

UQ & Data Analysis in Applied Sciences

# CS4: Optimization using Covariance Matrix Adaptation

Lecturer: Georgios Arampatzis

based on “**The CMA Evolution Strategy: A Tutorial**”  
by Nikolaus Hansen

# GOAL

- ◆ Given and objective function in continuous domain

$$f : \mathcal{X} \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$$

- ◆ find

$$x^* = \operatorname{argmax}_{x \in \mathcal{X}} f(x)$$

# BLACK BOX SCENARIO

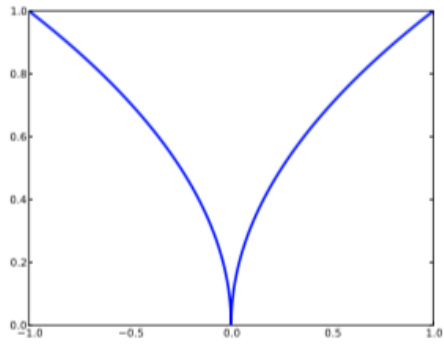


- gradients are not available
- non-convex
- non-smooth
- multimodal
- high dimensional
- noisy
- ...

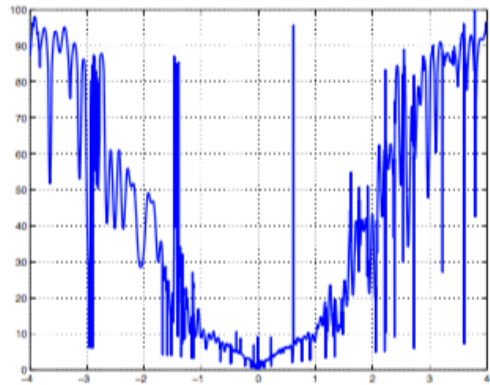


# DIFFICULT FUNCTIONS

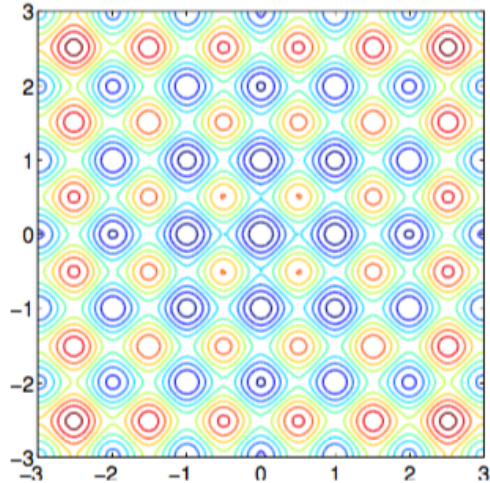
<https://www.lri.fr/~hansen/gecco2013-CMA-ES-tutorial.pdf>



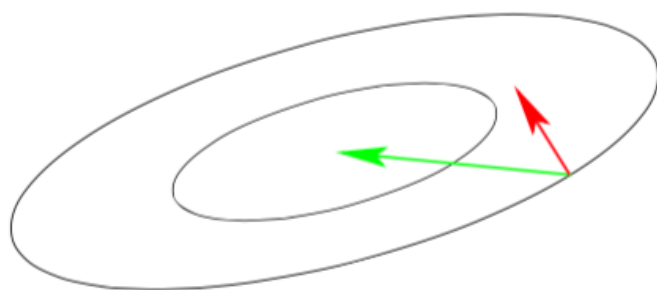
\* non-linear, non-quadratic, non-convex



\* non-smooth, multimodal, noisy



\* non-separability



\* ill conditioning

# EXAMPLES

- ◆ shape optimization
  - ◆ curve fitting
  - ◆ airfoils
- ◆ model calibration
  - ◆ biological
  - ◆ physical
- ◆ parameter calibration
  - ◆ controller

# RANDOMIZED BLACK BOX SEARCH

➔ initialize

➔ population size  $\lambda \in \mathbb{N}$

➔ distribution parameters  $\vartheta^{(0)}$

➔ until happy

➔ sample  $\mathbf{x}_i \sim P(\mathbf{x}|\vartheta^{(k)}), \quad i = 1, \dots, \lambda$

➔ evaluate  $f(\mathbf{x}_i), \quad i = 1, \dots, \lambda$

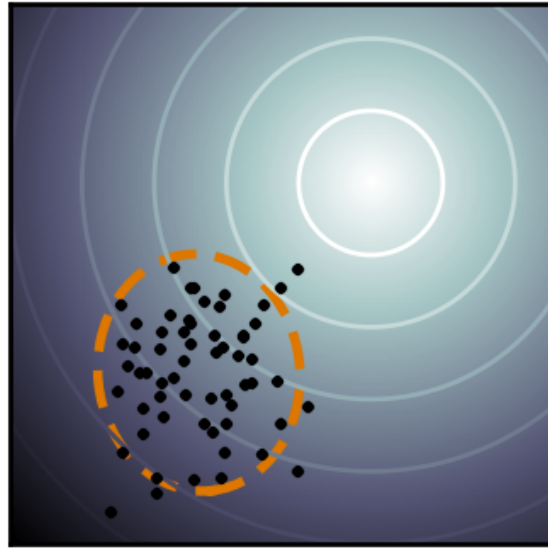
➔ update parameters

$$\vartheta^{(k+1)} = F(\vartheta, \mathbf{x}_1, \dots, \mathbf{x}_\lambda, f(\mathbf{x}_1), \dots, f(\mathbf{x}_\lambda))$$

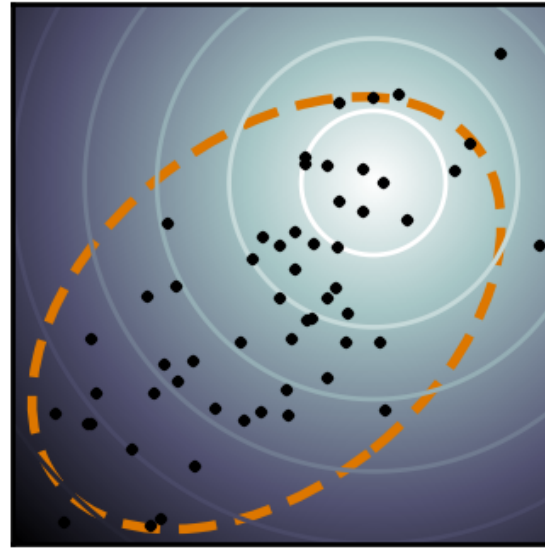
# RANDOMIZED BLACK BOX SEARCH

[https://upload.wikimedia.org/wikipedia/commons/d/d8/Concept\\_of\\_directional\\_optimization\\_in\\_CMA-ES\\_algorithm.png](https://upload.wikimedia.org/wikipedia/commons/d/d8/Concept_of_directional_optimization_in_CMA-ES_algorithm.png)

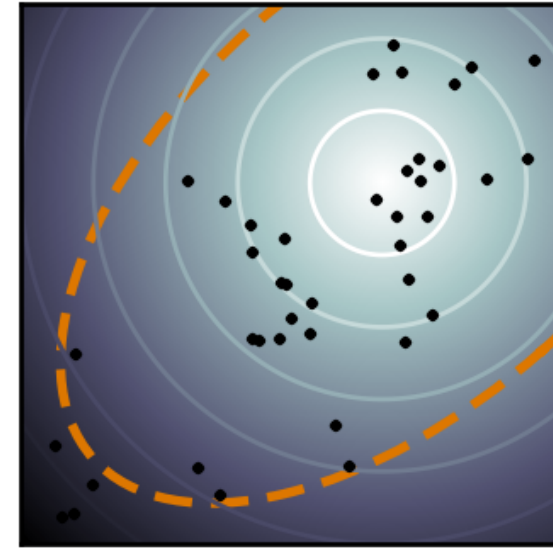
Generation 1



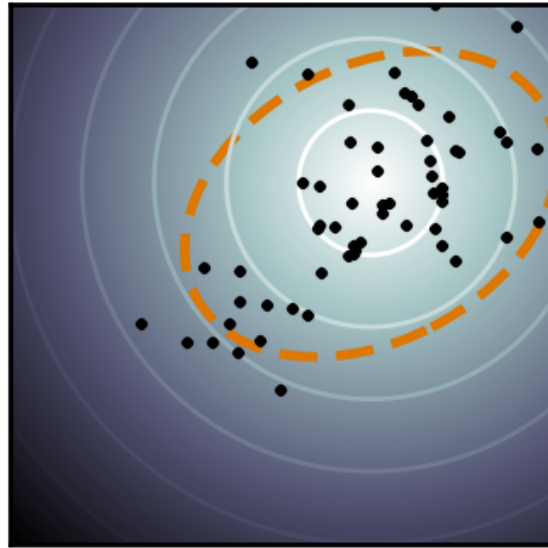
Generation 2



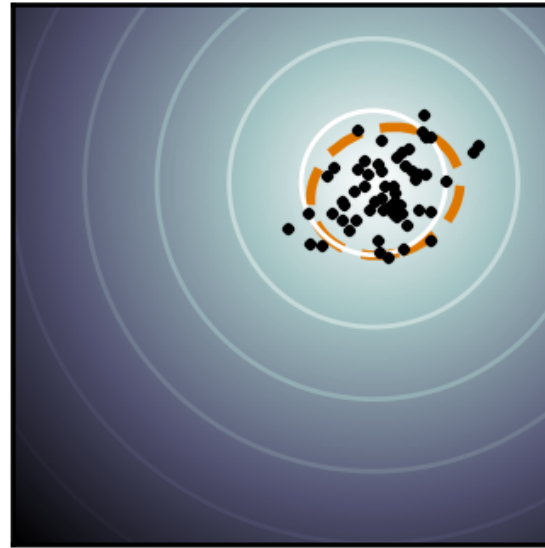
Generation 3



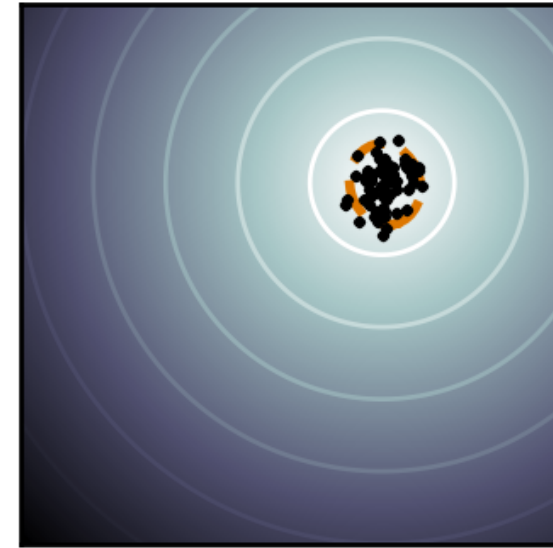
Generation 4



Generation 5



Generation 6



# THE SAMPLING DISTRIBUTION

❖ choose sampling distribution

\* isotropic  
\* maximum entropy

$$\mathbf{x}_i \sim P(\mathbf{x}|\vartheta^{(k)}) = \mathbf{m}^{(k)} + \sigma^{(k)} \mathcal{N}(\mathbf{0}, \mathbf{C}^{(k)})$$

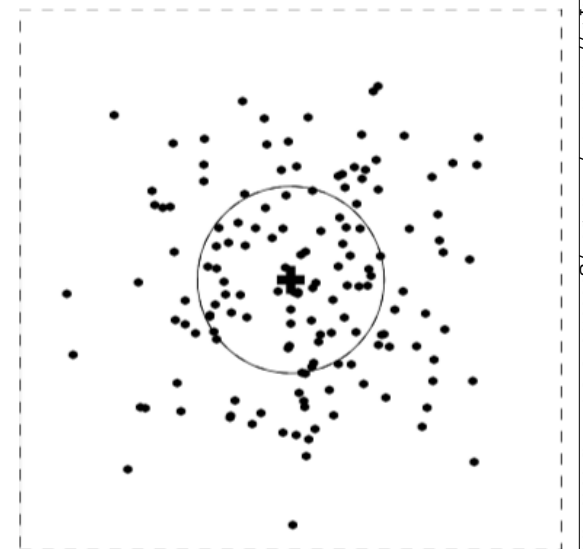
❖ choose how to update the parameters

$\mathbf{m}^{(k)}$  —  $\sigma^{(k)}$  —  $\mathbf{C}^{(k)}$

represents the favorite solution

controls the step-size

determines the shape of the distribution



# THE INGREDIENTS

$m^{(k)}$

Evolution Strategy

$C^{(k)}$

Covariance Matrix Adaptation

$\sigma^{(k)}$

Step Size Control

# EVOLUTION STRATEGIES

**A. SELECT AND RECOMBINE**

**B. COMPUTE THE MEAN**

# EVOLUTION STRATEGIES

\* # parents:  $\mu$

\* # children:  $\lambda$

---

\* elitist selection:  $(\mu + \lambda)$ -ES

\* non-elitist selection:  $(\mu, \lambda)$ -ES

---

$(1 + 1)$ -ES

\* sample one child from parent  $\mathbf{m}$

$$x \sim \mathbf{m} + \sigma \mathcal{N}(0, \mathbf{C})$$

\* if  $x$  is better than  $\mathbf{m}$  select

$$\mathbf{m} \leftarrow x$$



# EVOLUTION STRATEGIES

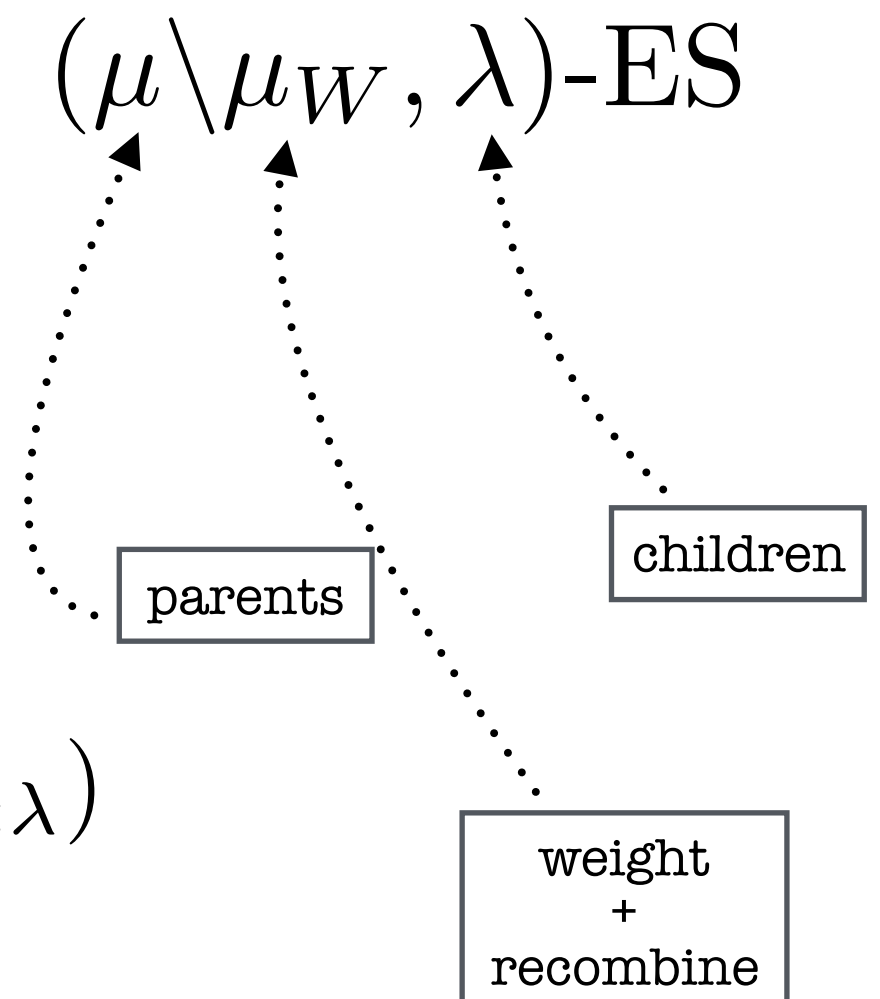
## Selection and weighted Recombination

$$\mathbf{m}^{(k+1)} = \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda}^{(k+1)}$$

$$w_i \propto \mu - i + 1$$

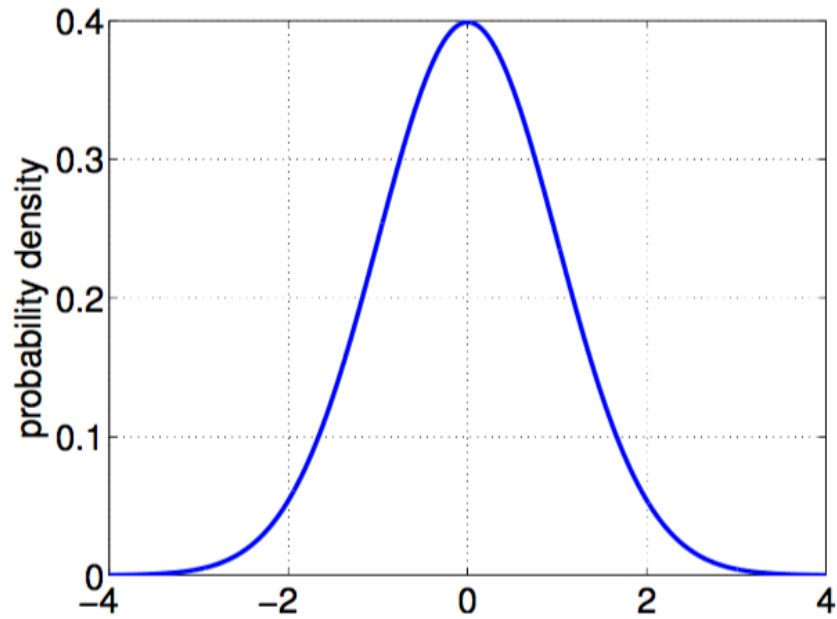
- $w_1 \geq w_2 \geq \dots \geq w_{\mu} > 0$
- $\sum_i w_i = 1$
- $f(\mathbf{x}_{1:\lambda}) \leq f(\mathbf{x}_{2:\lambda}) \leq \dots \leq f(\mathbf{x}_{\mu:\lambda})$

$(\mu \setminus \mu_w, \lambda)$ -ES

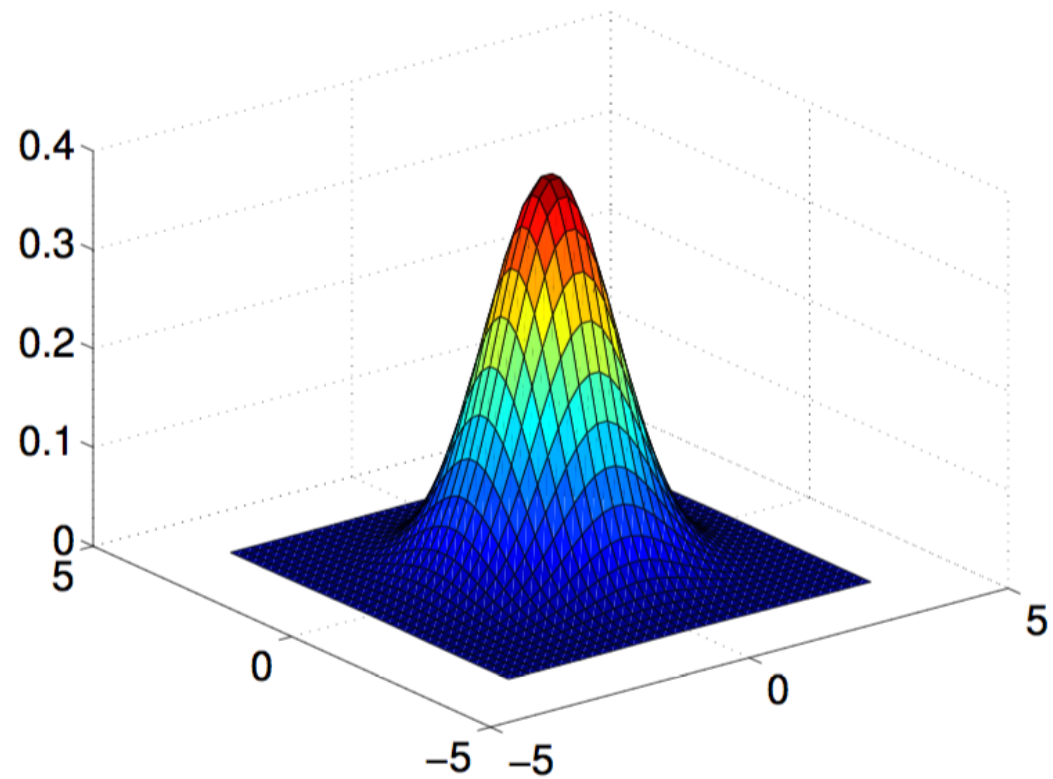


# THE NORMAL DISTRIBUTION

Standard Normal Distribution



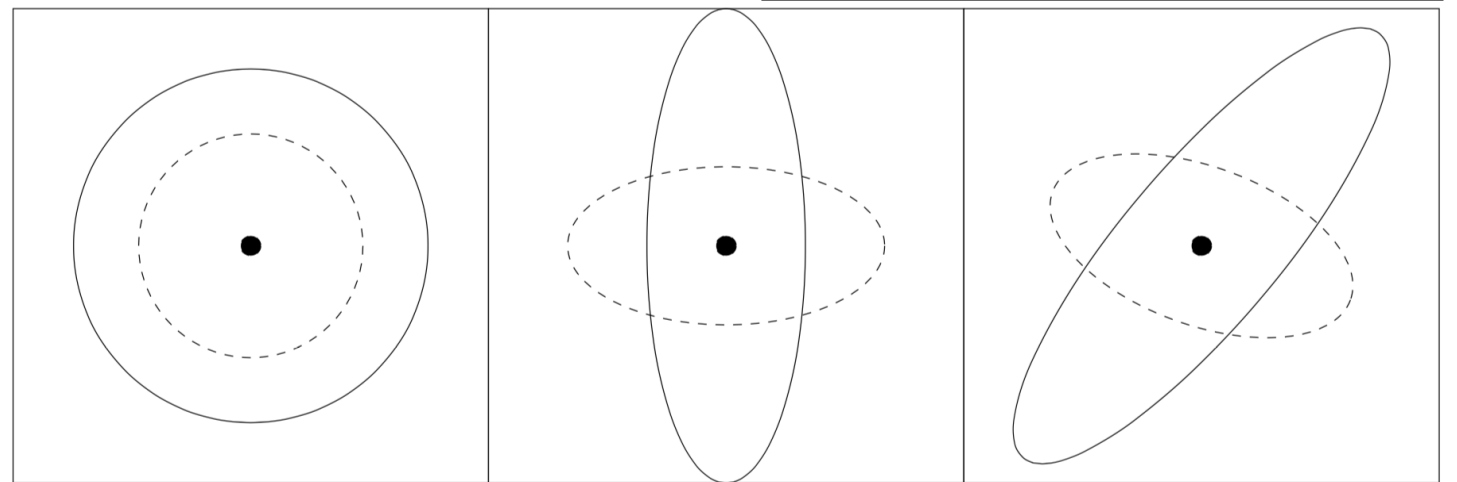
2-D Normal Distribution



⊠ geometrical interpretation of covariance matrix

$$(x - \mathbf{m})^\top \mathbf{C}^{-1} (x - \mathbf{m}) = \text{const.}$$

<https://www.lri.fr/~hansen/gecco2013-CMA-ES-tutorial.pdf>



$$\mathcal{N}(\mathbf{m}, \sigma^2 \mathbf{I}) \sim \mathbf{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathcal{N}(\mathbf{m}, \mathbf{D}^2) \sim \mathbf{m} + \mathbf{D} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathcal{N}(\mathbf{m}, \mathbf{C}^2) \sim \mathbf{m} + \mathbf{C}^{-\frac{1}{2}} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

# COVARIANCE MATRIX ADAPTATION

**A. ESTIMATION FROM SCRATCH**

**B. RANK-MU-UPDATE**

**C. RANK-ONE-UPDATE**

# COVARIANCE MATRIX ADAPTATION

## A. ESTIMATION FROM SCRATCH

$$\mathbf{x}_i \sim \mathbf{m}^{(k)} + \sigma^{(k)} \mathcal{N}(\mathbf{0}, \mathbf{C}^{(k)})$$

$$\mathbf{C}_{\text{emp}}^{(k+1)} = \frac{1}{\lambda - 1} \sum_{i=1}^{\lambda} (\mathbf{x}_i^{(k+1)} - \mathbf{m}^{(k+1)}) (\mathbf{x}_i^{(k+1)} - \mathbf{m}^{(k+1)})^\top$$

$$\mathbf{C}_{\lambda}^{(k+1)} = \frac{1}{\lambda} \sum_{i=1}^{\lambda} (\mathbf{x}_i^{(k+1)} - \mathbf{m}^{(k)}) (\mathbf{x}_i^{(k+1)} - \mathbf{m}^{(k)})^\top$$

$$\mathbb{E} \left[ \mathbf{C}_{\text{emp}}^{(k+1)} \mid \mathbf{C}^{(k)} \right] = \mathbf{C}^{(k)}$$

$$\mathbb{E} \left[ \mathbf{C}_{\lambda}^{(k+1)} \mid \mathbf{C}^{(k)} \right] = \mathbf{C}^{(k)}$$

unbiased  
estimators

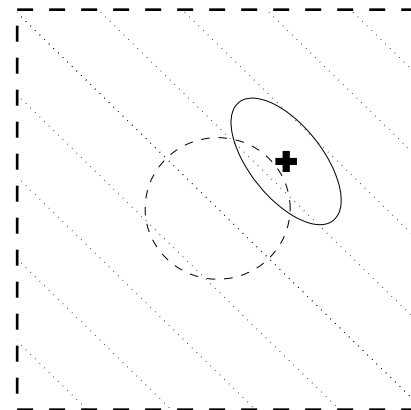
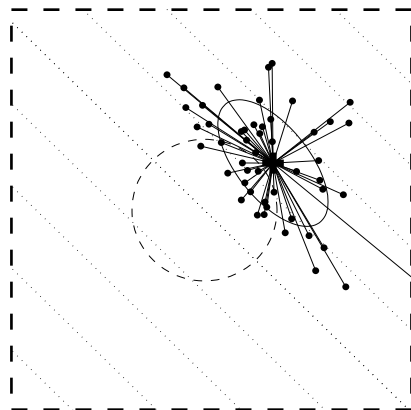
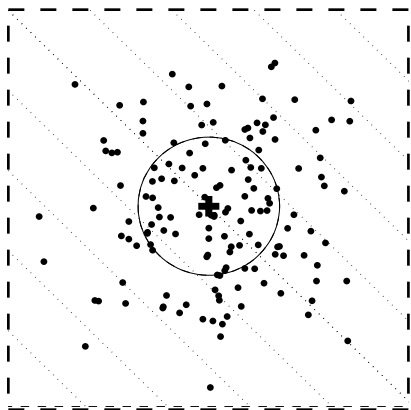
# COVARIANCE MATRIX ADAPTATION

## A. ESTIMATION FROM SCRATCH

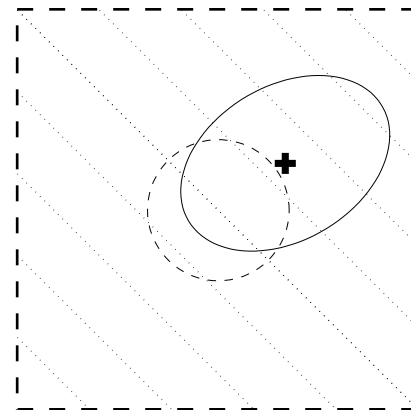
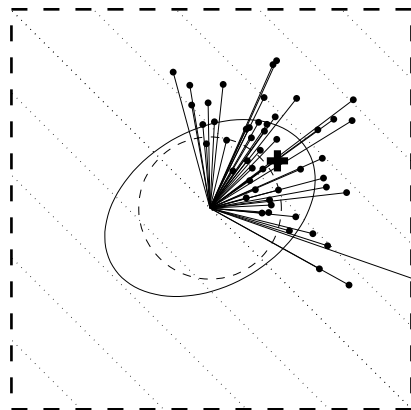
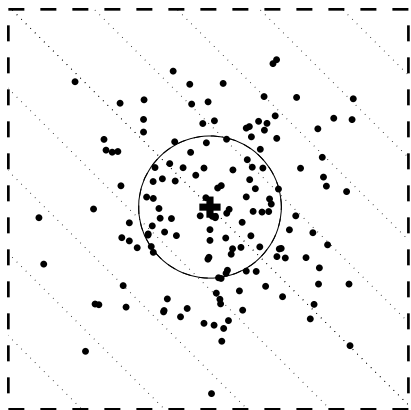
$$C_{EMNA}^{(k+1)} = \frac{1}{\mu} \sum_{i=1}^{\mu} (\mathbf{x}_i^{(k+1)} - \mathbf{m}^{(k+1)}) (\mathbf{x}_i^{(k+1)} - \mathbf{m}^{(k+1)})^\top$$

$$C_{\mu}^{(k+1)} = \frac{1}{\mu} \sum_{i=1}^{\mu} (\mathbf{x}_i^{(k+1)} - \mathbf{m}^{(k)}) (\mathbf{x}_i^{(k+1)} - \mathbf{m}^{(k)})^\top$$

$w_i$



- \* smaller variance
- \* increase geometrically fast
- \* premature convergence



# COVARIANCE MATRIX ADAPTATION

## B. RANK-MU UPDATE

- ✦ estimation from scratch works well for large populations
- ✦ in order to be fast population must be small
- ✦ use information from past

$$C^{k+1} = \frac{1}{k+1} \sum_{i=0}^k \frac{1}{\sigma^{(i)2}} C_{\mu}^{i+1}$$

# COVARIANCE MATRIX ADAPTATION

## B. RANK-MU UPDATE

- ❖ use information from past
- ❖ assign recent generations higher weight
- ❖ exponential smoothing

$$\begin{aligned} \mathbf{C}^{(k+1)} &= (1 - c_\mu) \mathbf{C}^{(k)} + c_\mu \frac{1}{\sigma^{(k)2}} \mathbf{C}_\mu^{(k+1)} \\ &= (1 - c_\mu) \mathbf{C}^{(k)} + c_\mu \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}^{(k+1)} \mathbf{y}_{i:\lambda}^{(k+1)\top} \end{aligned}$$

$$\mathbf{y}_{i:\lambda}^{(k+1)} = \frac{\mathbf{x}_{i:\lambda}^{(k+1)} - \mathbf{m}^{(k)}}{\sigma^{(k)}}$$

rank:  $\min\{\mu, n\}$

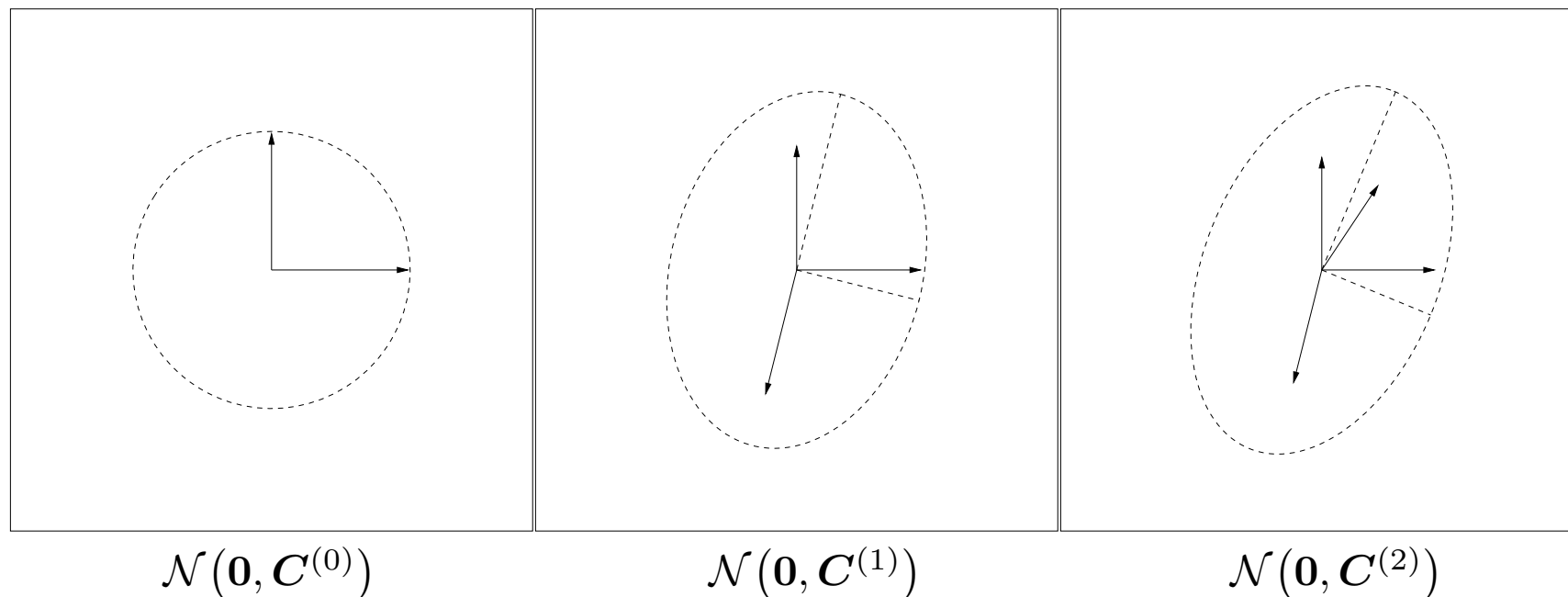
# COVARIANCE MATRIX ADAPTATION

## C. RANK-ONE UPDATE

$$\mathcal{N}(0, 1)\mathbf{y}_1 + \dots + \mathcal{N}(0, 1)\mathbf{y}_k \sim \mathcal{N}(\mathbf{0}, \sum_{i=1}^k \mathbf{y}_i \mathbf{y}_i^\top)$$

♣ the singular distribution  $\mathcal{N}(\mathbf{0}, \mathbf{y}_i \mathbf{y}_i^\top)$  generates the vector  $\mathbf{y}_i$  with maximum likelihood

<https://www.lri.fr/~hansen/gecco2013-CMA-ES-tutorial.pdf>





# COVARIANCE MATRIX ADAPTATION

## C. RANK-ONE UPDATE

♣ assume

$$\mathbf{y}^{(k+1)} = \frac{\mathbf{x}_{1:\lambda}^{(k+1)} - \mathbf{m}^{(k)}}{\sigma^{(k)}}$$

♣ then the covariance matrix

$$\mathbf{C}^{(k+1)} = (1 - c_1)\mathbf{C}^{(k)} + c_1 \mathbf{y}^{(k+1)} \mathbf{y}^{(k+1)\top}$$

♣ increases the probability of generating  $\mathbf{y}^{(k+1)}$  in the next generation

# COVARIANCE MATRIX ADAPTATION

## C. RANK-ONE UPDATE + CUMULATION

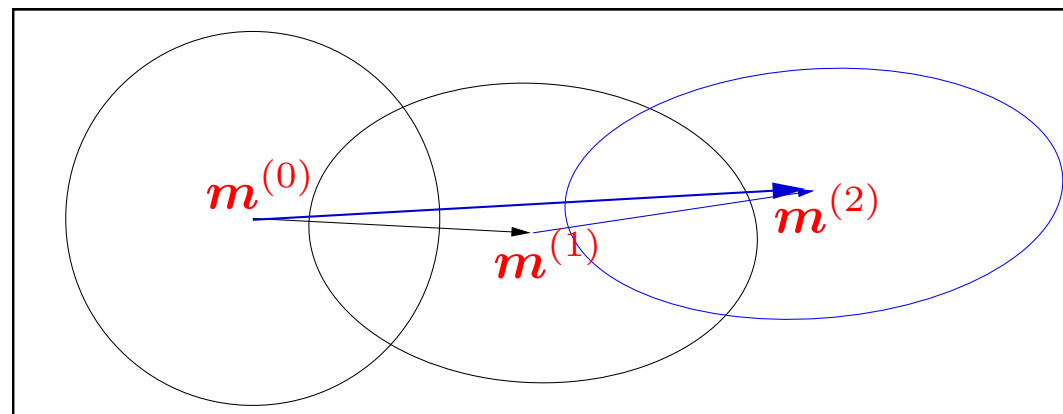
❖ loss of sign information  $\mathbf{y}\mathbf{y}^\top = -\mathbf{y}(-\mathbf{y})^\top$

❖ introduce the evolution path

$$\mathbf{p}_c^{(k+1)} = \sum_{i=1}^k \frac{\mathbf{m}^{(i+1)} - \mathbf{m}^{(i)}}{\sigma^{(i)}}$$

❖ or with exponential smoothing

$$\mathbf{p}_c^{(k+1)} = (1 - c_c)\mathbf{p}_c^{(k)} + c_c \frac{\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)}}{\sigma^{(i)}}$$



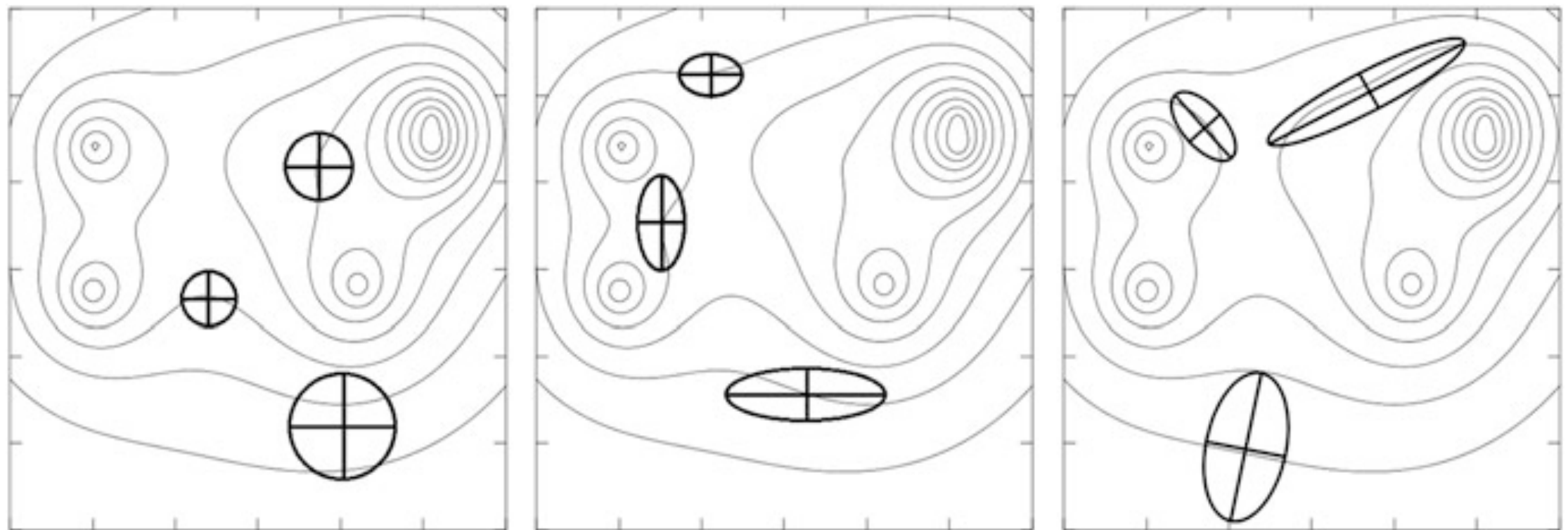
# COVARIANCE MATRIX ADAPTATION

$$\begin{aligned} \mathbf{C}^{(k+1)} &= (1 - c_\mu - c_1) \mathbf{C}^{(k)} \\ &+ c_\mu \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}^{(k+1)} \mathbf{y}_{i:\lambda}^{(k+1)\top} \\ &+ c_1 \mathbf{p}_c^{(k+1)} \mathbf{p}_c^{(k+1)} \end{aligned}$$

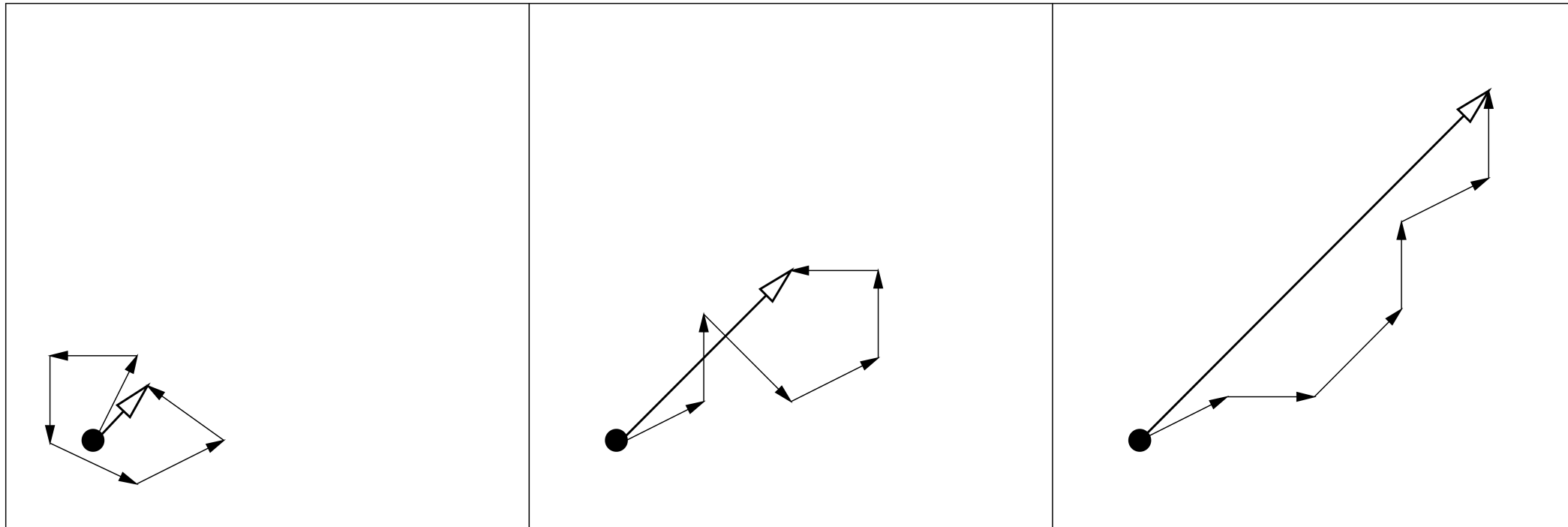
# COVARIANCE MATRIX ADAPTATION

- ❖ learns all pairwise dependencies between variables
- ❖ learns a pre rotated problem representation
- ❖ learns a new Mahalanobis metric
- ❖ approximates the inverse Hessian on quadratic functions

Contemporary Evolution Strategies, Bäck, T.,  
Foussette, C., Krause, P., Springer



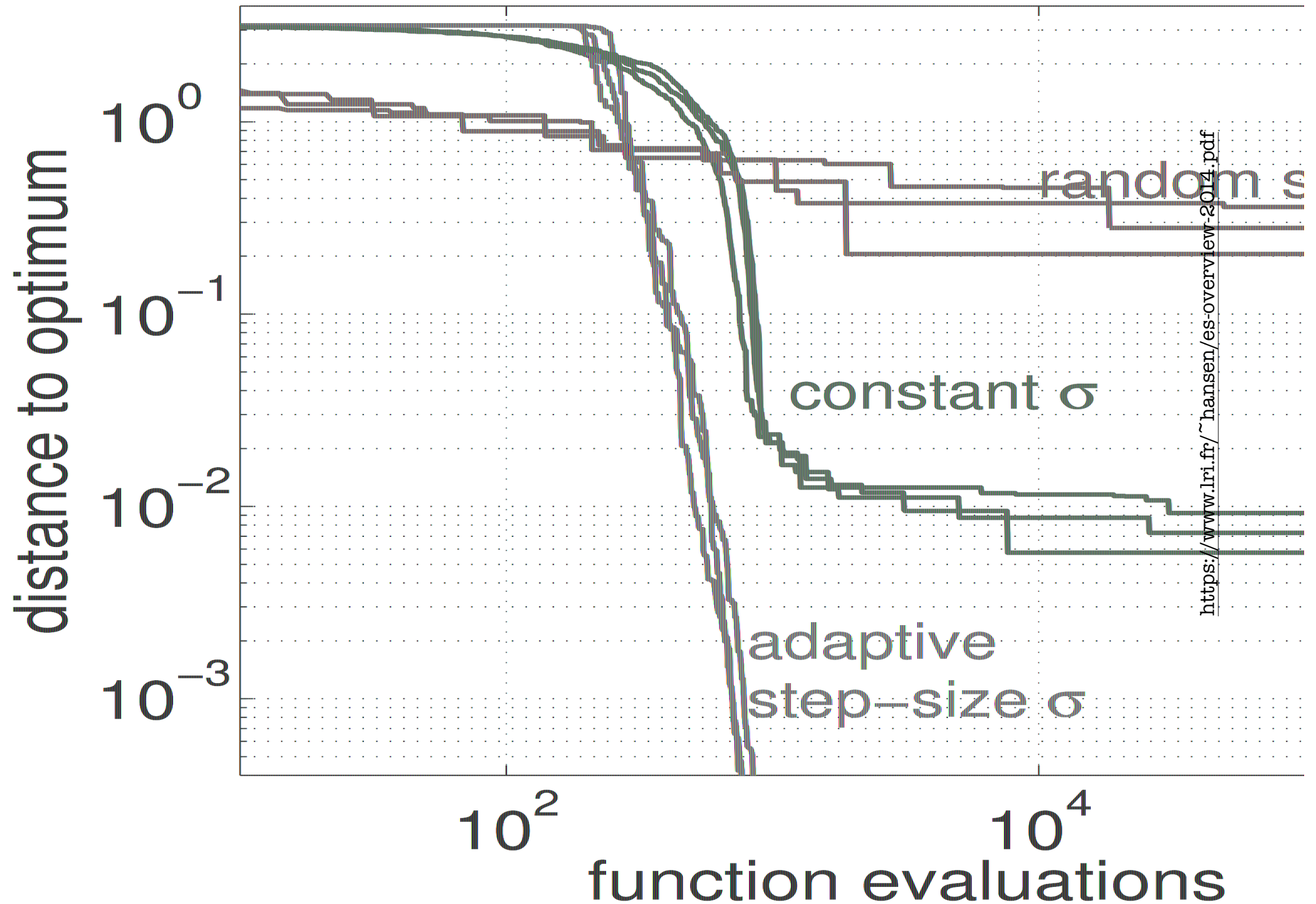
# STEP-SIZE CONTROL



<https://www.lri.fr/~hansen/gecco2013-CMA-ES-tutorial.pdf>

- ❖ evolution path is short: single steps cancel out (anti-correlated)
- ❖ evolution path is long: steps point to the same direction (correlated)
- ❖ evolution path is OK: **uncorrelated** steps

# STEP-SIZE CONTROL



# THE ALGORITHM

## Set parameters

Set parameters  $\lambda$ ,  $\mu$ ,  $w_{i=1\dots\mu}$ ,  $c_\sigma$ ,  $d_\sigma$ ,  $c_c$ ,  $c_1$ , and  $c_\mu$  to their default values according to Table 1.

## Initialization

Set evolution paths  $\mathbf{p}_\sigma = \mathbf{0}$ ,  $\mathbf{p}_c = \mathbf{0}$ , covariance matrix  $\mathbf{C} = \mathbf{I}$ , and  $g = 0$ .

Choose distribution mean  $\mathbf{m} \in \mathbb{R}^n$  and step-size  $\sigma \in \mathbb{R}_+$  problem dependent.<sup>1</sup>

## Until termination criterion met, $g \leftarrow g + 1$

Sample new population of search points, for  $k = 1, \dots, \lambda$

$$\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (35)$$

$$\mathbf{y}_k = \mathbf{B}\mathbf{D}\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad (36)$$

$$\mathbf{x}_k = \mathbf{m} + \sigma\mathbf{y}_k \sim \mathcal{N}(\mathbf{m}, \sigma^2\mathbf{C}) \quad (37)$$

## Selection and recombination

$$\langle \mathbf{y} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \quad \text{where } \sum_{i=1}^{\mu} w_i = 1, w_i > 0 \quad (38)$$

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{y} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda} \quad (39)$$

## Step-size control

$$\mathbf{p}_\sigma \leftarrow (1 - c_\sigma)\mathbf{p}_\sigma + \sqrt{c_\sigma(2 - c_\sigma)\mu_{\text{eff}}} \mathbf{C}^{-\frac{1}{2}} \langle \mathbf{y} \rangle_w \quad (40)$$

$$\sigma \leftarrow \sigma \times \exp\left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1\right)\right) \quad (41)$$

## Covariance matrix adaptation

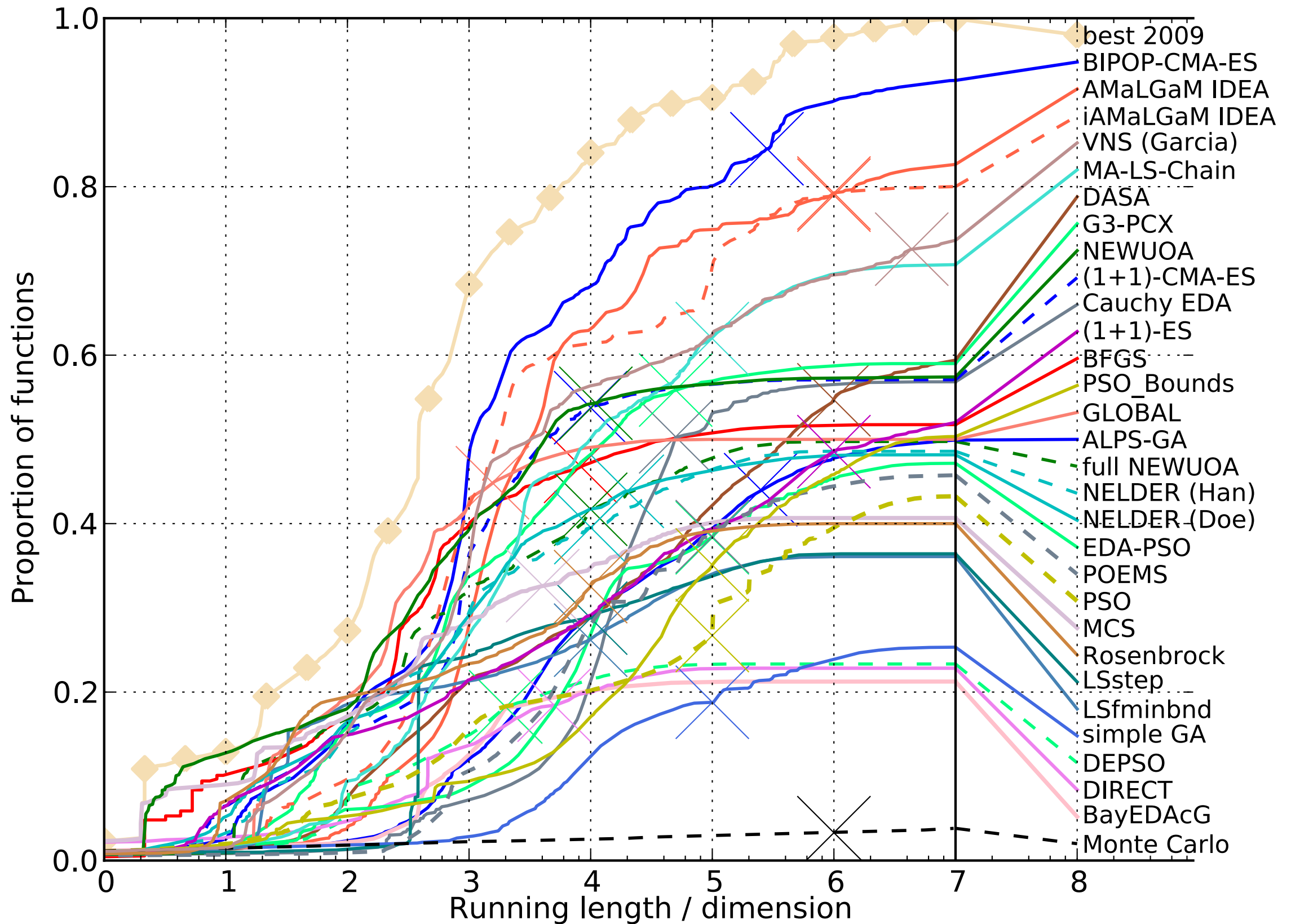
$$\mathbf{p}_c \leftarrow (1 - c_c)\mathbf{p}_c + h_\sigma \sqrt{c_c(2 - c_c)\mu_{\text{eff}}} \langle \mathbf{y} \rangle_w \quad (42)$$

$$\mathbf{C} \leftarrow (1 - c_1 - c_\mu)\mathbf{C} + c_1(\mathbf{p}_c\mathbf{p}_c^T + \delta(h_\sigma)\mathbf{C}) + c_\mu \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}\mathbf{y}_{i:\lambda}^T \quad (43)$$

<sup>1</sup>The optimum should presumably be within the initial cube  $\mathbf{m} \pm 3\sigma(1, \dots, 1)^T$ . If the optimum is expected to be in the initial search interval  $[a, b]^n$  we may choose the initial search point,  $\mathbf{m}$ , uniformly randomly in  $[a, b]^n$ , and  $\sigma = 0.3(b - a)$ . Different search intervals  $\Delta s_i$  for different variables can be reflected by a different initialization of  $\mathbf{C}$ , in that the diagonal elements of  $\mathbf{C}$  obey  $c_{ii} = (\Delta s_i)^2$ . Remark that the  $\Delta s_i$  should not disagree by several orders of magnitude. Otherwise a scaling of the variables should be applied.

# BENCHMARKS

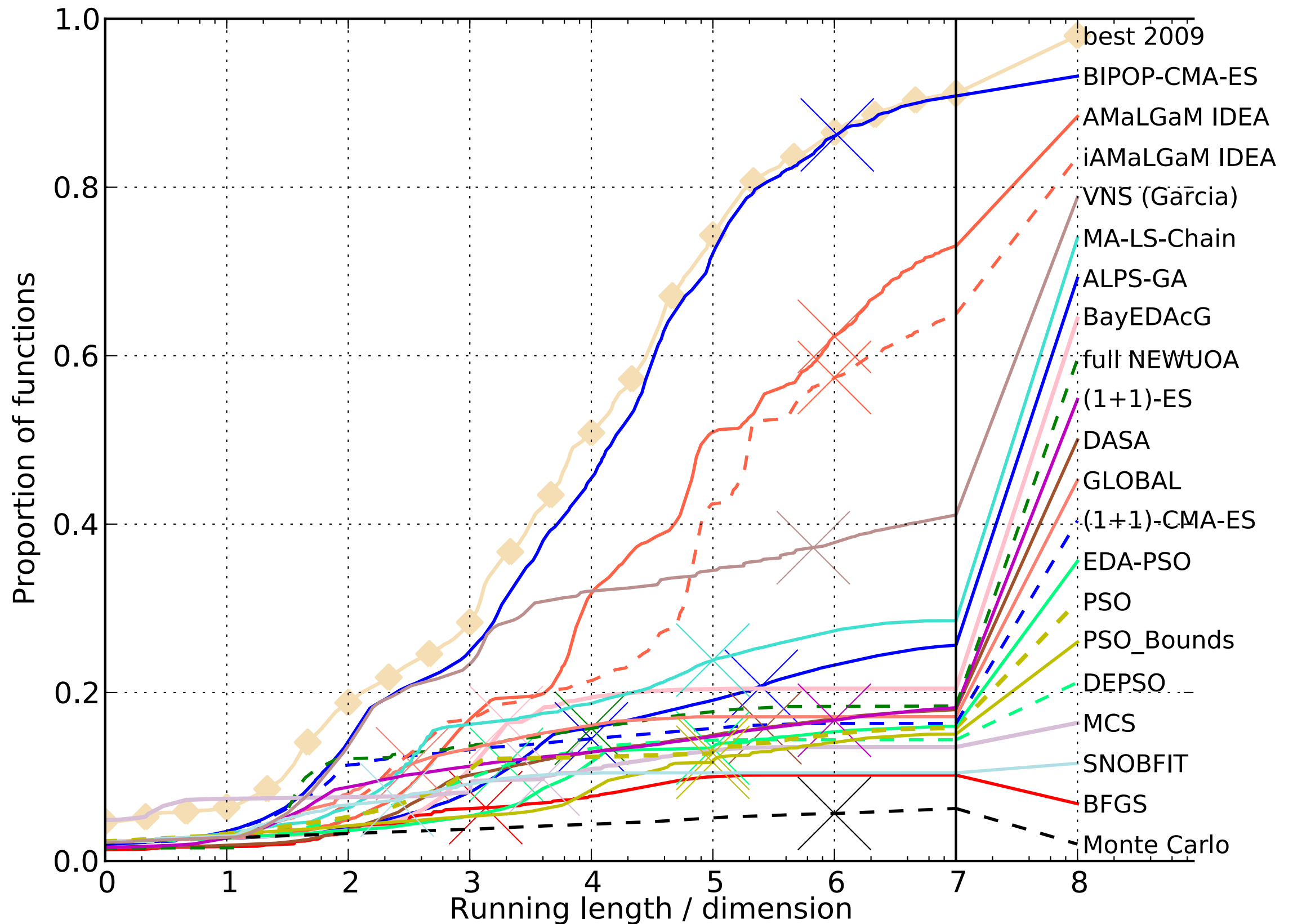
## deterministic functions



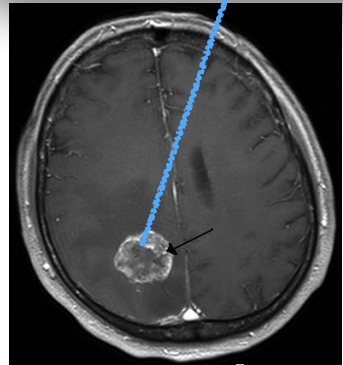
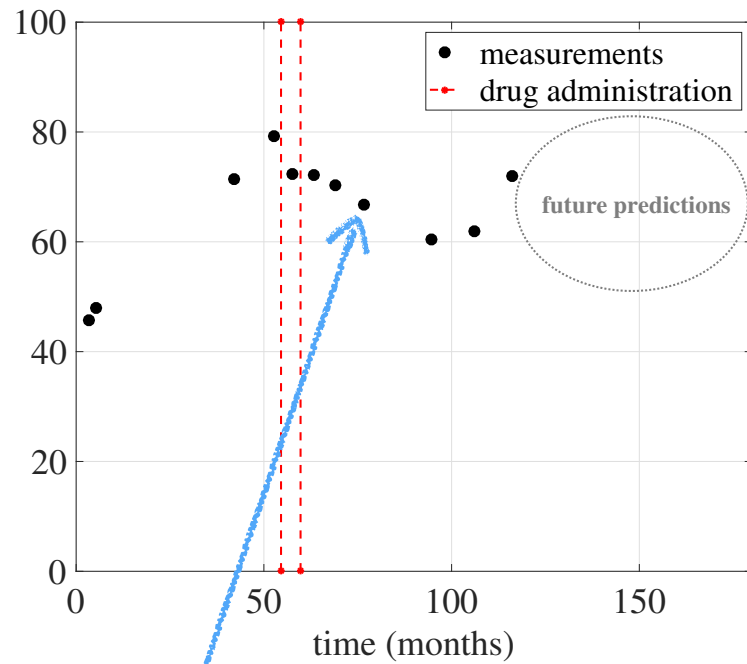


# BENCHMARKS

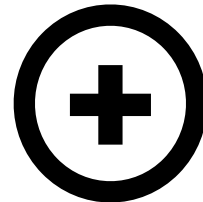
## noisy functions



# APPLICATION



[https://en.wikipedia.org/wiki/Brain\\_tumor#/media/File:Hirnmetastase\\_MRT-T1\\_KM.jpg](https://en.wikipedia.org/wiki/Brain_tumor#/media/File:Hirnmetastase_MRT-T1_KM.jpg)



$$\frac{dC}{dt} = -\vartheta_1 C$$

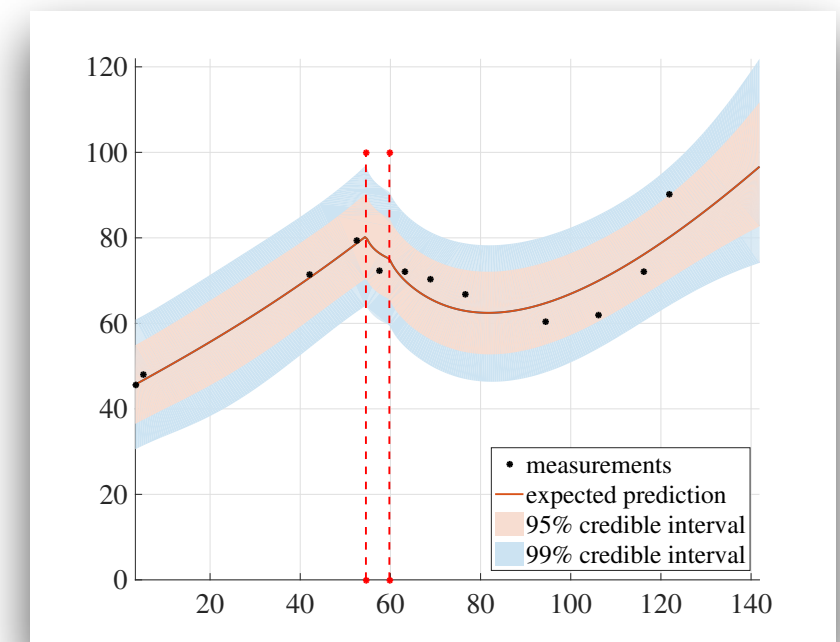
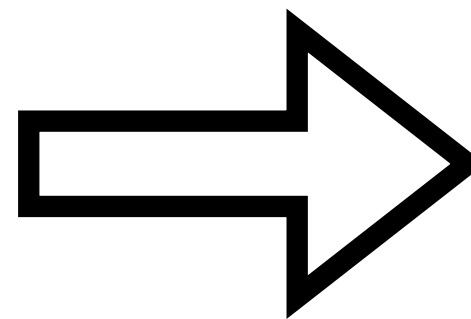
$$\frac{dP}{dt} = \vartheta_4 P \left(1 - \frac{P^*}{K}\right) + \vartheta_5 Q_P - \vartheta_3 P - \vartheta_1 \vartheta_2 C P$$

$$\frac{dQ}{dt} = \vartheta_3 P + \vartheta_1 \vartheta_2 C Q$$

$$\frac{dQ_P}{dt} = \vartheta_1 \vartheta_2 C Q - \vartheta_5 Q_P - \vartheta_6 Q_P$$

$$C(0) = 0, P(0) = \vartheta_7, Q(0) = \vartheta_8, Q_P(0) = 0$$

$$P^*(t) = P(t) + Q(t) + Q_P(t)$$



# APPLICATION

★ data-model assumption

$$d_i = f(t_i; \vartheta) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma)$$

★ likelihood

$$p(\mathbf{d}|\vartheta) = \mathcal{N}(\mathbf{f}(\cdot; \vartheta), \sigma)$$

★ optimization

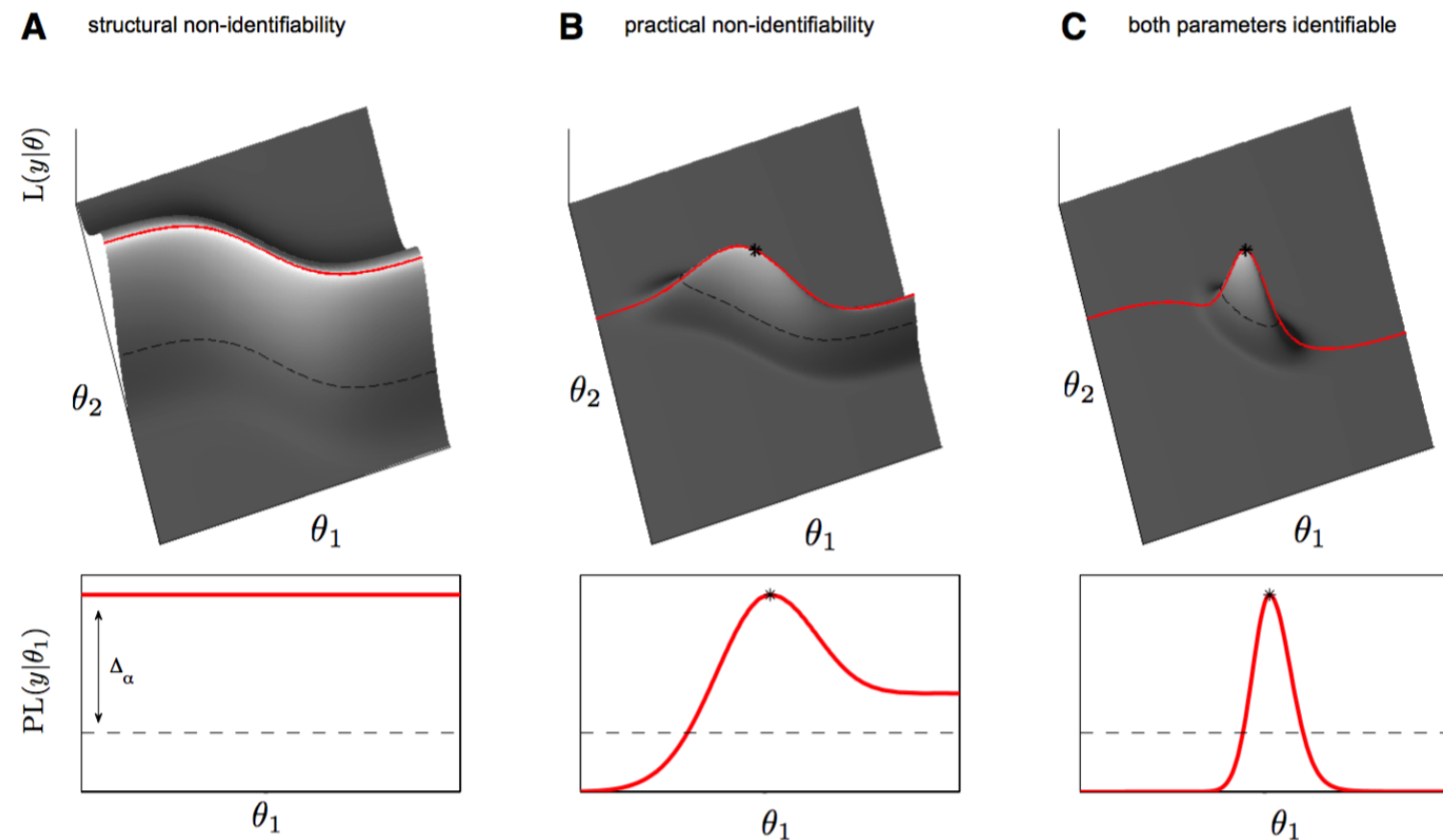
$$\vartheta^* = \operatorname{argmax}_{\vartheta} p(\mathbf{d}|\vartheta)$$

# APPLICATION

with P. Chadjidoukas

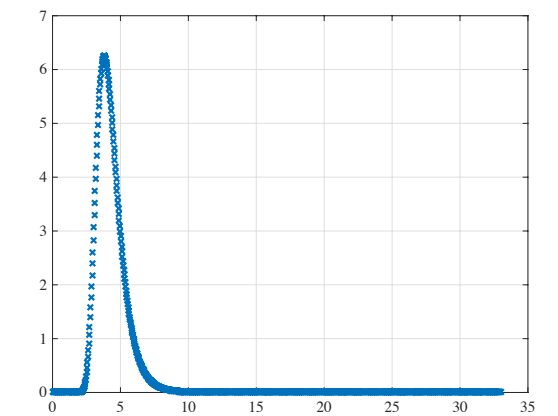
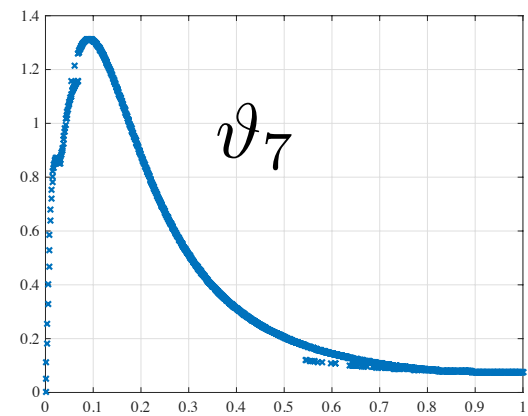
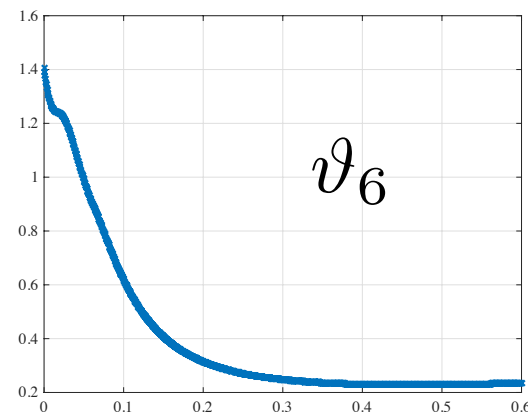
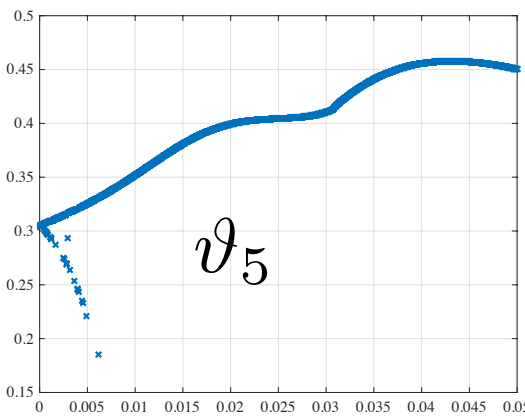
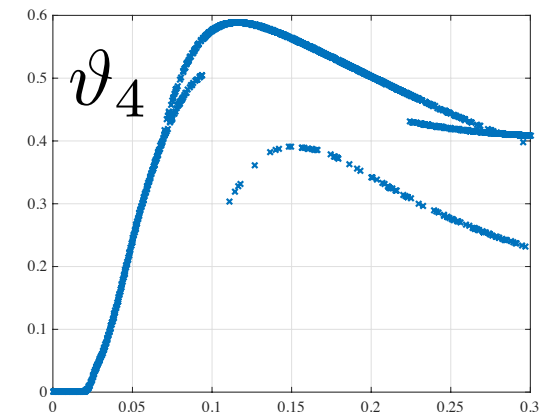
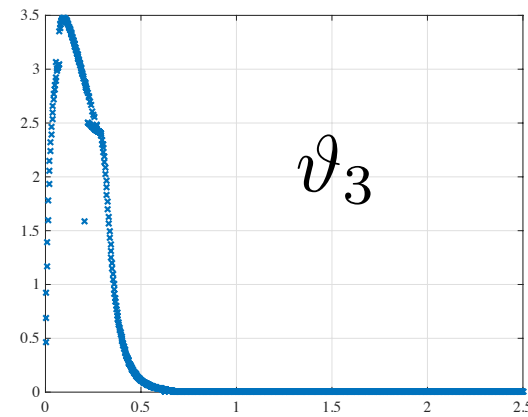
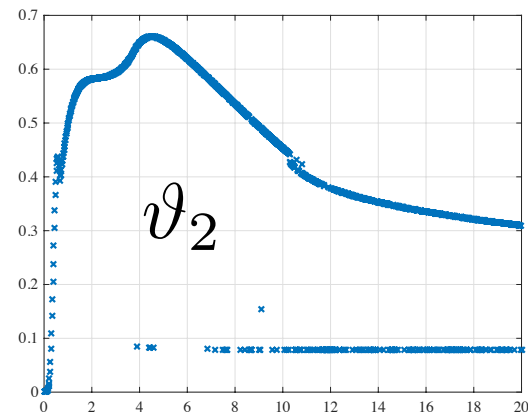
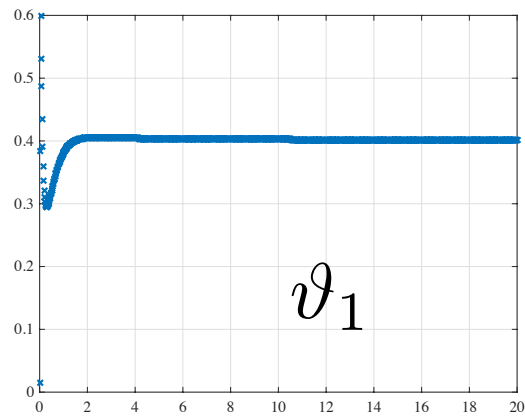
## ★ profile likelihood function

$$q(\mathbf{d}|\vartheta_i) = \max_{\vartheta_{j \neq i}} p(\mathbf{d}|\vartheta)$$



# APPLICATION

with P. Chadjidoukas



# IMPLEMENTATION

\* Π4U

\* High Performance framework for Uncertainty Quantification

\* See the next talk of P. Chadjidoukas

# BIBLIOGRAPHY

- Hansen N. and Ostermeier A. 2001. **Completely Derandomized Self-Adaptation in Evolution Strategies**
- Hansen N., Müller S. D., Koumoutsakos P. 2003. **Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES)**
- Hansen N., **The CMA Evolution Strategy: A Tutorial**

**THANK YOU**