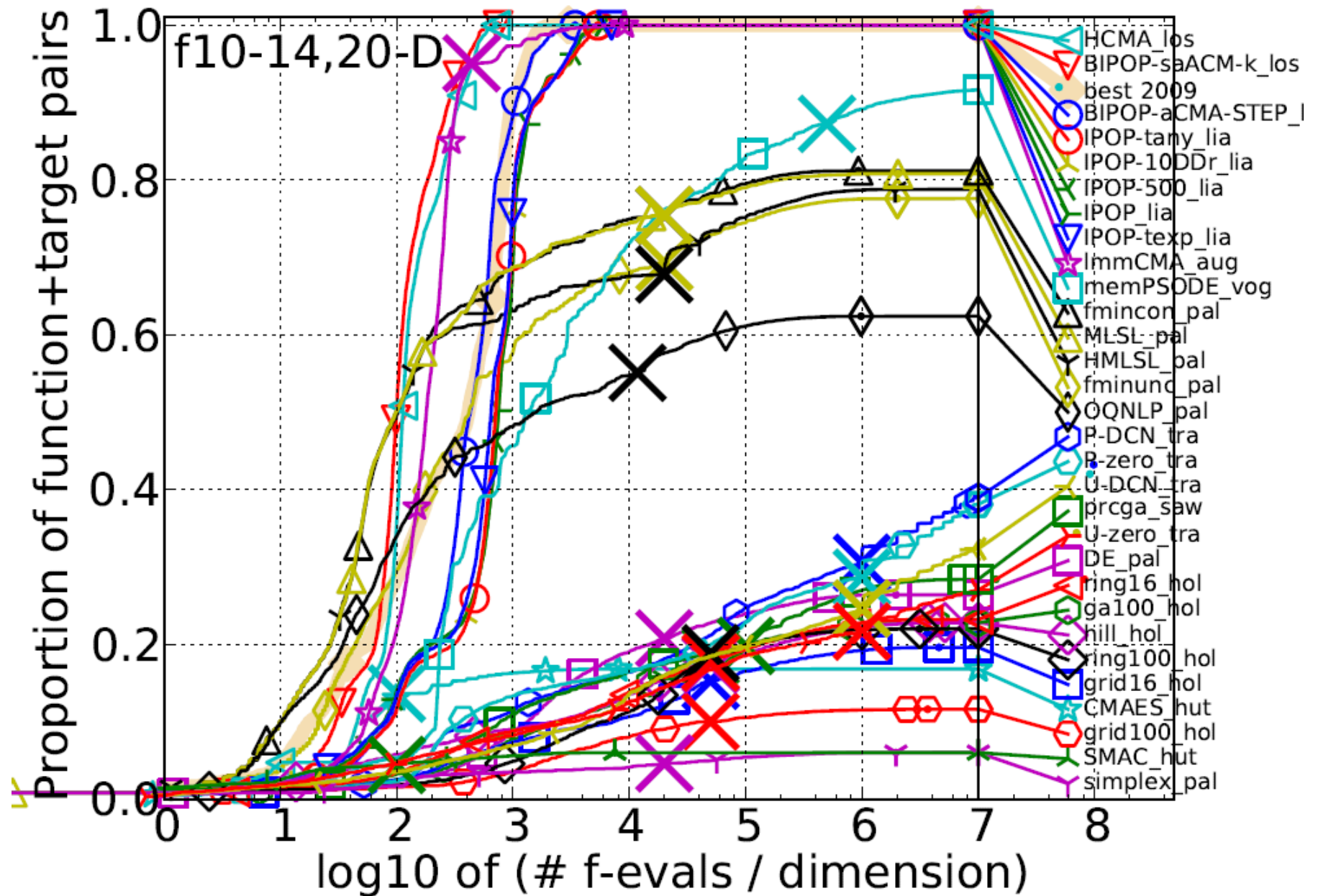  with the least amount of information loss into a single curve

- in BBOB:

  - 50 target values (log-uniform in [1e-8,100]) and 15 trials per function = 750 runlength values per function

  - aggregate of one to 30 functions

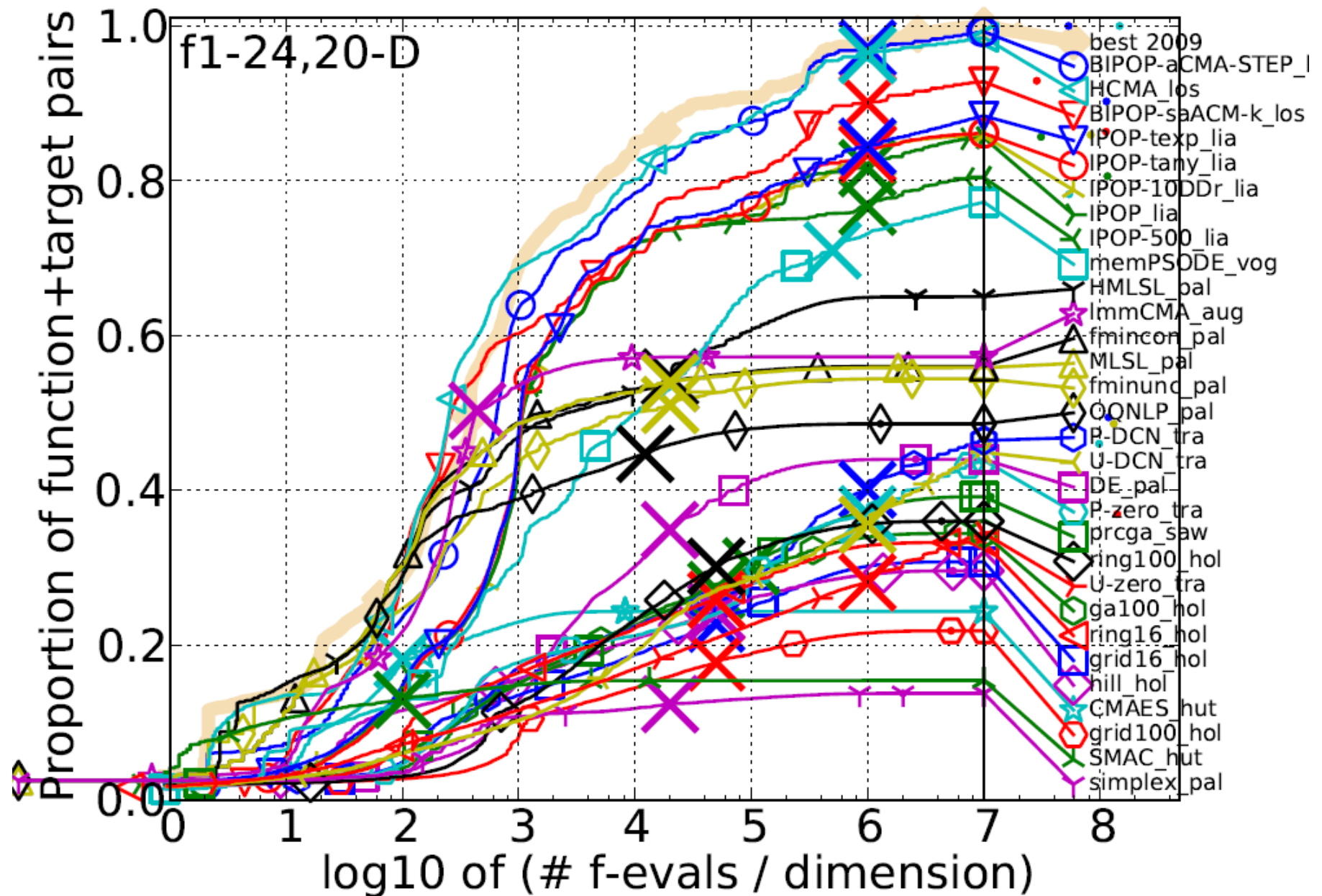  - for unsuccessful runs: simulated restart within 15 instances

# Evaluation of Search Algorithms

a performance should be

- **quantitative** on the ratio scale (highest possible)

  + "algorithm A is two *times* better than algorithm B"
  is a meaningful statement
  + can assume a wide range of values

- **meaningful (interpretable)** with regard to the real world

  possible to transfer from benchmarking to real world

  runtime is the prime candidate (we don't have many choices anyway)

# other plots for single functions

# Scaling Behaviour with Dimension



13 Sharp ridge

# Scaling Behaviour with Dimension



- slanted grid lines: quadratic scaling

- horizontal lines: linear scaling

- light brown: artificial best 2009

# Example: Scaling Behaviour



- slanted grid lines: quadratic scaling

- horizontal lines: linear scaling

- light brown: artificial best 2009

$\Longrightarrow$ Experiments in >40-D are more often than not virtually superfluous

# ERT scatter plot, all dimensions & targets



- estimated Expected Run Time (ERT), two algorithms

- 2-10 D: first algorithm "dominates"

- 20 & 40 D: second algorithm "dominates"

dimension: 2:+, 3:▽, 5:⋆, 10:∘, 20:□, 40:◇

$[\log_{10}(\text{function evaluations})]$

12 Bent cigar

# Questions?

- "two objectives":
  - fast
  - successful
- overfitting?

- "two objectives":
  - fast
  - successful
- overfitting?

f1-24,20-D

- best 2010
- MOS
- IPOP-ACTCMA-ES
- IPOP-CMA-ES
- DE-F-AUC
- DEuniform
- PM-AdapSS-DE
- CMAEGS
- 1plus2mirser
- 1komma2mir
- 1plus1
- 1komma4mirser
- 1komma2mirser
- 1komma4mir
- AVGNEWUOA
- 1komma4
- 1komma4ser
- NBC-CMA
- 1komma2
- 1komma2ser
- ABC
- RCGA
- SPSA

f1-24,20-D

best 2009
BIPOP-CMA-ES
AMALGAM
iAMALGAM
IPOP-SEP-CMA-E
VNS
MA-LS-CHAIN
DASA
G3PCX
NEWUOA
CMA-ESPLUSSEL
Cauchy-EDA
ONEFIFTH
BFGS
PSO_Bounds
GLOBAL
ALPS
FULLNEWUOA
NELDER
NELDERDOERR
EDA-PSO
POEMS
PSO
MCS
Rosenbrock
LSstep
LSfminbnd
GA
DE-PSO
DIRECT
BAYEDA
RANDOMSEARCI
SNOBFIT

Proportion of functions

log10 of (# f-evals / dimension)

# All data 2012 (noisy)



f101-130,20-D

Legend (right side, top to bottom):
best 2012, IPOPsaACM, xNESas, xNES, SNES, ACOR, BIPOPaCMA, BIPOPsaACM, aCMA, CMAES, aCMAa, aCMAm, aCMAma, aCMAmah, aCMAmh, DBRCGA, DE, DEAE, DEb, DEctpb, JADE, JADEb, JADEctpb, NBIPOPaCMA, NIPOPaCMA, DE-AUTO, DE-BFGS, DE-ROLL, DE-SIMPLEX, MVDE, PSO-BFGS

X-axis: log10 of (# f-evals / dimension)
Y-axis: Proportion of functions

# All data 2010 (noisy)



f101-130,20-D

Proportion of functions

log10 of (# f-evals / dimension)

- best 2010
- IPOP-ACTCMA-ES
- IPOP-CMA-ES
- MOS
- CMAEGS
- RCGA
- 1komma4mirser
- 1komma2mirser
- 1komma4mir
- 1komma4
- 1komma2mir
- 1komma4ser
- 1komma2
- 1komma2ser
- AVGNEWUOA
- NEWUOA
- SPSA
- 1plus1
- 1plus2mirser
- ABC
- DE-F-AUC
- DEuniform
- NBC-CMA
- PM-AdapSS-DE
- oPOEMS
- pPOEMS

# All data 2009 (noisy)



f101-130,20-D

best 2009
BIPOP-CMA-ES
AMALGAM
iAMALGAM
IPOP-SEP-CMA-E
VNS
MA-LS-CHAIN
ALPS
BAYEDA
FULLNEWUOA
ONEFIFTH
DASA
GLOBAL
CMA-ESPLUSSEL
EDA-PSO
PSO
NEWUOA
PSO_Bounds
DE-PSO
MCS
SNOBFIT
BFGS
RANDOMSEARCI
Cauchy-EDA
DIRECT
G3PCX
GA
LSfminbnd
LSstep
NELDER
NELDERDOERR
POEMS
Rosenbrock

Proportion of functions

log10 of (# f-evals / dimension)