

# Mirrored Sampling and Sequential Selection for Evolution Strategies

Dimo Brockhoff<sup>1</sup>, Anne Auger<sup>1</sup>, Nikolaus Hansen<sup>1</sup>, Dirk V. Arnold<sup>2</sup>, and Tim Hohm<sup>3</sup>

## Motivation

### BBOB'2009

- (1+1)-CMA surprisingly good on some functions
- even outperforms BIPOP-CMA-ES on Gallagher functions by factor of about 10 (in 20D)
- but elitism bad for noise

### Derandomization

- can improve convergence of ES
- but might introduce certain bias towards smaller step sizes

**Question:** How to design fast "local" non-elitist ES? with derandomization without bias

## Conclusions

### Proposed:

Mirroring and sequential selection are two independent ways to improve local non-elitist ESs

### Results:

- Improved convergence on sphere function
- Implementation within CMA-ES results in better performance on some BBOB'10 functions
- (1,4<sup>s</sup><sub>m</sub>)-ES even faster than (1+1)-ES on sphere w.r.t. the theoretical convergence rates: 0.202 (1+1)-ES ≤ 0.223 (1,4<sup>s</sup><sub>m</sub>)-ES ≤ 0.235 (1+1<sup>s</sup><sub>m</sub>)-ES

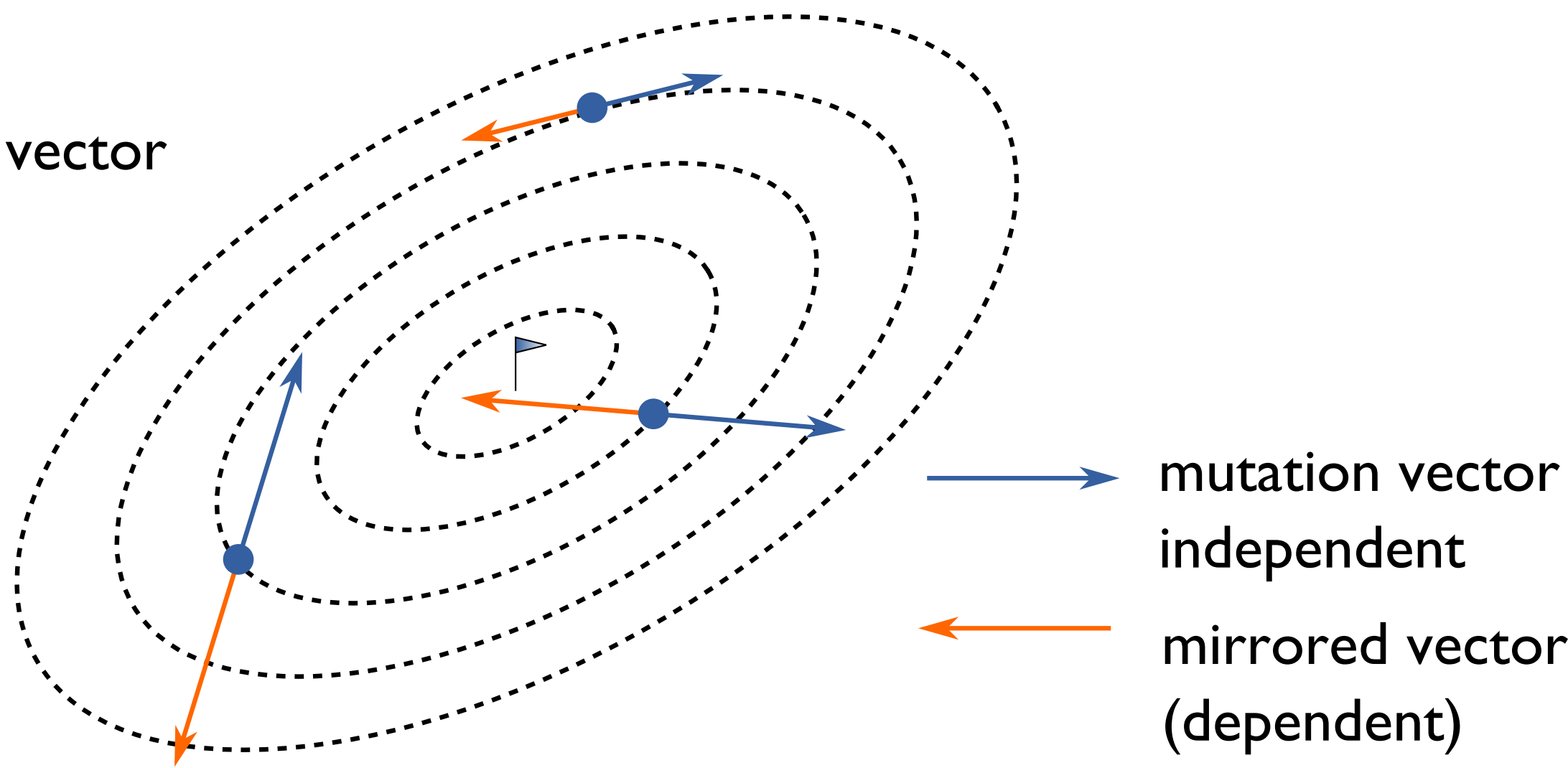
## Idea 1: Mirrored Mutation

### Idea

use one random vector to generate two offspring

### Reasoning

often "good" and "bad" in opposite directions



→ mutation vector independent  
← mirrored vector (dependent)

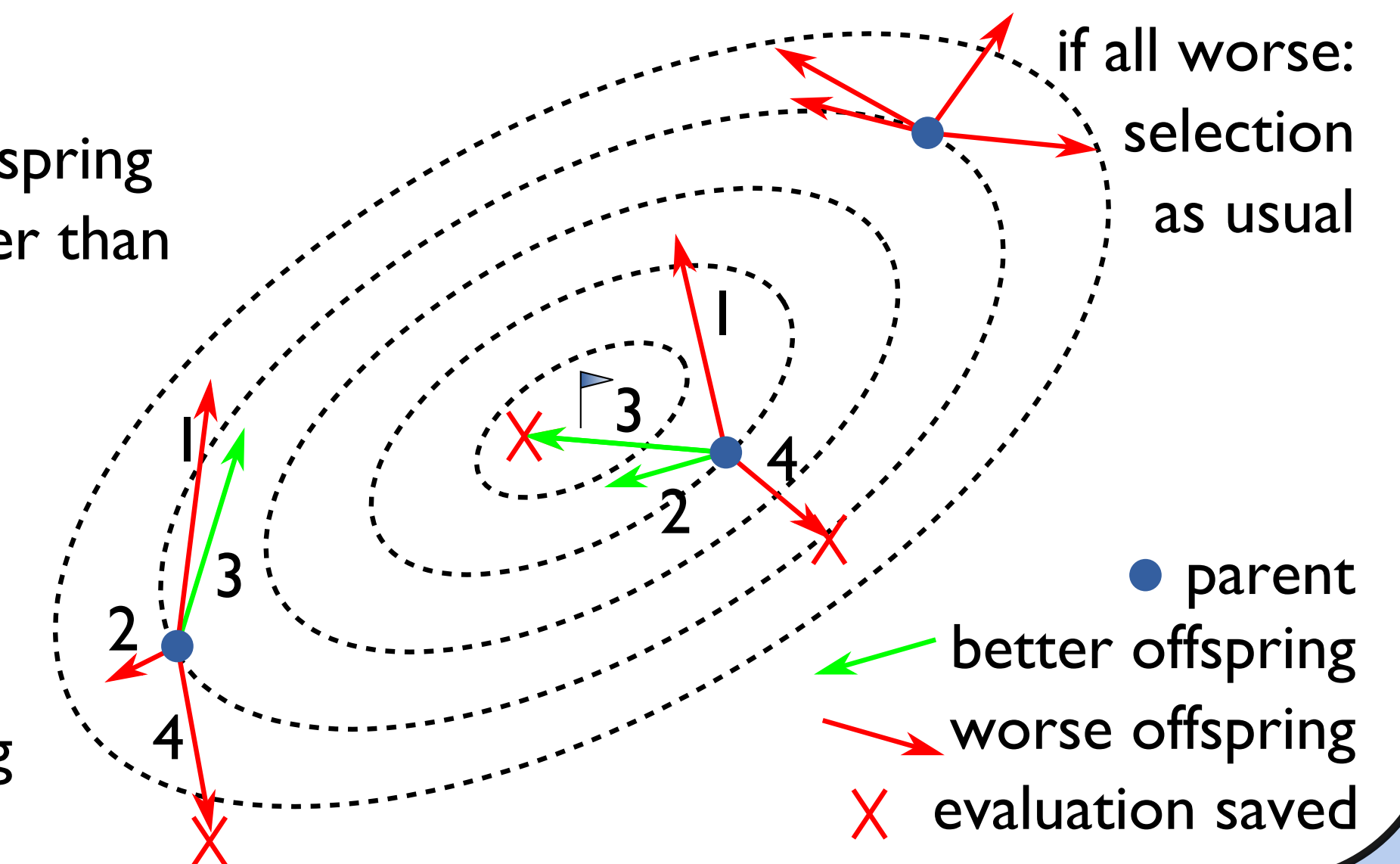
## Idea 2: Sequential Selection

### Idea

stop generation of new offspring as soon as a solution, better than the parent, is found

### Reasoning

if sublevel sets convex one better is enough in particular with mirroring



if all worse: selection as usual

● parent  
← better offspring  
→ worse offspring  
X evaluation saved

## Resulting Variants of the (1,λ)-ES

given:  $X_k \in \mathbb{R}^d, \sigma_k \in \mathbb{R}_{>0}, j \in \mathbb{N}, \lambda \in \mathbb{N}^+, f: \mathbb{R}^d \rightarrow \mathbb{R}$   
 $i \leftarrow 0$

while  $i < \lambda$  do

$i \leftarrow i + 1, j \leftarrow j + 1$

if mirrored sampling and  $j \equiv 0 \pmod{2}$  then

$X_k^i = X_k - \sigma_k \mathcal{N}_k^{i-1}$  use previous sample

else

$X_k^i = X_k + \sigma_k \mathcal{N}_k^i$

if sequential selection and  $f(X_k^i) < f(X_k)$  then

$j \leftarrow 0$  start with a new sample in the next iteration

break;

end while

return  $X_{k+1} = \operatorname{argmin}\{f(X_k^1), \dots, f(X_k^i)\}$

mirroring and sequentialism are independent

standard: (1,λ)-ES

mirroring only: (1,λ<sub>m</sub>)-ES

sequentialism only: (1,λ<sup>s</sup>)-ES

mirroring and sequentialism: (1,λ<sup>s</sup><sub>m</sub>)-ES

## Some Algorithms' Properties

If  $f$  has convex sublevel sets:

never both mirrored offspring are better than parent

⇒ accepted parent after one iteration of (1,2<sub>m</sub>)-ES and (1,2<sup>s</sup><sub>m</sub>)-ES are the same (assuming the same random numbers are used)

⇒ (1,2<sup>s</sup><sub>m</sub>)-ES converges faster than (1,2<sub>m</sub>)-ES on those functions (and diverges faster as well)

(1,∞<sup>s</sup>)-ES = (1+1)-ES if  $\sigma_k$  scale-invariant or constant

[proofs submitted to FOGA]

## Theoretical Convergence Rates of (1,λ)-ES Variants

### Idea

explaining theoretically what can be gained by mirrored mutation and sequential selection in terms of convergence rate

### Scenario

- scale-invariant (1,λ)-ES variants:  $\sigma_k = \sigma \|X_k\|$ : gives optimal convergence rate among all step size adaptive (1,λ)-ES for  $\sigma$  well chosen
- sphere function:  $f(x) = g(\|x\|)$

### Linear Convergence

$\frac{1}{T_k} \ln \frac{\|X_k\|}{\|X_0\|} \rightarrow c$  a.s.  
number of evaluations until iteration  $k$

### Results

**Theorem 4.** For a (1,2<sub>m</sub>)-ES with scale-invariant step-size ( $\sigma_k = \sigma \|X_k\| > 0$ ) on the sphere function  $g(\|x\|)$ , for  $g \in \mathcal{M}$ , linear convergence holds and

$$\frac{1}{T_k} \ln \frac{\|X_k\|}{\|X_0\|} \xrightarrow{k \rightarrow \infty} \frac{1}{2} \frac{1}{2 - p_s(\sigma)} \times E[\ln(1 - 2\sigma|\mathcal{N}|_1 + \sigma^2\|\mathcal{N}\|^2)] \text{ a.s.}$$

where  $T_k$  is the random variable for the number of function evaluations until iteration  $k$ ,  $\mathcal{N}$  is a random vector following a multivariate normal distribution, and  $p_s(\sigma) = \Pr(2|\mathcal{N}|_1 + \sigma\|\mathcal{N}\|^2 < 0)$  is the probability that the first offspring is successful.

### Proof Ideas

$$\frac{1}{T_k} \ln \frac{\|X_k\|}{\|X_0\|} = \frac{1}{T_k} \sum_{i=0}^{k-1} \ln \frac{\|X_{i+1}\|}{\|X_i\|}$$

independent and identically distributed:

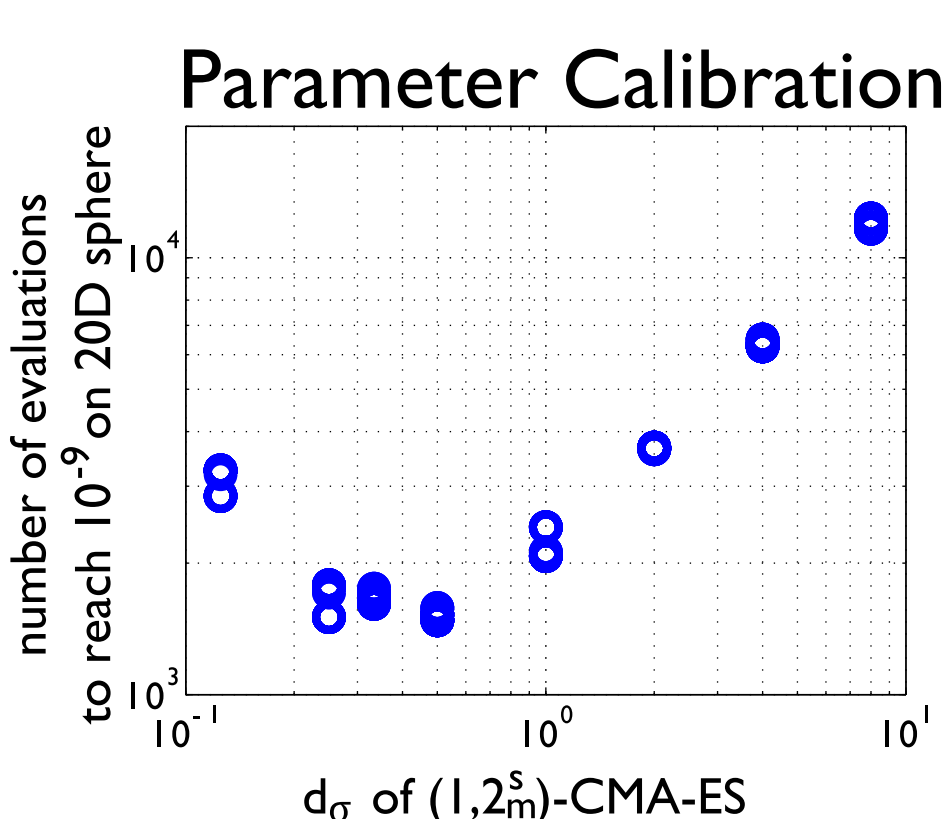
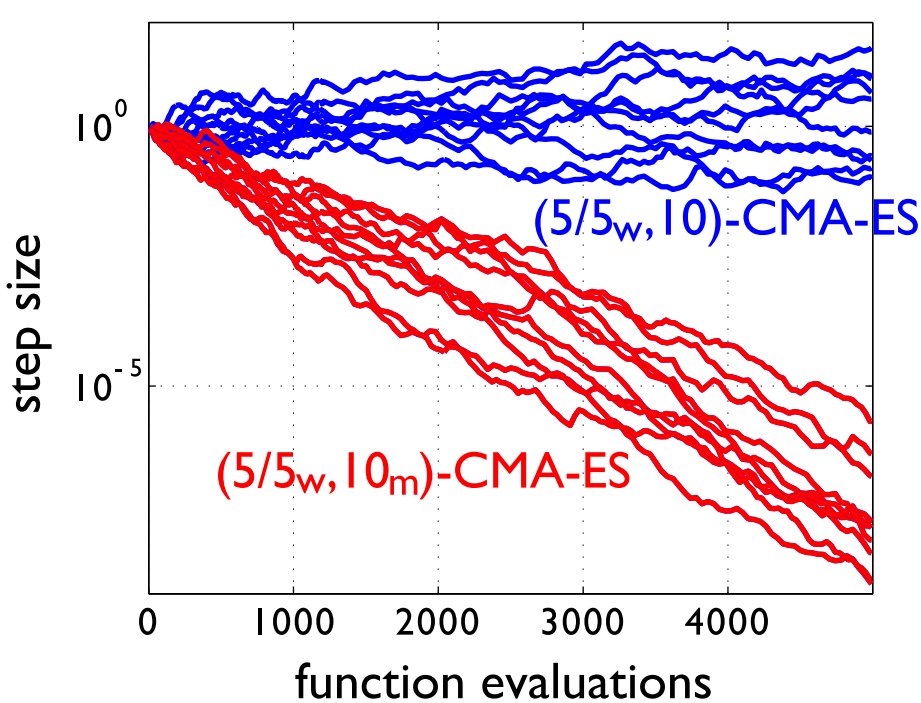
isotropy of sphere function and normal distribution  
+ scale-invariant step size rule

LLN shows that RHS converges to  $E(\ln(\|X_{i+1}\|/\|X_i\|))$

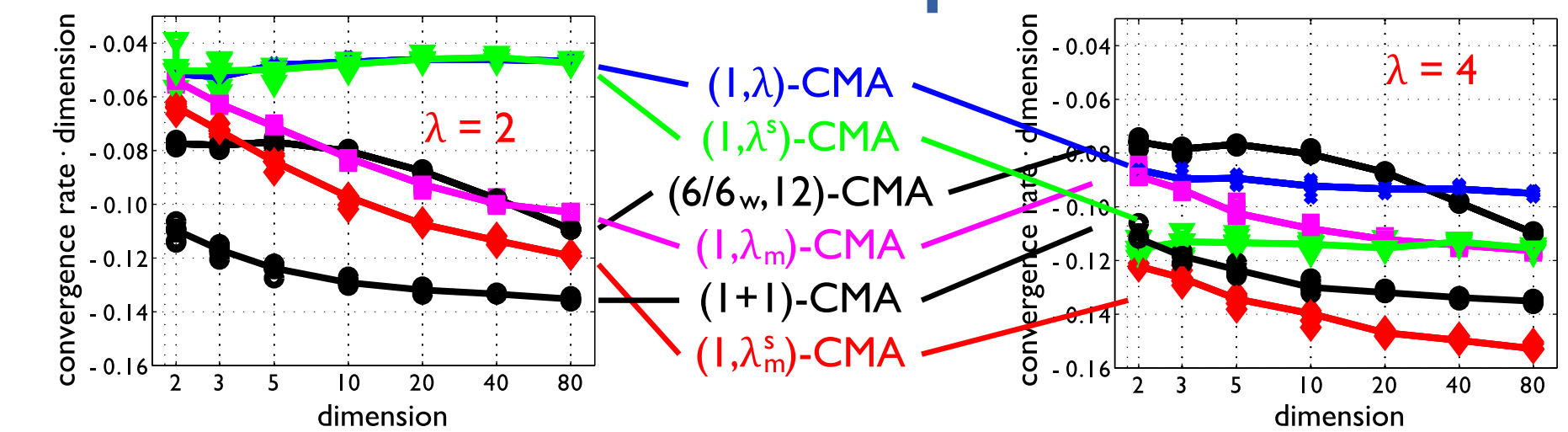
## Application to the CMA-ES

implementation straight forward, but with bias if  $\mu > 1$

⇒ only  $\mu=1$  here

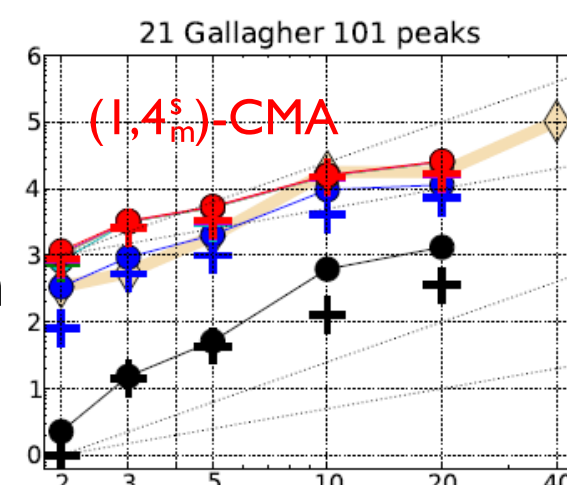


### results on sphere



### results on BBOB'2010

- (1,4<sup>s</sup><sub>m</sub>)-CMA-ES turned out to be fastest local non-elitist strategy tested
- 3rd best of BBOB'2009/10 on Gallagher with 101 peaks (3x faster than (1+1)-CMA-ES)
- even more competitive on noisy functions



## Numerical Results

### Idea

Estimate convergence rates by exploiting theoretical results via Monte Carlo sampling

### Procedure

- 1) Exploitation of theoretical expression of convergence rate: perform Monte Carlo estimate of probability of success and expected value for different  $\sigma$  ( $0.01 \leq \sigma \leq 3$ ); averaged over  $10^6$  samples
- 2) extract smallest value (corresponds to optimal conv. rate though bias towards smaller values)

### Simulated Convergence Rates

