

# Covariance Matrix Adaptation (CMA) Evolution Strategy



*Nikolaus Hansen*

# Motivation::

**Optimization** refers to finding the **best element** from a set of available alternatives

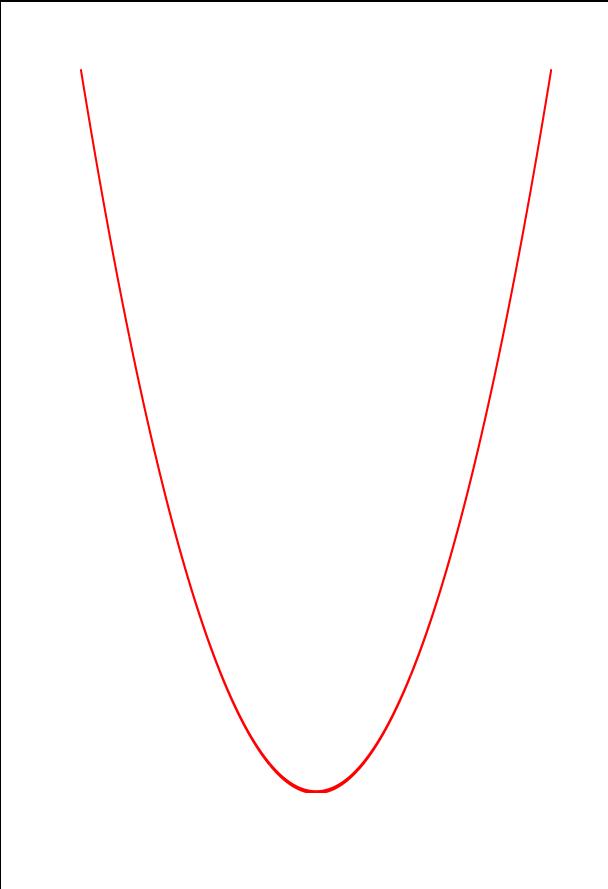
“element”  $\mathbf{x} \in S \subseteq \mathbb{R}^n$

“best”  $f : S \subseteq \mathbb{R}^n \longrightarrow \mathbb{R}$

“finding the best”  $\min_{\mathbf{x} \subseteq S} f(\mathbf{x})$

# Motivation::

Cost function “well behaved” (*analytical definition, differentiable, convex, separable, etc*)

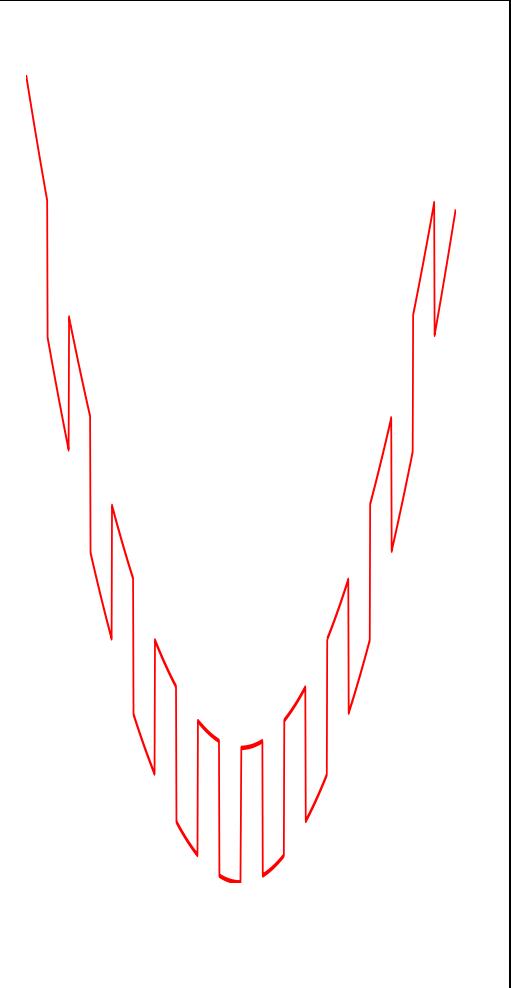
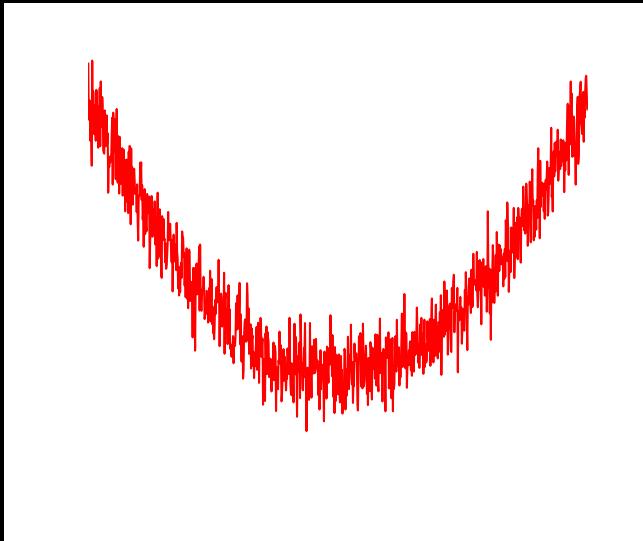
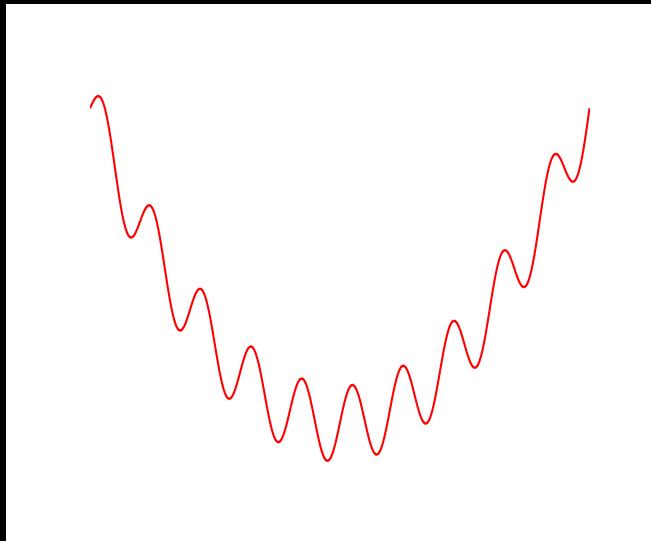


- analytically
- gradient descent methods
- convex programming
- quadratic programming
- nonlinear programming
- etc

<http://www.flickr.com/photos/narnonic/>

# Motivation::

Real world problems (*no derivatives, non-convex, non-separable, non-linear, noisy, etc*)



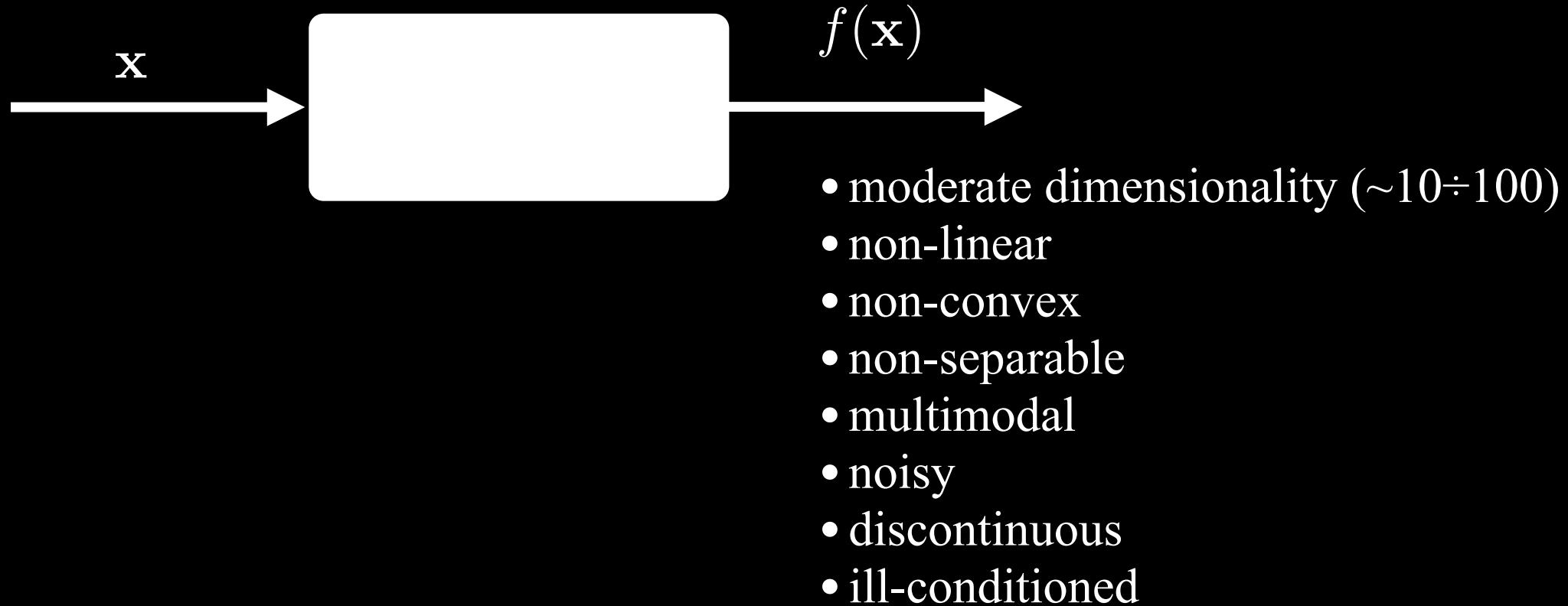
<http://www.flickr.com/photos/pattylagera/>

<http://www.flickr.com/photos/70148269@N00/>

<http://www.flickr.com/photos/castaspella/>

# Motivation::

**Goal:** cope with a real-world black box scenario

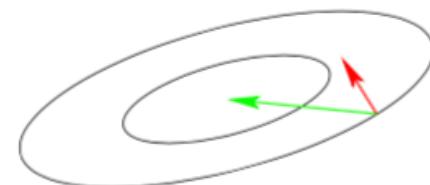
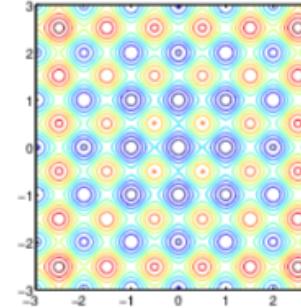
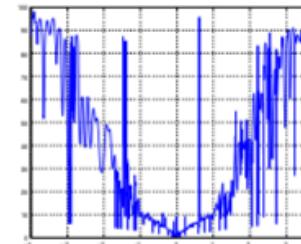
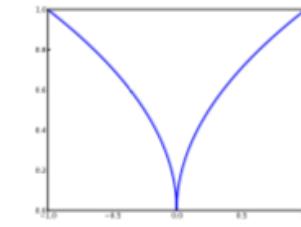


**Approach:** stochastic search, evolutionary strategies

# What Makes a Function Difficult to Solve?

Why stochastic search?

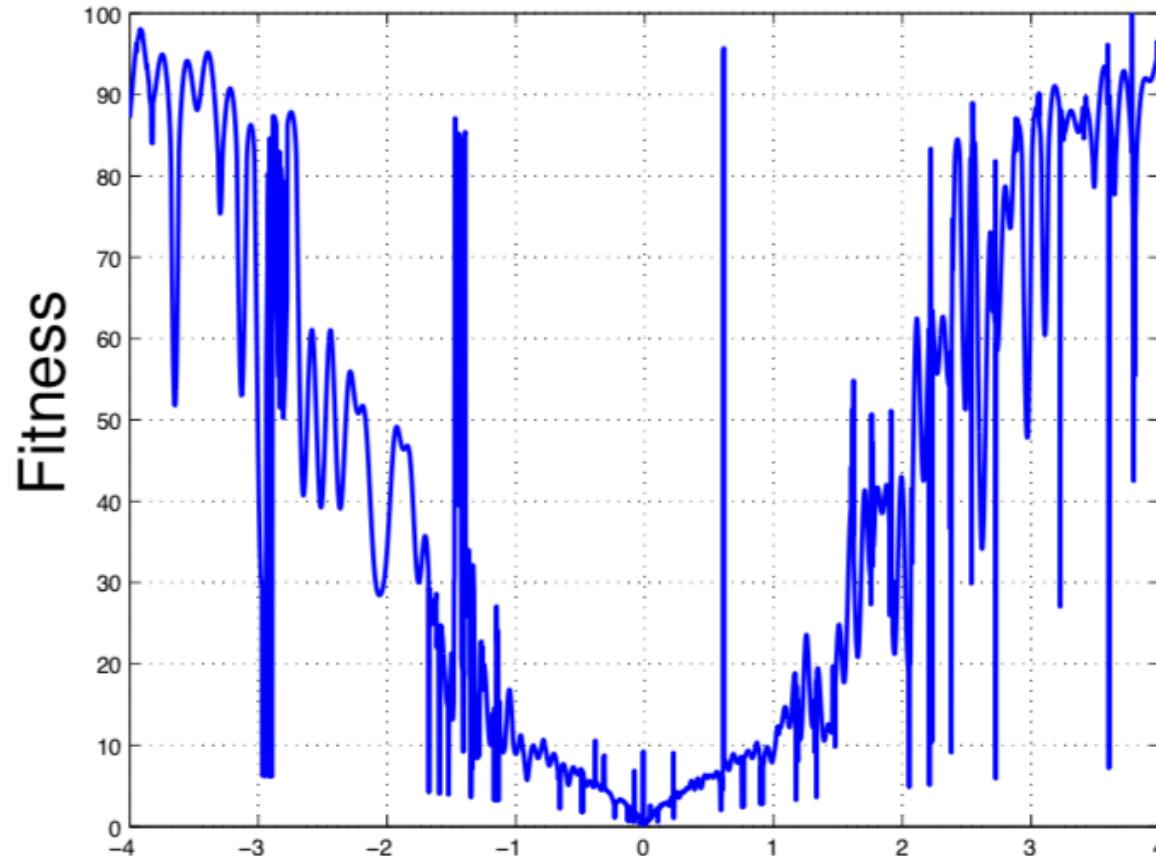
- non-linear, non-quadratic, non-convex  
on linear and quadratic functions much better  
search policies are available
- ruggedness  
non-smooth, discontinuous, multimodal, and/or  
noisy function
- dimensionality (size of search space)  
(considerably) larger than three
- non-separability  
dependencies between the objective variables
- ill-conditioning



gradient direction Newton direction

# Ruggedness

non-smooth, discontinuous, multimodal, and/or noisy



cut from a 5-D example, (easily) solvable with evolution strategies

# Curse of Dimensionality

The term *Curse of dimensionality* (Richard Bellman) refers to problems caused by the **rapid increase in volume** associated with adding extra dimensions to a (mathematical) space.

Example: Consider placing 100 points onto a real interval, say  $[0, 1]$ . Now consider the 10-dimensional space  $[0, 1]^{10}$ . To get **similar coverage** in terms of distance between adjacent points would require  $100^{10} = 10^{20}$  points. A 100 points appear now as isolated points in a vast empty space.

Remark: **distance measures** break down in higher dimensionalities (the central limit theorem kicks in)

Consequence: a **search policy** that is valuable in small dimensions **might be useless** in moderate or large dimensional search spaces.

Example: exhaustive search.

# Separable Problems

## Definition (Separable Problem)

A function  $f$  is separable if

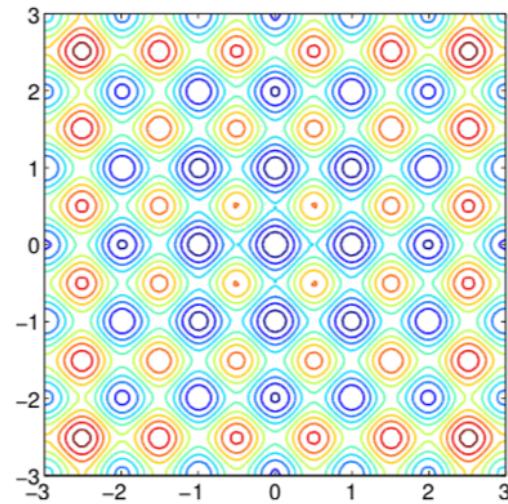
$$\arg \min_{(x_1, \dots, x_n)} f(x_1, \dots, x_n) = \left( \arg \min_{x_1} f(x_1, \dots), \dots, \arg \min_{x_n} f(\dots, x_n) \right)$$

⇒ it follows that  $f$  can be optimized in a sequence of  $n$  independent 1-D optimization processes

## Example: Additively decomposable functions

$$f(x_1, \dots, x_n) = \sum_{i=1}^n f_i(x_i)$$

Rastrigin function



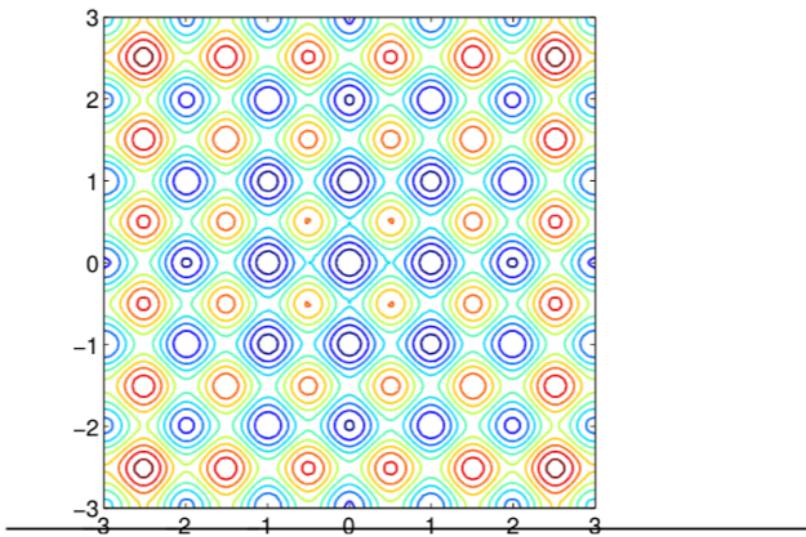
# Non-Separable Problems

Building a non-separable problem from a separable one <sup>(1,2)</sup>

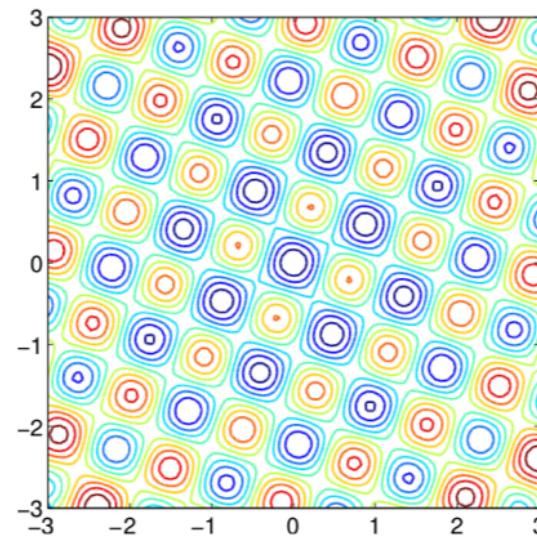
## Rotating the coordinate system

- $f : \mathbf{x} \mapsto f(\mathbf{x})$  separable
- $f : \mathbf{x} \mapsto f(\mathbf{Rx})$  **non-separable**

$\mathbf{R}$  rotation matrix



$\mathbf{R}$   
→



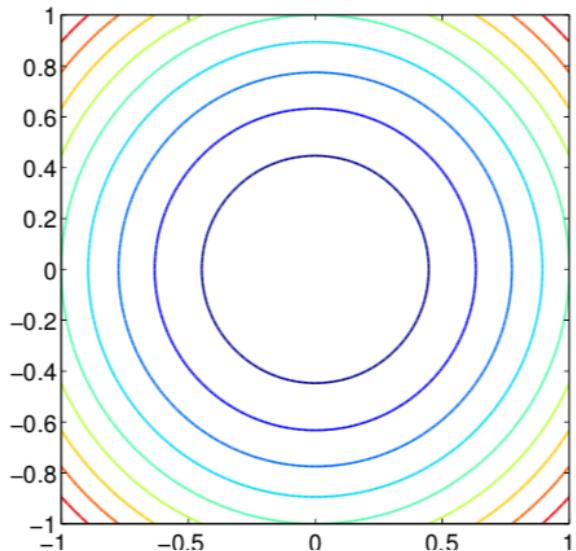
<sup>1</sup> Hansen, Ostermeier, Gawelczyk (1995). On the adaptation of arbitrary normal mutation distributions in evolution strategies: The generating set adaptation. Sixth ICGA, pp. 57-64, Morgan Kaufmann

<sup>2</sup> Salomon (1996). "Reevaluating Genetic Algorithm Performance under Coordinate Rotation of Benchmark Functions; A survey of some theoretical and practical aspects of genetic algorithms." BioSystems, 39(3):263-278

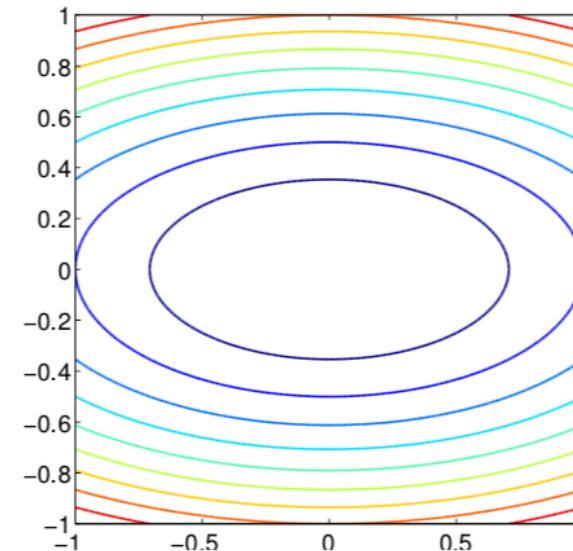
## III-Conditioned Problems

If  $f$  is quadratic,  $f : \mathbf{x} \mapsto \mathbf{x}^T \mathbf{H} \mathbf{x}$ , ill-conditioned means a high condition number of Hessian Matrix  $\mathbf{H}$

ill-conditioned means “**squeezed**” lines of equal function value



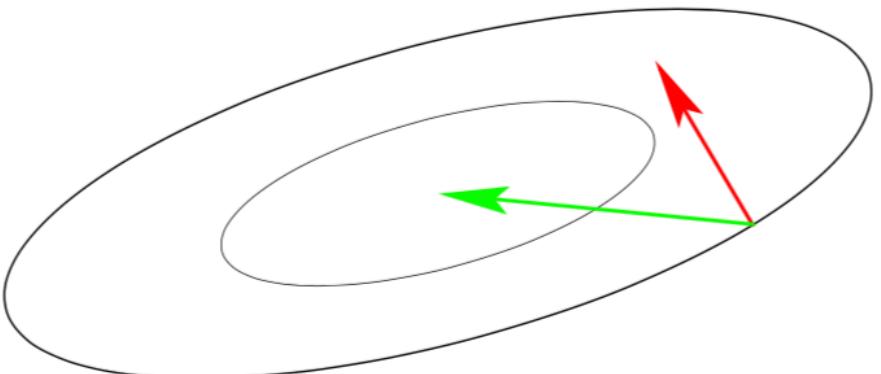
Increased  
→  
condition  
number



consider the curvature of iso-fitness lines

# The Benefit of Second Order Information

Consider the convex quadratic function  $f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T \mathbf{H}(\mathbf{x} - \mathbf{x}^*)$



gradient direction  $-f'(\mathbf{x})^T$

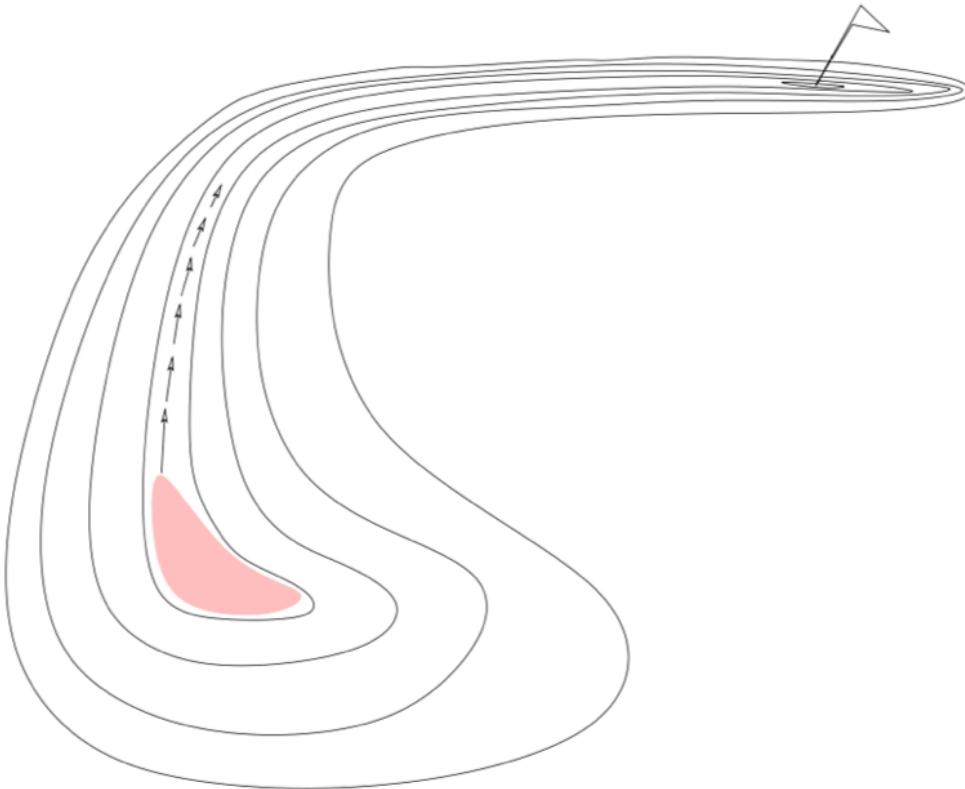
Newton direction  $-\mathbf{H}^{-1}f'(\mathbf{x})^T$

Condition number equals nine here. Condition numbers between 100 and even  $10^6$  can be observed in real world problems.

If  $\mathbf{H} \approx \mathbf{I}$  (small condition number of  $\mathbf{H}$ ) first order information (e.g. the gradient) is sufficient. Otherwise **second order information** (estimation of  $\mathbf{H}^{-1}$ ) **is required**.

# III-Conditioned Problems

Example: A Narrow Ridge



Volume oriented search ends up in the pink area.

To approach the optimum an ill-conditioned problem needs to be solved (e.g. by following the narrow bent ridge).<sup>3</sup>

---

<sup>3</sup> Whitley, Lunacek, Knight 2004. Ruffled by Ridges: How Evolutionary Algorithms Can Fail, GECCO ▶ ◀ ⏪ ⏩ ⏴ ⏵ ⏵ ⏵

# What Makes a Function Difficult to Solve?

... and what can be done

The Problem	Possible Approaches
Dimensionality	exploiting the problem structure separability, locality/neighborhood, encoding
Ill-conditioning	second order approach changes the neighborhood metric
Ruggedness	<b>non-local</b> policy, large sampling width (step-size) as large as possible while preserving a reasonable convergence speed  <b>population-based</b> method, stochastic, non-elitistic recombination operator serves as repair mechanism
	restarts

... metaphors

# Stochastic search & CMA::

Black box template to minimize  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

set variable strategy parameters  $\theta$

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

loop

1) sample distribution  $P(\mathbf{x}|\theta) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_\lambda \in \mathbb{R}^n$

2) evaluate  $f(\mathbf{x}_1), \dots, f(\mathbf{x}_\lambda)$

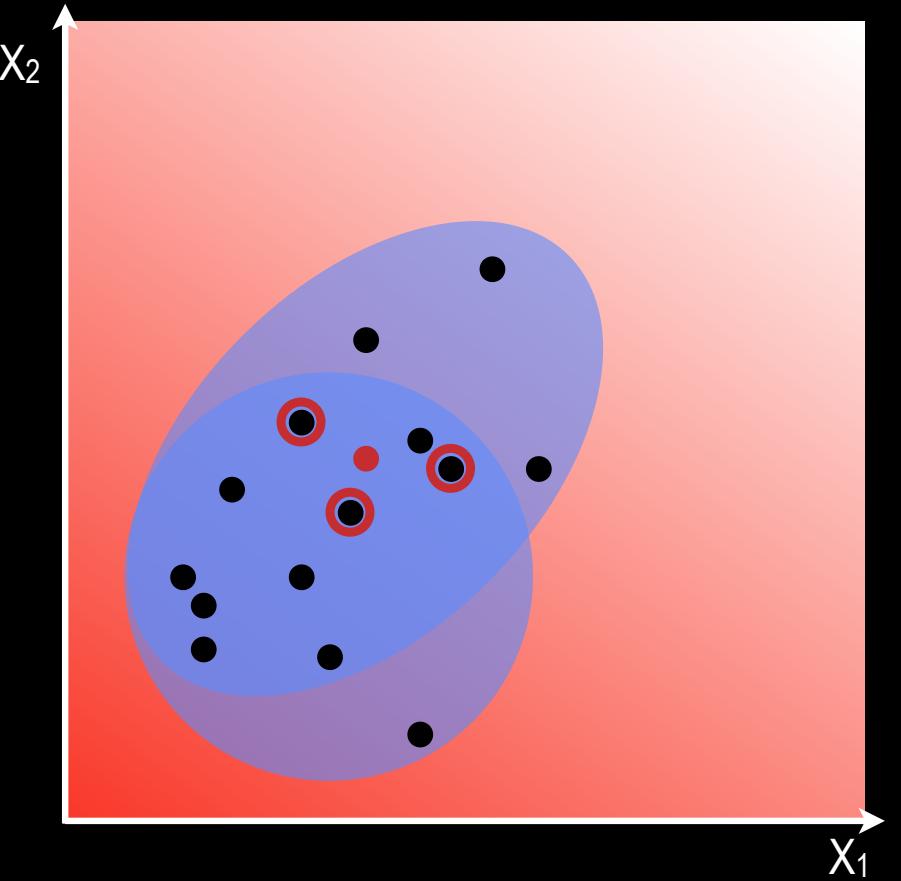
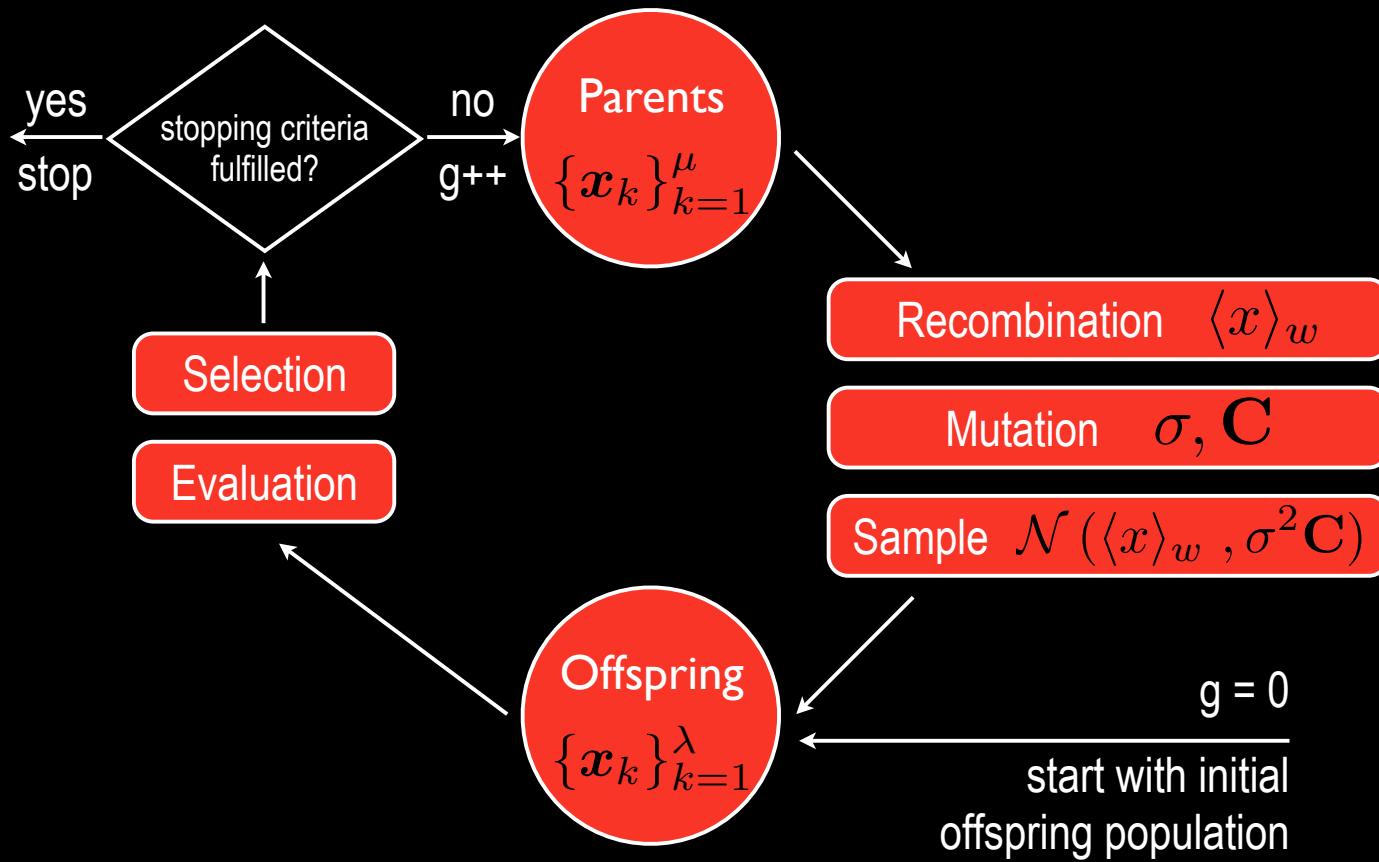
3) update  $\theta \leftarrow F_\theta(\theta, \mathbf{x}_1, \dots, \mathbf{x}_\lambda, f(\mathbf{x}_1), \dots, f(\mathbf{x}_\lambda))$

## CMA-ES

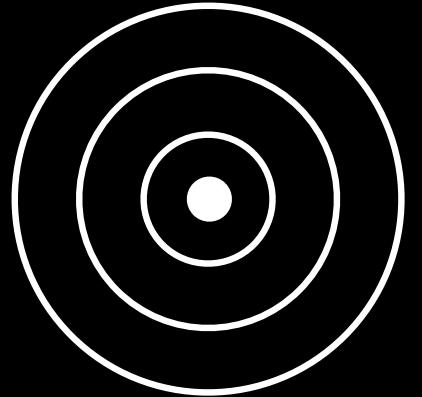
- P is a multivariate normal distribution  $\mathcal{N}(\mathbf{m}, \sigma^2 \mathbf{C})$
- $\theta = \{\mathbf{m}, \sigma, \mathbf{C}\} \in \mathbb{R}^n \times \mathbb{R}^{n \times n} \times \mathbb{R}_+$
- $F_\theta = F_\theta(\theta, \mathbf{x}_{1:\lambda}, \dots, \mathbf{x}_{\mu:\lambda})$  where  $\mathbf{x}_{i:\lambda}$  is the i-th best point and  $\mu \leq \lambda$

# Covariance Matrix Adaptation - ES::

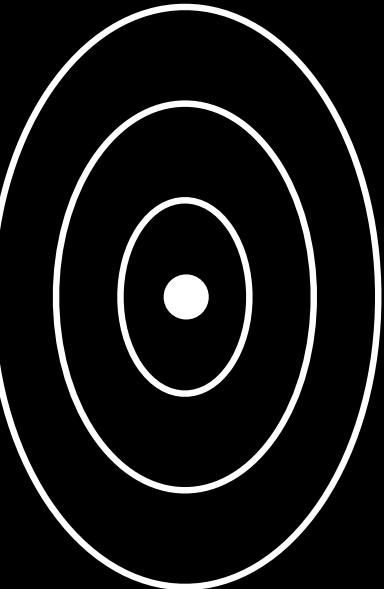
- black box approach (no  $\nabla\Phi$ ,  $\Phi$  = cost function)
- iterative methods operating with **populations** of candidate solutions



# Sampling multivariate normal distribution::

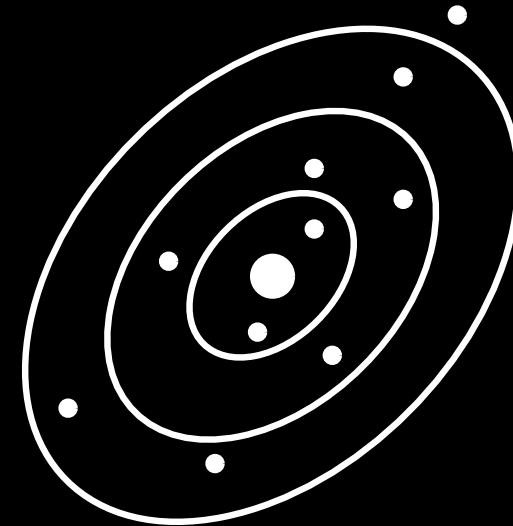


$$\mathcal{N}(\mathbf{m}, \sigma^2 \mathbf{I}) \sim \mathbf{m} + \sigma \mathcal{N}(\mathbf{0}, \mathbf{I})$$



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

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

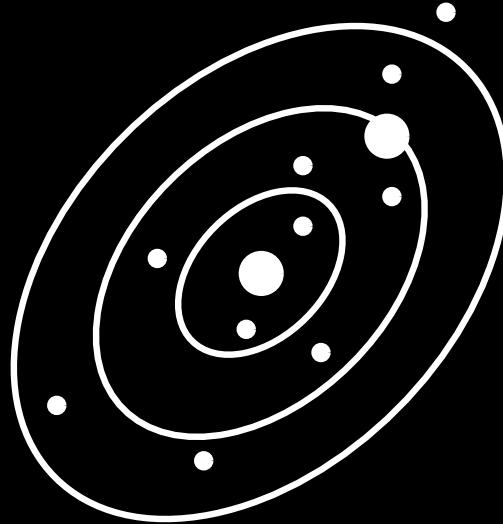


# Recombination:: Update $m$ ::

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

$$f(\mathbf{x}_{1:\lambda}) \leq \dots \leq f(\mathbf{x}_{\lambda:\lambda})$$

⋮

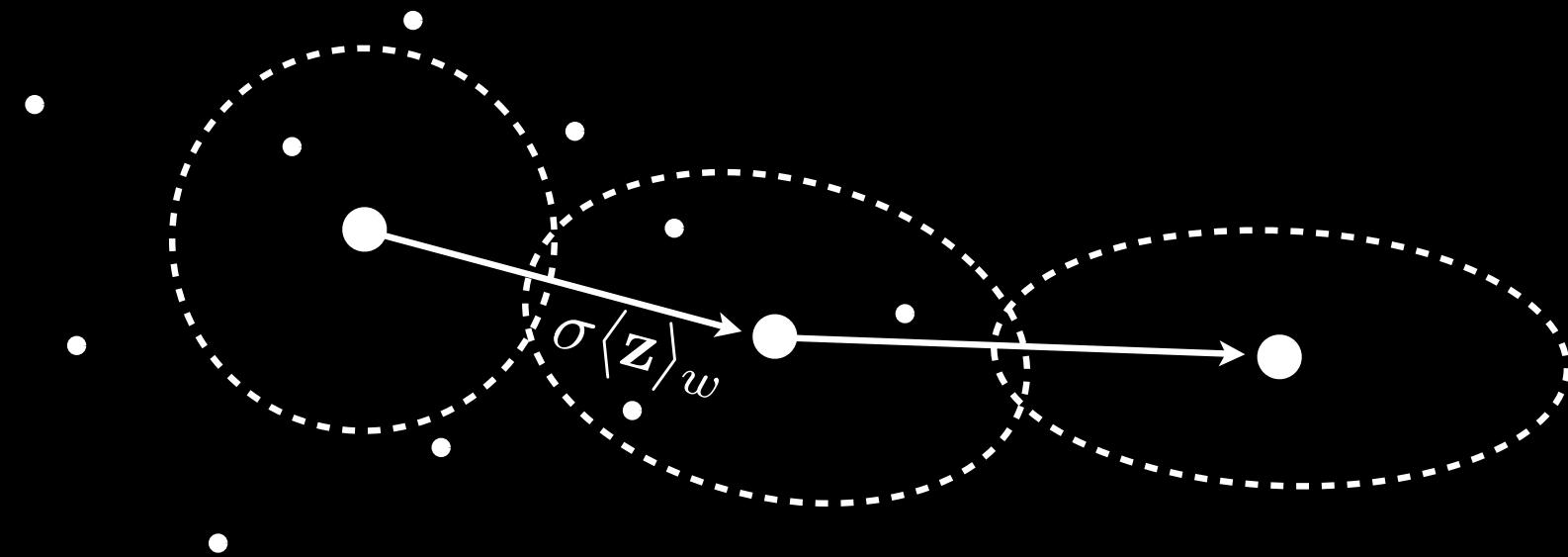


$$\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda} = \mathbf{m} + \sigma \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda} = \mathbf{m} + \sigma \langle \mathbf{z} \rangle_w$$

$$\sum_{i=1}^{\mu} w_i = 1 \quad w_1 \geq \dots \geq w_{\mu} \geq 1$$

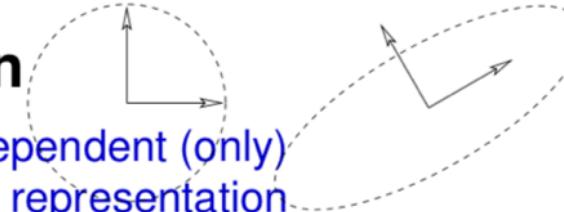
# Mutation:: Update C:: Rank-one update

Mutate **C** to increase the probability of successful steps to appear again



$$\mathbf{C} \leftarrow (1 - LR_{cov})\mathbf{C} + LR_{cov}\langle \mathbf{z} \rangle_w \langle \mathbf{z} \rangle_w^T$$

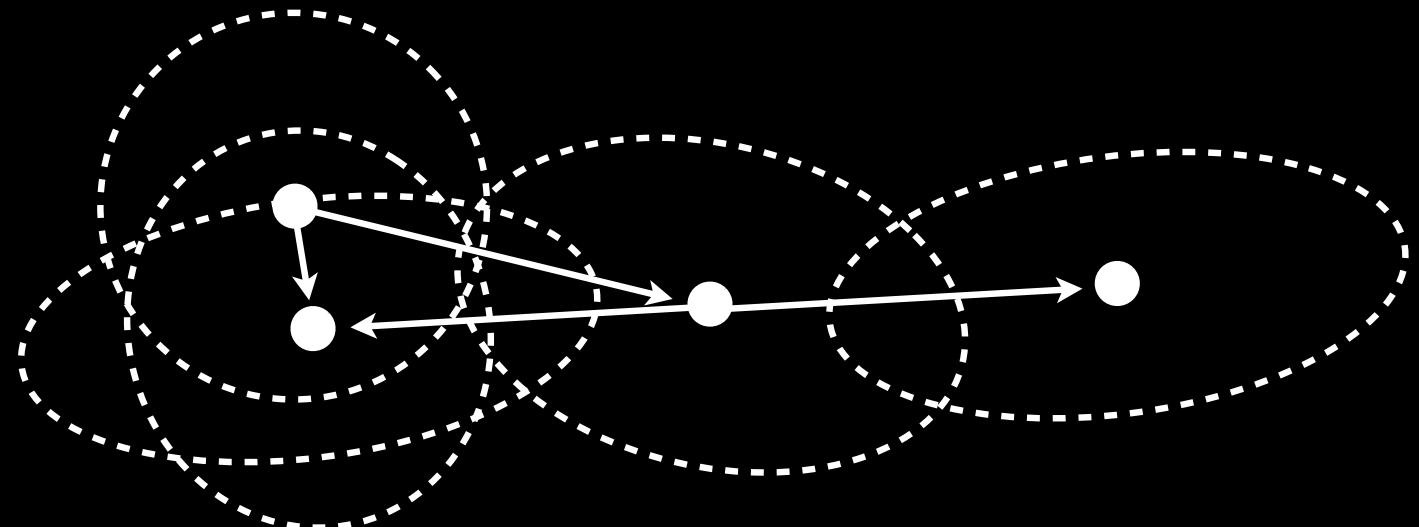
## covariance matrix adaptation

- learns all **pairwise dependencies** between variables  
off-diagonal entries in the covariance matrix reflect the dependencies
- conducts a **principle component analysis** (PCA) of steps  $y_w$ , sequentially in time and space  
eigenvectors of the covariance matrix  $\mathbf{C}$  are the principle components / the principle axes of the mutation ellipsoid
- learns a new **rotated problem representation**  
components are independent (only) in the new representation
- learns a **new (Mahalanobis) metric**  
variable metric method
- approximates the **inverse Hessian** on quadratic functions  
transformation into the sphere function
- for  $\mu = 1$ : conducts a **natural gradient ascent** on the distribution  $\mathcal{N}$   
entirely independent of the given coordinate system

**Mahalanobis metric:** It is a multi-dimensional generalization of the idea of measuring how many standard deviations away  $P$  is from the mean of  $D$ . This distance is zero if  $P$  is at the mean of  $D$ , and grows as  $P$  moves away from the mean along each principal component axis.

# Mutation:: Update C:: Evolution path

$$\langle \mathbf{z} \rangle_w \langle \mathbf{z} \rangle_w^T = -\langle \mathbf{z} \rangle_w (-\langle \mathbf{z} \rangle_w)^T$$



$$\mathbf{p}_c \leftarrow (1 - LR_c)\mathbf{p}_c + LR_c \langle \mathbf{z} \rangle_w$$

# Recap::

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

$$\mathbf{m} \longleftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_w \quad \langle \mathbf{z} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}$$

$$\mathbf{p}_c \longleftarrow (1 - c_c) \mathbf{p}_c + \sqrt{1 - (1 - c_c)^2} \sqrt{\frac{1}{\sum_{i=1}^{\mu} w_i^2}} \langle \mathbf{z} \rangle_w$$

$$\mathbf{C} \longleftarrow (1 - c_{cov}) \mathbf{C} + c_{cov} \mathbf{p}_c \mathbf{p}_c^T$$

# Recap:: Rank-one + cumulation

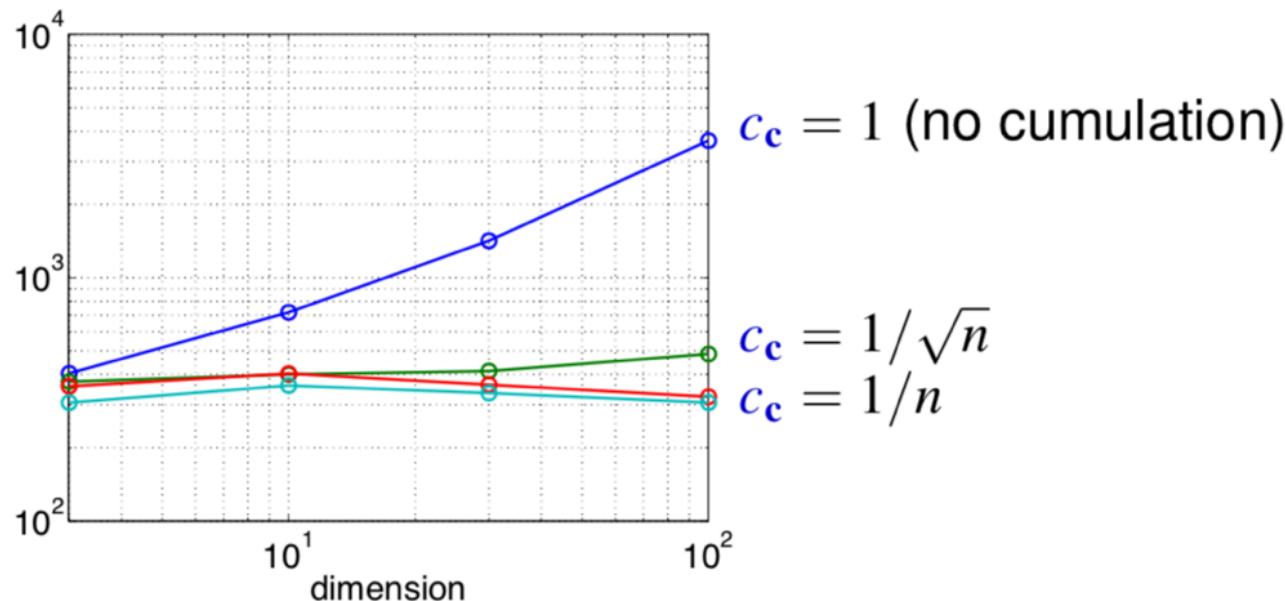
Evolution path + rank-one update reduces the number of function evaluations from  
 $O(n^2)$  to  $O(n)$

*model complexity is  $n^2$  but an important part of it can be learnt in  $n$  function evaluations*

Using an **evolution path** for the **rank-one update** of the covariance matrix reduces the number of function evaluations to adapt to a straight ridge **from about**  $\mathcal{O}(n^2)$  **to**  $\mathcal{O}(n)$ .<sup>(a)</sup>

<sup>a</sup>Hansen & Auger 2013. Principled design of continuous stochastic search: From theory to practice.

Number of  $f$ -evaluations divided by dimension on the cigar function  $f(\mathbf{x}) = x_1^2 + 10^6 \sum_{i=2}^n x_i^2$



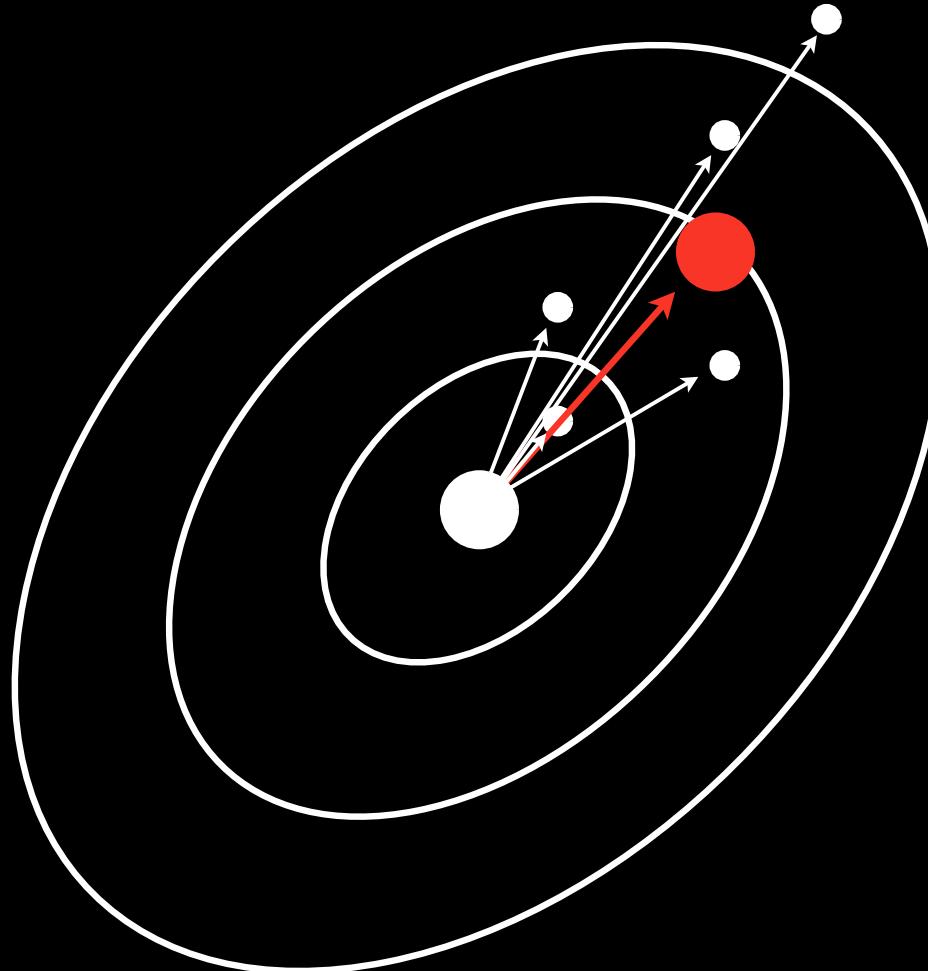
The overall model complexity is  $n^2$  but important parts of the model can be learned in time of order  $n$

# Mutation:: Update C:: Rank- $\mu$ update

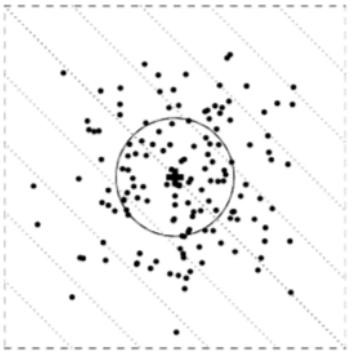
There is more information out there

$$\mathbf{Z} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda} \mathbf{z}_{i:\lambda}^T$$

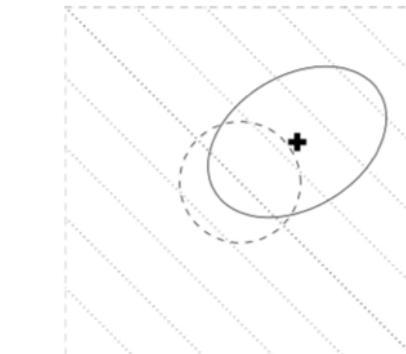
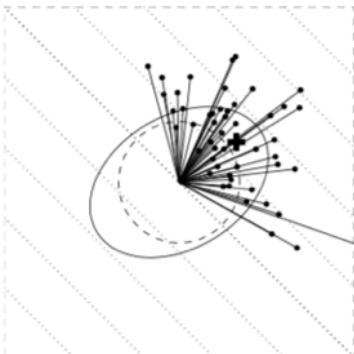
$$\mathbf{C} \leftarrow (1 - c_{cov})\mathbf{C} + c_{cov}\mathbf{Z}$$



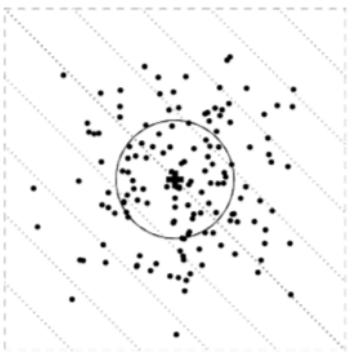
## Rank- $\mu$ CMA versus Estimation of Multivariate Normal Algorithm EMNA<sub>global</sub><sup>11</sup>



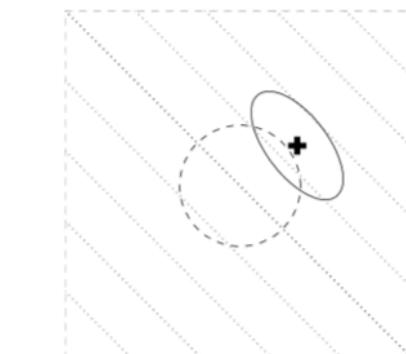
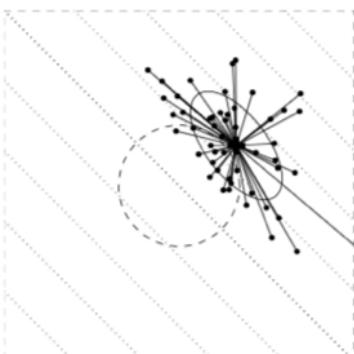
$$x_i = \mathbf{m}_{\text{old}} + y_i, \quad y_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad \mathbf{C} \leftarrow \frac{1}{\mu} \sum (x_{i:\lambda} - \mathbf{m}_{\text{old}})(x_{i:\lambda} - \mathbf{m}_{\text{old}})^T$$



$$\mathbf{m}_{\text{new}} = \mathbf{m}_{\text{old}} + \frac{1}{\mu} \sum y_{i:\lambda}$$



$$x_i = \mathbf{m}_{\text{old}} + y_i, \quad y_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$$



$$\mathbf{m}_{\text{new}} = \mathbf{m}_{\text{old}} + \frac{1}{\mu} \sum y_{i:\lambda}$$

sampling of  $\lambda = 150$  solutions (dots)

calculating  $\mathbf{C}$  from  $\mu = 50$  solutions

new distribution

$\mathbf{m}_{\text{new}}$  is the minimizer for the variances when calculating  $\mathbf{C}$

rank- $\mu$  CMA  
conducts a  
**PCA of steps**

EMNA<sub>global</sub>  
conducts a  
**PCA of points**

<sup>11</sup> Hansen, N. (2006). The CMA Evolution Strategy: A Comparing Review. In J.A. Lozano, P. Larranga, I. Inza and E. Bengoetxea (Eds.). Towards a new evolutionary computation. Advances in estimation of distribution algorithms. pp. 75-102

# Recap::

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

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_w \quad \langle \mathbf{z} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}$$

$$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \sqrt{1 - (1 - c_c)^2} \sqrt{\frac{1}{\sum_{i=1}^{\mu} w_i^2}} \langle \mathbf{z} \rangle_w$$

$$\mathbf{C} \leftarrow (1 - c_{cov}) \mathbf{C}$$

$$+ c_{cov} \frac{1}{\mu_{cov}} \mathbf{p}_c \mathbf{p}_c^T$$

$$+ c_{cov} \left( 1 - \frac{1}{\mu_{cov}} \right) \mathbf{Z}$$

$$\mu_{cov} = \sqrt{\frac{1}{\sum_{i=1}^{\mu} w_i^2}}$$

$$\mathbf{Z} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda} \mathbf{z}_{i:\lambda}^T$$

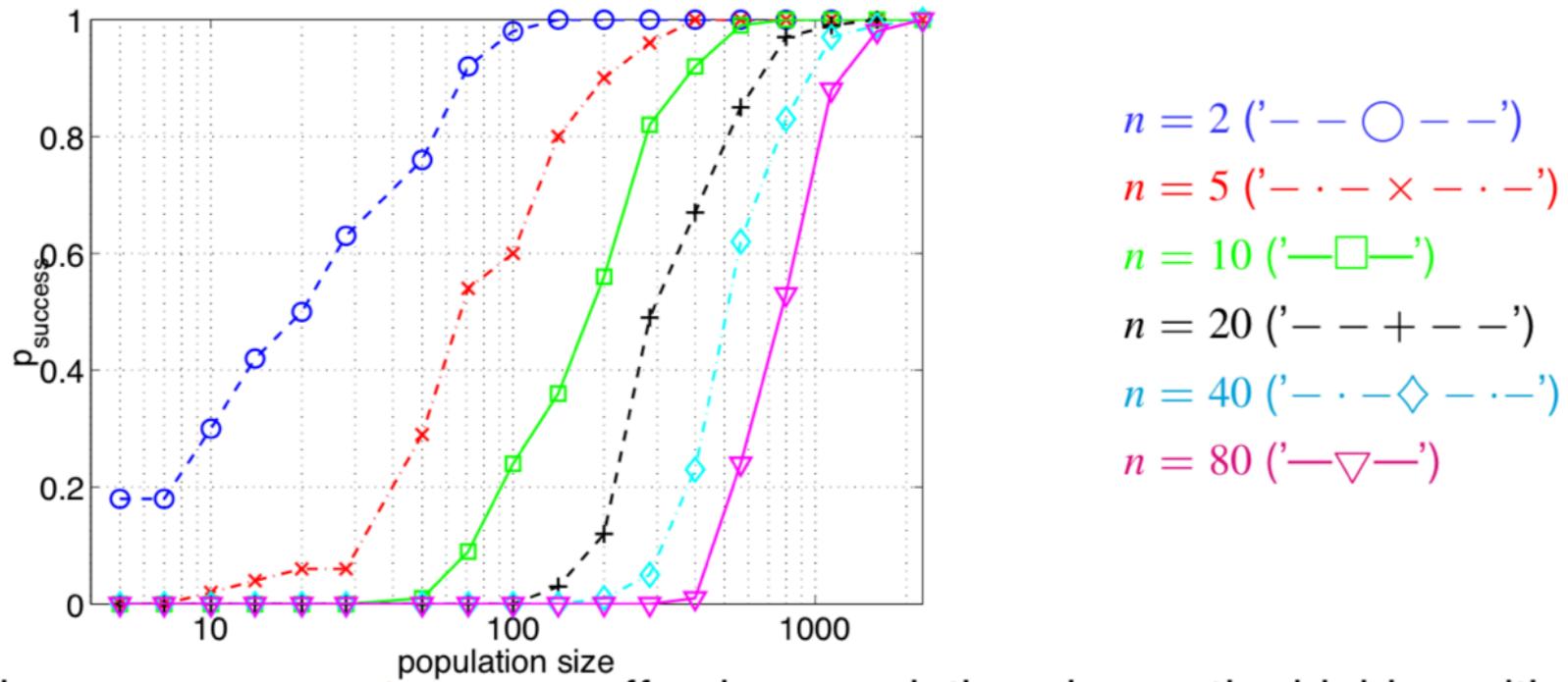
# Recap:: Rank- $\mu$ update

Rank- $\mu$  increases possible learning rates

Rank- $\mu$  reduces the number of **generations** roughly from  $O(n^2)$  to  $O(n)$

# Population Size on Multi-Modal Functions

Success Probability to Find the Global Optimum

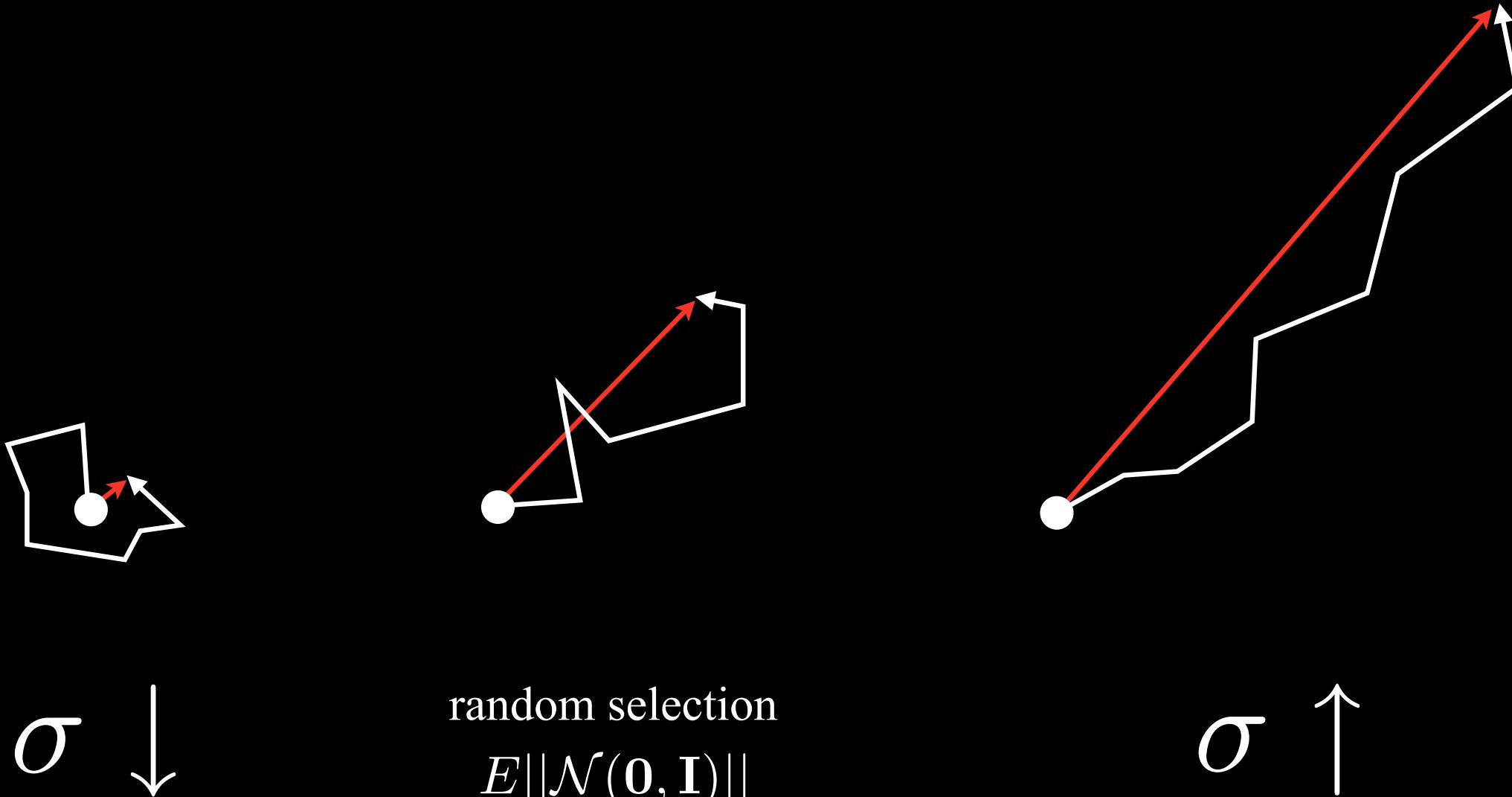


Shown: **success rate** versus offspring population size on the highly multi-modal Rastrigins function<sup>7</sup>

On multi-modal functions increasing the population size can sharply increase the success probability to find the global optimum

<sup>7</sup> Hansen & Kern 2004. Evaluating the CMA Evolution Strategy on Multimodal Test Functions. PPSN VIII, Springer-Verlag, pp. 282-291.

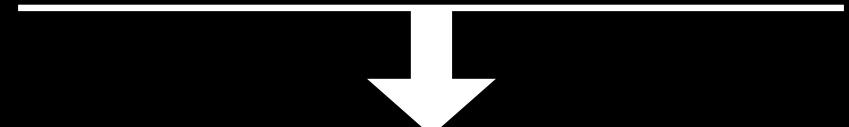
# Mutation:: Update $\sigma$ :: Evolution path



# Mutation:: Update $\sigma$ :: Evolution path

$$\mathbf{p}_\sigma \leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \sqrt{1 - (1 - c_\sigma)^2} \sqrt{\sum_{i=1}^{\mu} w_i^2} \mathbf{C}^{-\frac{1}{2}} \langle \mathbf{z} \rangle_w$$

$$\sigma \leftarrow \sigma e^{\left( \frac{1}{d_\sigma} \left( \frac{||p_\sigma||}{E ||\mathcal{N}(\mathbf{0}, \mathbf{I})||} - 1 \right) \right)}$$



$$> 1 \iff ||p_\sigma|| > E ||\mathcal{N}(\mathbf{0}, \mathbf{I})||$$

# Recap::

$$\mathbf{x}_i = m + \sigma \mathbf{z}_i$$

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_w$$

$$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \sqrt{1 - (1 - c_c)^2} \sqrt{\sum_{i=1}^{\mu} w_i^2} \langle \mathbf{z} \rangle_w$$

$$\mathbf{C} \leftarrow (1 - c_{cov}) \mathbf{C} + c_{cov} \frac{1}{\mu_{cov}} \mathbf{p}_c \mathbf{p}_c^T + c_{cov} \left( 1 - \frac{1}{\mu_{cov}} \right) \mathbf{Z}$$

$$\mathbf{p}_{\sigma} \leftarrow (1 - c_{\sigma}) \mathbf{p}_{\sigma} + \sqrt{1 - (1 - c_{\sigma})^2} \sqrt{\sum_{i=1}^{\mu} w_i^2} \mathbf{C}^{-\frac{1}{2}} \langle \mathbf{z} \rangle_w$$

$$\sigma \leftarrow \sigma e \left( \frac{1}{d_{\sigma}} \left( \frac{||p_{\sigma}||}{E ||\mathcal{N}(\mathbf{0}, \mathbf{I})||} - 1 \right) \right)$$

$$\mathbf{z}_i \sim_{\mu} \mathcal{N}(\mathbf{0}, \mathbf{C})$$

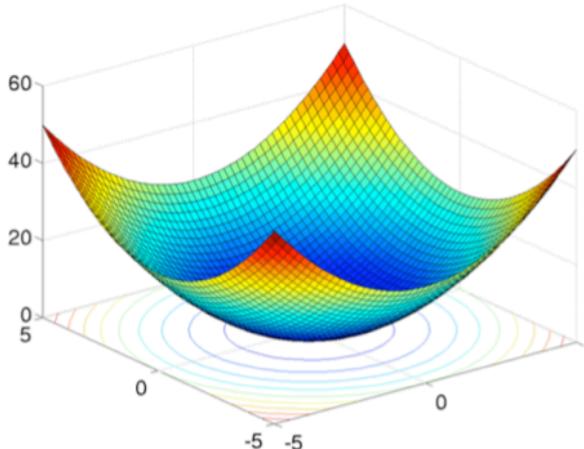
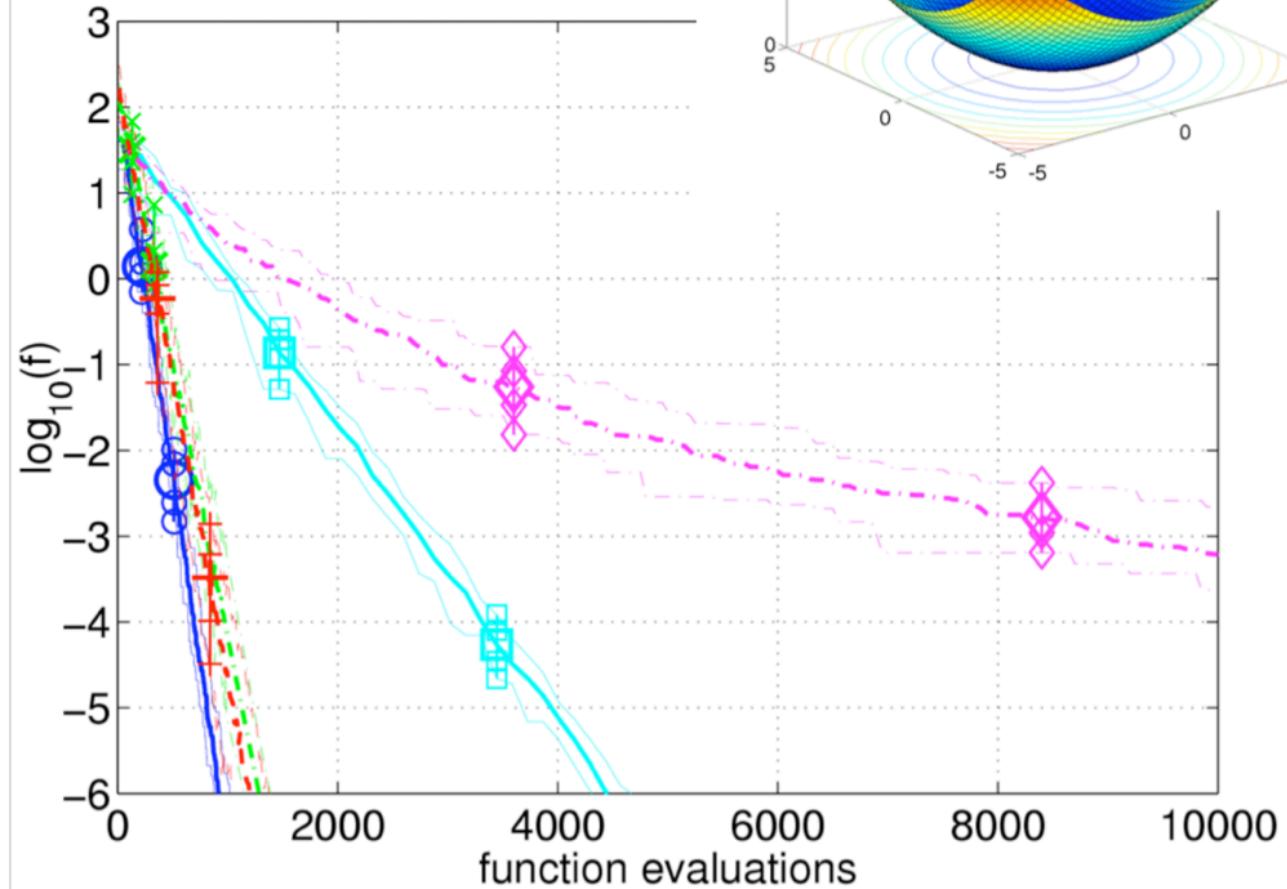
$$\langle \mathbf{z} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}$$

$$\mathbf{Z} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda} \mathbf{z}_{i:\lambda}^T$$

$$\mu_{cov} = \sqrt{\frac{1}{\sum_{i=1}^{\mu} w_i^2}}$$

# Sphere function, 10D

$$f_{Sphere}(\mathbf{x}) = \sum_{i=1}^n x_i^2$$



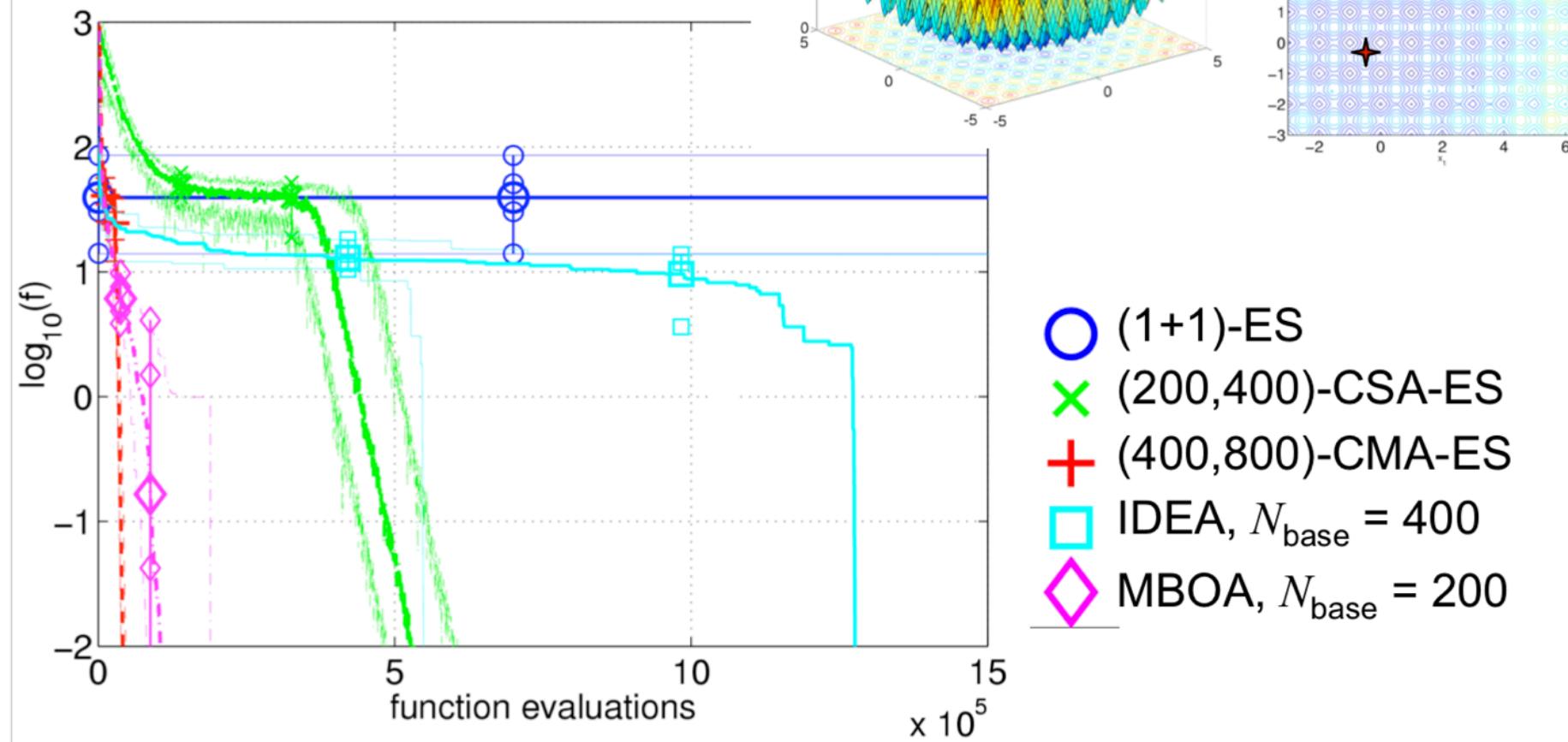
- Separable
- Unimodal
- Well scaled

- (1+1)-ES
- (5,10)-CSA-ES
- (5,10)-CMA-ES
- IDEA,  $N_{base} = 200$
- MBOA,  $N_{base} = 100$

# Rastrigin function, 10D

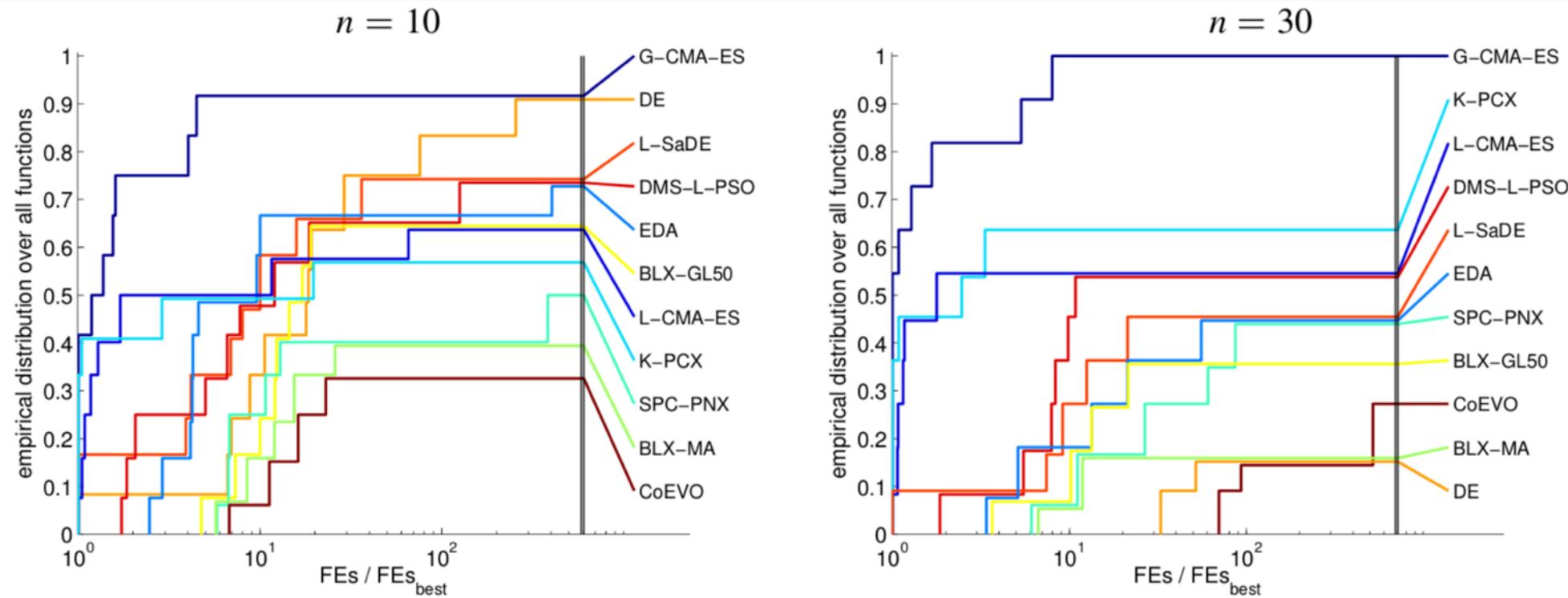
$$f_{Rastrigin}(\mathbf{x}) = 10n + \sum_{i=1}^n (x_i^2 - 10\cos(2\pi \times x_i))$$

■ Multimodal - Separable



# Summarized Results

## Empirical Distribution of Normalized Success Performance



$FEs = \text{mean}(\#fevals) \times \frac{\#\text{all runs (25)}}{\#\text{successful runs}}$ , where `#fevals` includes only successful runs.

Shown: **empirical distribution function** of the Success Performance  $FEs$  divided by  $FEs$  of the best algorithm on the respective function.

Results of all functions are used where at least one algorithm was successful at least once, i.e. where the target function value was reached in at least one experiment (out of  $11 \times 25$  experiments).

Small values for  $FEs$  and therefore large (cumulative frequency) values in the graphs are preferable.

# CMA internal parameters::

Lots of pain, testing and benchmarking

Just believe

# Recap::

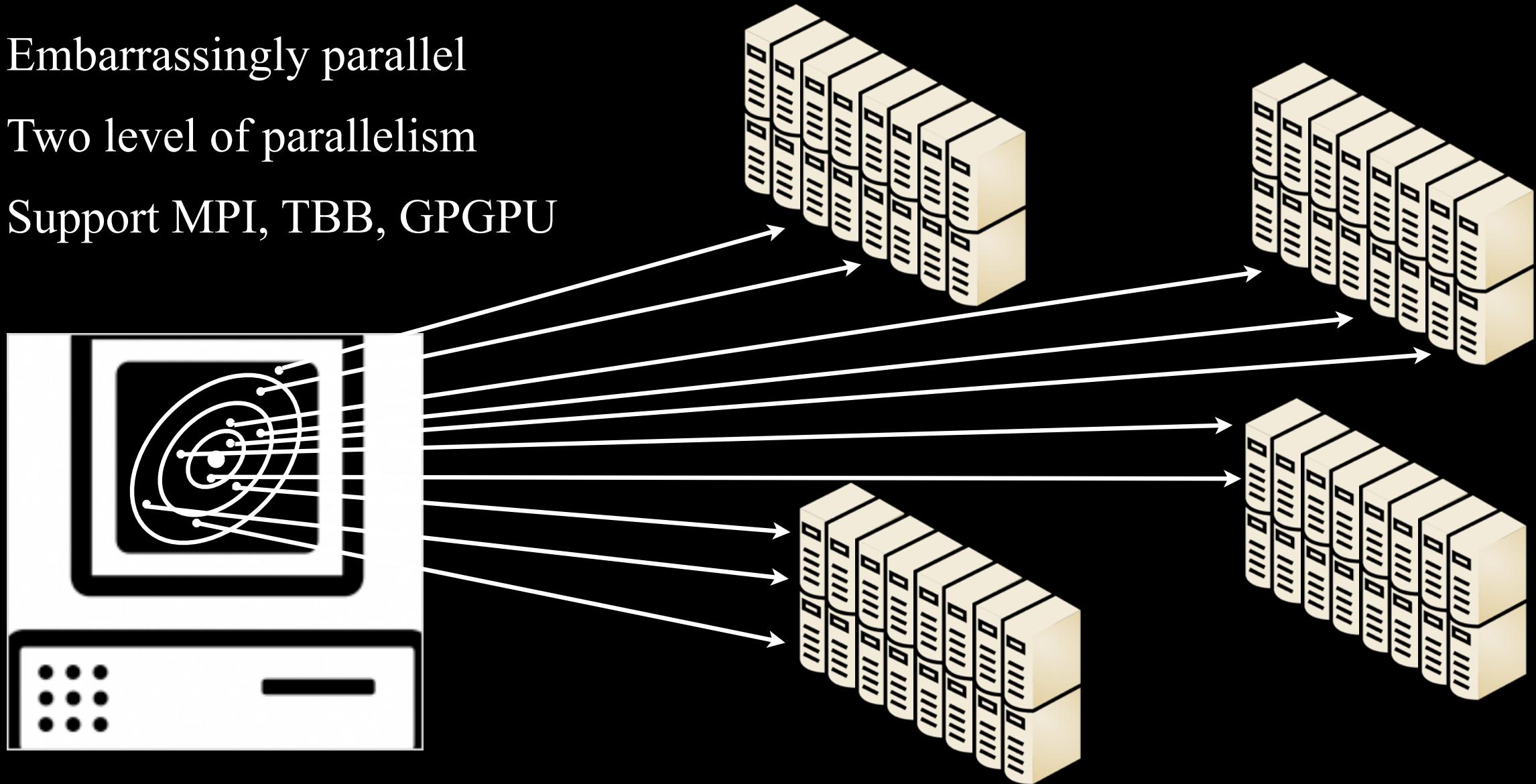
Overall control over on normal multivariate scaling

Prevents premature convergence

On quadratic functions step length is close to optimal

# Multihost CMA-ES::

- Embarrassingly parallel
- Two level of parallelism
- Support MPI, TBB, GPGPU



# Conclusions::

- Multivariate normal distribution sampling
- Rank based selection
- Covariance matrix adaptation to maximize probability of repeating successful steps
- Step length control
- Allows massively parallel computation

CMA-ES is the algorithm of choice for moderate dimensional non-linear, non-convex, non-separable, multimodal, noisy and discontinuous real valued functions

# References::

- We credit Nikolaus Hansen for many of the slides used in the CMA discussion, and much of the material is based on his publications and software page
- Hansen, Nikolaus, and Andreas Ostermeier. "Completely derandomized self-adaptation in evolution strategies." *Evolutionary computation* 9.2 (2001): 159-195.
- Hansen, Nikolaus. "The CMA evolution strategy: a comparing review." *Towards a new evolutionary computation*. Springer, Berlin, Heidelberg, 2006. 75-102.
- Hansen, Nikolaus. "The CMA evolution strategy: A tutorial." arXiv preprint arXiv: 1604.00772 (2016).