

# Design of a Surrogate Model Assisted $(\mu/\mu, \lambda)$ -ES

[Honours Thesis]

JINGYUN YANG, Faculty of Computer Science, Dalhousie University

---

Surrogate models have been widely used to assist evolutionary algorithms (EAs) to avoid unnecessary objective function evaluations. But those surrogate assisted EAs are usually complicated and the behaviours of the algorithms are not well understood. A recent analysis of a surrogate model assisted  $(1+1)$ -ES has helped understand the behaviour of the algorithm and resulted in a step size adaptation mechanism. The goal of this thesis is to conduct a similar analysis for  $(\mu/\mu, \lambda)$ -ES that potentially more fully exploits the surrogate model in the sense that a population of candidate solutions are evaluated by the surrogate in each iteration. It is unclear whether any additional performance advantage can be derived from this.

Additional Key Words and Phrases:  $(\mu/\mu, \lambda)$ -ES, Surrogate Model, Evolutionary algorithms (EAs), Gaussian Process

---

## 1 INTRODUCTION

*Evolution strategies (ESs)* have been widely utilized to solve optimization problems where the true objective function evaluation is computationally-intensive. ES is flexible and able to solve many optimization problems from two aspects, variation and selection. Firstly, using a stochastic variation from mutation (random sampling of new directions) and recombination (combine the selected mutations) can introduce new unbiased information that may help explore the search space via generating new offspring. Secondly, search using a population of candidate solutions is more robust under moderate noise in multi-modal optimizations, as opposed to some classical search methods like quasi-Newton. Besides, applying a selection on the population can extract potential good step information that may help solve the optimization problem.

Various attempts have been made to reduce the cost by extracting information obtained from points evaluated in previous iterations. Such information yields insights into an efficient selection that help generate potential promising offspring. One way is to use a surrogate model; an approximation model trained based on the candidate solutions evaluated by the true objective function in previous iterations. The surrogate model acts as a substitution of the true objective function that gives an inaccurate estimate of the objective function value at a much lower cost compared to using the exact objective function. Despite the computation saving of applying surrogate modelling, the estimated objective function value may contain a model bias that can affect both the step size being adapted and the direction selected. Therefore, surrogate modelling is helpful if the computational saving in using the true objective function outshines the potential poor step size and biased direction resulting from the inaccurate surrogate estimation of the candidate solution.

Some of the commonly used surrogate models include, but are not limited to, Polynomial Regression (PR, Response surface), Gaussian Process (GP, Kriging), Neural Networks, and Support Vector Machine (SVM), and a comprehensive survey can be found by Jin [10] and Loshchilov [14]. Most recent works on surrogate model assisted ES considers sophisticated algorithms. These algorithms are heuristic in nature, and the step behaviours of the algorithm are not always well interpreted. In this context, a simple model for surrogate models can be helpful in understanding the surrogate behaviour, leading to potential modification to surrogate update or parameter-setting. A recent paper in

surrogate assisted EAs proposed by Kayhani and Arnold [11] analyzes surrogate assisted (1+1)-ES using a simple model for surrogate models on simple test functions where the surrogate estimate is modelled using a noisy estimate of the true objective function. The step size behaviour of the strategy on the test function is clearly interpreted. As a natural sequence, we investigate the surrogate assisted  $(\mu/\mu, \lambda)$ -ES using the same surrogate model and following a similar analysis. Since the  $(\mu/\mu, \lambda)$ -ES generates a population of candidate solutions where the surrogate model can be potentially more fully exploited compared with the (1+1)-ES, it is interesting how the model error would be affected and how much the ES is to benefit from the surrogate and the resulting step behaviour.

This thesis intends to analyze and understand the surrogate-assisted  $(\mu/\mu, \lambda)$ -ES on simple test functions following the analysis of surrogate model-assisted (1+1)-ES [11] and to exploit the potential benefit of using an extensive sampling with surrogate model assistance. The thesis is organized as follows: In Section 2 we give a brief review of the related background and previous analysis that is needed later, in Section 3 we analyze the proposed local surrogate model-assisted  $(\mu/\mu, \lambda)$ -ES and study its behaviour on sphere functions. Based on the result, in Section 4, we first apply the well established cumulative step size adaptation (CSA) to the algorithm and report the results. Given the experimental results, we propose an algorithm that is a cross between (1+1)-ES and  $(\mu/\mu, \lambda)$ -ES where the performance on several test functions are recorded followed by a discussion and future work in Section 5.

## 2 RELATED WORK

### 2.1 Evolution Strategies

Evolution strategies (ESs), a category of Evolutionary Algorithms (EAs), is a nature-inspired direct search method that addresses optimization problems by using stochastic variation and selection. In each iteration, new offspring are generated from the parental population through mutation, followed by a selection based on the fitness of the offspring. A subset of selected offspring is referred to as the parental population for the next iteration.

ESs are commonly used in black-box optimization where the search space  $\mathbb{R}^N$  is  $N$ -dimensional, whereas the objective function value is 1-dimensional (in  $\mathbb{R}$ ). We consider minimization of an objective function  $f : \mathbb{R}^N \rightarrow \mathbb{R}$  that maps the search space to the space for objective function values, i.e., maps a point (individual) in the search space to a value (its fitness) in the fitness space. It is worth noting that an individual with a larger fitness (larger value in the fitness space) has a smaller objective function value. There is no assumption on the objective function. Such optimization problems are referred to as black box optimization.

#### 2.1.1 $(\mu/\rho^+ \lambda)$ – ES .

In this Section, we use the formulation and description of the generalized ES from Hansen et al [16]. A single iteration of the general ES is shown in Alg. 1. Assume a parental population  $X$  with size  $\mu$ , the number of parents for recombination (in offspring generation)  $\rho$ , and the offspring generated in each iteration  $Y$  with size  $\lambda$ , where  $\mu, \rho, \lambda$  are positive integers with  $\rho \leq \mu$ .

In generation  $(g + 1)$ , we denote the parental population  $X^{(g)} = \{x_1^g, x_2^g, \dots, x_\mu^g\}$  where  $x_i^{(g)} \in \mathbb{R}^N$  for  $i = 1, \dots, \mu$  and  $j = 1, \dots, \lambda$ . In each inner iteration of generation  $(g)$ , the recombination is performed, meaning  $\rho$  individuals are randomly chosen from the parental population  $X^{(g)}$  and recombined for offspring generation described later. Here, we only introduce intermediate recombination that simply takes the arithmetic average of the  $\rho$  randomly selected individuals from  $X^{(g)}$  where the point obtained in inner iteration  $[i]$  after recombination is referred to as the centroid  $x_{\text{centroid}[i]}^{(g)}$ . The centroid  $x_{\text{centroid}[i]}^{(g)}$  obtained from intermediate recombination in inner iteration  $[i]$  of generation  $(g)$  is

---

**Algorithm 1** The  $(\mu/\rho^+\lambda) - ES$ 


---

```

1: Initialize  $N, \rho, \mu, \lambda \in N_+, \sigma \in \mathbb{R}, g \leftarrow 1$ 
2: Initialize parental population  $X^{(1)} \leftarrow \{x_i^{(1)} : i = 1, 2, \dots, \mu\}$ 
3: Evaluate  $X^{(1)}$  using objective function, yielding  $fX^{(1)} \leftarrow \{f(x_i^{(1)}) : i = 1, 2, \dots, \mu\}$ 
4: while not terminate() do
5:   for  $i = 1, 2, \dots, \lambda$  do
6:     Generate standard normally distributed  $z_i^{(g)} \in \mathbb{R}^N$ 
7:      $x_{\text{centroid}[i]}^{(g)} \leftarrow \text{recombine}(\text{select\_random}(\rho, X^{(g)}))$ 
8:      $y_i^{(g)} \leftarrow x_{\text{centroid}[i]}^{(g)} + \sigma z_i^{(g)}$ 
9:     Evaluate  $y_i^{(g)}$ , yielding  $f(y_i^{(g)})$ 
10:   end for
11:    $Y^{(g)} \leftarrow \{y_i^{(g)} : i = 1, 2, \dots, \lambda\}$ 
12:    $fY^{(g)} \leftarrow \{f(y_i^{(g)}) : i = 1, 2, \dots, \lambda\}$ 
13:   if comma-selection then
14:      $X^{(g+1)} \leftarrow \text{select\_best}(\mu, Y^{(g)}, fY^{(g)})$ 
15:   else if plus-selection then
16:      $X^{(g+1)} \leftarrow \text{select\_best}(\mu, X^{(g)} \cup Y^{(g)}, fX^{(g)} \cup fY^{(g)})$ 
17:   end if
18:   Update step size  $\sigma$ 
19:    $g \leftarrow g + 1$ 
20: end while

```

---

defined as

$$x_{\text{centroid}[i]}^{(g)} = \frac{1}{\rho} \sum_{i=1}^{\rho} x_i^{(g)}, x_i^{(g)} \in \text{select\_random}(\rho, X^{(g)}) \quad (1)$$

where  $\text{select\_random}(\rho, X^{(g+1)})$  randomly selects  $\rho$  individuals from  $X^{(g)}$  without replacement.

In offspring generation, a standard normally distributed mutation vector  $z_i^{(g)} \in \mathbb{R}^N$  is generated in each inner iteration  $[i]$  of generation  $(g)$  right after recombination. The mutation vector is added to the centroid with a step size parameter  $\sigma \in \mathbb{R}$  and we have

$$y_i^{(g)} = x_{\text{centroid}[i]}^{(g)} + \sigma z_i^{(g+1)}, 1 \leq i \leq \lambda \quad (2)$$

where  $z_i^{(g)}$  represents the mutation,  $x_{\text{centroid}[i]}^{(g+1)}$  is the centroid obtained after recombination in inner iteration  $[i]$  of generation  $(g)$ .

After the whole inner iteration in generation  $(g)$ , the offspring generated in generation  $(g)$  are denoted  $Y^{(g)} = \{y_1^g, y_2^g, \dots, y_\lambda^g\}$ , where  $y_j^{(g)} \in \mathbb{R}^N$  for  $j = 1, \dots, \lambda$ . A selection comes after that refers to how the parental population is updated; two selection techniques will be introduced, namely plus- or comma-selection ( $^+$ ). If a plus-selection is applied, only the best  $\mu$  individuals are chosen considering both the parental population and the offspring generated in this iteration (i.e., totally  $\mu + \lambda$  individuals are considered for selection). Whereas a comma-selection only chooses individuals from offspring population  $Y$  to update the parental population, no individual from past parental populations can be chosen (i.e., only  $\lambda$  individuals are considered for selection).  $X^{(g+1)}$  is defined as the parental population in

generation  $(g + 1)$  that follows

$$X^{(g+1)} = \begin{cases} \text{select\_best}(\mu, X^{(g)} \cup Y^{(g)}, fX^{(g)} \cup fY^{(g)}), & \text{plus-selection} \\ \text{select\_best}(\mu, Y^{(g)}, fY^{(g)}), & \text{comma-selection,} \end{cases} \quad (3)$$

where the  $(g)$  on the top right denotes generation,  $fX^{(g)}$  and  $fY^{(g)}$  are objective function values for each individual in population  $X^{(g)}$  and  $Y^{(g)}$ ,  $\text{select\_best}(\mu, X^{(g)}, fX^{(g)})$  selects the best  $\mu$  individuals from  $X^{(g)}$  according to their fitness recorded in  $fX^{(g)}$  (i.e.,  $\text{select\_best}(\mu, X^{(g)}, fX^{(g)}) = \{x_{i;\lambda}^{(g)} : 1 \leq i \leq \mu\}$  where  $f(x_{i;\lambda}^{(g)}) < f(x_{j;\lambda}^{(g)})$ ,  $1 \leq i < j \leq \lambda$ ).

Here we consider two special cases of the general ES, namely (1+1)-ES and  $(\mu/\mu, \lambda)$ -ES. The (1+1)-ES ( $\mu = \rho = \lambda = 1$  with plus-selection) generates a single offspring  $y = x + \sigma z$  in each generation, and the fitness of the offspring  $y$  is evaluated and compared to its parent  $x$ . The parent  $x$  is updated iff. the offspring is superior to its parent i.e.,  $f(y) < f(x)$ . Whereas the  $(\mu/\mu, \lambda)$ -ES ( $\mu = \rho$  with comma-selection) generates  $\lambda$  offspring with offspring population  $Y = \{y_i : y_i = \text{recombine}(X) + \sigma z_i, 1 \leq i \leq \lambda\}$ , the parental population  $X$  is updated by selecting the best  $\mu$  individuals from  $Y$  (i.e.,  $X = \text{select\_best}(\mu, Y, fY)$  where  $fY = \{f(y_i) : y_i \in Y\}$ ).

### 2.1.2 Step size adaptation.

*The 1/5th Success Rule.* The 1/5th success rule is a basic step size control for (1+1)-ES. The step size is adapted according to the success rate of generating a good offspring i.e., an offspring  $y$  with  $f(y) < f(x)$  in the case of (1+1)-ES  $x_{\text{centroid}} = x$ . If the success rate is lower than 1/5, the step size is decreased, otherwise increased. The 1/5 is chosen by Rechenberg [20] after obtaining the optimal success rate (i.e., achieving the largest fitness gain per iteration) for corridor function and quadratic sphere function to be  $\approx 0.184$  and  $\approx 0.270$  respectively for  $N \rightarrow \infty$ . The implementation of that rule suggested by Kern et al. can be found in [13]. **Why not include it here? not clear about this comment, what to include**

*Cumulative Step-Size Adaptation.* The step size of  $(\mu/\mu, \lambda)$ -ES is commonly adapted using cumulative step size adaptation (CSA) proposed by Ostermeier et al. [17]. For a strategy with ideally adapted step size, each step should be uncorrelated. If the consecutive steps are negatively correlated, the step size should be decreased. In contrast, if the consecutive steps are positively correlated, meaning the steps are pointing to the same direction. Then a number of small steps can be replaced by fewer large steps and therefore, the step size should increase.

To decide the correlation, information from previous steps and mutations are cumulated. By comparing the search path with the expected step length under random selection, the search path is adapted according to the expected length. Step size decreases if the length is less than expected and increases otherwise.

Define the search path as

$$p^{(g+1)} \leftarrow (1 - c)p^{(g)} + \sqrt{\mu c(2 - c)}z_{\text{step}}^{(g)}, \quad (4)$$

where  $0 < c \leq 1$  helps retain the history information (in generation  $(g)$ ) and pass that to the evolution path in the next generation  $(g + 1)$ ,  $\sqrt{\mu c(2 - c)}$  is a normalization constant denotes the proportion of the information obtained in generation  $(g)$  used to updates the evolution path and  $z_{\text{step}}^{(g)}$  is the direction vector obtained by averaging the direction

vectors from the best  $\mu$  individuals from selection. The  $z_{\text{step}}^{(g)}$  follows

$$z_{\text{step}}^{(g)} = \frac{1}{\mu} \sum_{i=1}^{\mu} z_{i;\lambda}^{(g)}, \text{comma-selection} \quad (5)$$

$$(6)$$

Note that the centroid in each generation does not change, so the inner iteration  $[i]$  is eliminated for simplicity and the centroid in generation  $(g)$  is simply  $x_{\text{centroid}}^{(g)}$ .

The step size is adapted

$$\sigma^{(g+1)} \leftarrow \sigma \exp \left( \frac{c}{d} \left( \frac{\|p^{(g+1)}\|}{E\|N(0, I)\|} - 1 \right) \right), \quad (7)$$

where under random selection and given  $p^{(0)} \sim N(0, I)$ , the expected length of the search path  $p^{(g+1)}$  can be approximated as  $E\|N(0, I)\| \approx \sqrt{N}(1 - 1/4N + 1/21N^2)$ . In Section 4, we use the well established parameters for CSA from Hansen [8] that follows

$$\begin{cases} c = (\mu + 2)/(N + \mu + 5) \\ d = 1 + 2 \max \left( 0, \sqrt{(\mu - 1)/(N + 1) - 1} \right) + c. \end{cases} \quad (8)$$

### 2.1.3 Analyzing ES.

To understand the behaviour of EAs, we first introduce analyzing ES on simple test functions where the step behaviours of the algorithm are more likely to be understood, and then proceed to the analysis on noisy sphere where the same analysis can be used to model the surrogate assisted  $(\mu/\mu.\lambda)$ -ES. Specifically, the  $(\mu/\mu.\lambda)$ -ES is first analyzed on quadratic sphere defined as  $f(x) = \sum_{i=1}^N x_i^2$  and then noisy sphere that models the ideal performance of surrogate model assisted  $(\mu/\mu.\lambda)$ -ES.

*On Sphere Function.* Decomposition of  $z$ , first proposed by Rechenberg [20] can be used to study the expected step of the strategy. Vector  $z$  can be decomposed as a vector sum  $z = z_1 + z_2$ , where  $z_1$  is in the direction of the negative gradient of the objective function  $\nabla f(x)$  with  $z_2$  orthogonal to  $z_1$ . We have  $z_1$  standard normally distributed,  $\|z_2\|^2 \chi^2$ -distributed with  $N - 1$  degree of freedom and  $\|z_2\|^2/N \xrightarrow{N \rightarrow \infty} 1$  (see [6]). Denote  $\delta^* = N(f(x) - f(y))/(2R^2)$ , where  $R = \|x\|$  is the euclidean distance to the optimal, we further introduce normalized step size  $\sigma^* = N\sigma/R$ . The normalized fitness gain of  $y$  over  $x$  given mutation  $z$  follows

$$\begin{aligned} \delta^*(z) &= \frac{N}{2R^2} (f(x) - f(y)) \\ &= \frac{N}{2R^2} (x^T x - (x + \sigma z)^T (x + \sigma z)) \\ &= \frac{N}{2R^2} (-2\sigma x^T z - \sigma^2 \|z\|^2) \\ &\stackrel{N \rightarrow \infty}{=} \sigma^* z_1 - \frac{(\sigma^*)^2}{2}, \end{aligned} \quad (9)$$

where  $z_1$  is the component of  $z$  pointing to the negative gradient of  $f(x)$  and  $\stackrel{N \rightarrow \infty}{=}$  denotes the convergence in distribution  $\|z\|^2/N = 1$  **fixed?**.

In the case of  $(\mu/\mu, \lambda)$ -ES, the progress vector  $z_{\text{step}} = 1/\mu \sum_{i=1}^{\mu} z_{i;\lambda}$  **fixed?** is the averaged  $z$  taken by the best  $\mu$  candidate solutions. The  $z_{i;\lambda,1}$  components of the selected  $\mu$  mutation vectors are correlated, while the  $z_{i;\lambda,2}$  components are uncorrelated. So that the length of  $z_{i;\lambda,2}, i = 1, 2, \dots, \mu$  component is reduced while the similarities in the selected  $\mu$  mutation vectors are persevered in  $z_{\text{step},1}$ . This is referred to as the genetic repair by Beyer [5].

So the normalized fitness gain given  $z_{\text{step}}$  is

$$\begin{aligned}\delta^*(z_{\text{step}}) &= \frac{N}{2R^2} (x^T x - (x + \sigma z_{\text{step}})^T (x + \sigma z_{\text{step}})) \\ &= \frac{N}{2R^2} (-2\sigma x^T z_{\text{step}} - \sigma^2 \|z_{\text{step}}\|^2) \\ &\stackrel{N \rightarrow \infty}{=} \sigma^* z_{\text{step},1} - \frac{(\sigma^*)^2}{2\mu},\end{aligned}\tag{10}$$

where  $z_{\text{step},1}$  is the component of  $z_{\text{step}}$  pointing to the negative gradient of  $f(x)$  and by assuming  $\mu \ll N$ ,  $\stackrel{N \rightarrow \infty}{=}$  denotes the convergence of the distribution  $\|z_{\text{step},1}\|^2/N = 1/\mu$  **fixed?**.

*On Noisy Sphere Function.* The sphere is considered noisy when the fitness evaluation is inaccurate and the objective function on a fixed point may vary within the noise strength in different objective function calls. **fixed?** The following uses the analysis and modelling proposed by Arnold and Beyer [2].

The objective function value on noisy sphere can be modelled by adding a Gaussian random variable with mean equals the true objective function value and some standard deviation referred to as noise strength  $\sigma_\epsilon$ . The noisy estimate of a candidate solution  $y$  follows  $f_\epsilon(y) = f(y) + \sigma_\epsilon z_\epsilon$  where  $z_\epsilon \in \mathbb{R}$  is a standard normally distributed random variable that randomize the noise generated. By further introduces  $\sigma_\epsilon^* = N\sigma_\epsilon/(2R^2)$ , the normalized fitness noise [3] and replace the accurate objective function evaluation with the noisy estimate, the normalized fitness gain of  $y$  on noisy sphere given mutation  $z$  when  $N \rightarrow \infty$  in Eq. (10) is

$$\begin{aligned}\delta_\epsilon^*(z) &= \frac{N}{2R^2} (f(x) - f_\epsilon(y)) \\ &= \frac{N}{2R^2} (x^T x - (x + \sigma z)^T (x + \sigma z) + \sigma_\epsilon z_\epsilon) \\ &= \frac{N}{2R^2} (-2\sigma x^T z - \sigma^2 \|z\|^2 - \sigma_\epsilon z_\epsilon) \\ &= \delta(z) + \sigma_\epsilon z_\epsilon \\ &\stackrel{N \rightarrow \infty}{=} \sigma^*(z_1 + \vartheta z_\epsilon) - \frac{(\sigma^*)^2}{2},\end{aligned}\tag{11}$$

where  $\vartheta = \sigma_\epsilon^*/\sigma^*$  is the noise-to-signal ratio, defined to measure the noise level relative to the algorithm's step size, the term  $+\sigma_\epsilon z_\epsilon$  denotes the added noise. Since  $z_\epsilon$  is standard normally distributed, so is  $-z_\epsilon$ , by substituting  $z_\epsilon = -z_\epsilon$  and adding the noise term to Eq. (10), we get the simplified Eq. (11).

For  $(\mu/\mu, \lambda)$ -ES, the normalized fitness gain in each generation is measured between two consecutive parents over a progress vector  $z_{\text{step}}$  where the fitness of the two are evaluated using the true objective function. We want to compute on average, how much benefit one objective function call can bring in each generation, therefore the fitness gain obtained over progress vector  $z_{\text{step}}$  is divided by the number of objective function calls made in this generation. The

normalized fitness gain when dimensionality  $N \rightarrow \infty$  follows

$$\eta = \frac{1}{\lambda} E[\delta^*(z_{\text{step}})] \approx \frac{1}{\lambda} E \left[ \sigma^* z_{\text{step},1} - \frac{(\sigma^*)^2}{2\mu} \right] \quad (12)$$

where the expected value of the normalized fitness gain over the progress vector is divided by  $\lambda$ , the number of offspring evaluated in each iteration i.e., the the objective function evaluation per iteration for  $(\mu/\mu, \lambda)$ -ES is  $\lambda$ .

The expected value of  $z_{\text{step},1}$  derived by Arnold [2] is

$$\begin{aligned} E[z_{\text{step},1}] &= E \left[ \frac{1}{\mu} \sum_{i=1}^{\mu} z_{(i;\lambda),1} \right] \\ &= \frac{1}{\mu} \sum_{i=1}^{\mu} E[z_{(i;\lambda),1}] \\ &= \frac{1}{\mu} \sum_{i=1}^{\mu} \int_{-\infty}^{\infty} x p_{i;\lambda}(x) dx \\ &= \frac{c_{\mu/\mu,\lambda}}{\sqrt{1+\vartheta^2}} \end{aligned} \quad (13)$$

where  $z_{(i;\lambda),1}$  is the component of  $z_{(i;\lambda)}$  in the negative gradient direction of the objective function  $f(x)$ ,  $p_{i;\lambda}$  corresponds to the probability density function of  $z_{(i;\lambda),1}$  that the generated offspring using this direction has the  $i^{\text{th}}$  largest fitness (i.e.,  $i^{\text{th}}$  smallest objective function value) **fixed?**,  $c_{\mu/\mu,\lambda}$  is the  $(\mu/\mu, \lambda)$ -progress coefficient derived by Arnold and Beyer [1] that follows

$$c_{\mu/\mu,\lambda} = \frac{\lambda - \mu}{2\pi} \left( \frac{\lambda}{\mu} \right) \int_{-\infty}^{\infty} e^{-x^2} [\Phi(x)]^{\lambda-\mu-1} [1 - \Phi(x)]^{\mu-1} dx, \quad (14)$$

where  $\Phi$  is the normal cumulative distribution function. The integral can be solved numerically.

Therefore by Eq. (12) (13), when dimensionality  $N \rightarrow \infty$  the normalized fitness gain can be simplified as

$$\eta = \frac{1}{\lambda} E[\delta^*(z_{\text{step}})] \approx \frac{1}{\lambda} \left( \frac{\sigma^* c_{\mu/\mu,\lambda}}{\sqrt{1+\vartheta^2}} - \frac{(\sigma^*)^2}{2\mu} \right) \quad (15)$$

## 2.2 Surrogate Model

Surrogate models are computational models constructed based on the data evaluated using true objective function. The surrogate acts as an approximation to the true objective function that is costly in most cases. We make the assumption that the objective function value estimated using the surrogate model, although inaccurate, can be achieved at vanishing cost.

The surrogate model can be applied to EAs as an approximate fitness to accelerate the evolution process [19]. Despite the computational saving when using a surrogate model, issues can occur when the surrogate built leads to a false optimum (the optimum does not exist in the true objective function), leading to potential divergence and unstable optimization path where the convergence property of the ES may not be well preserved [10]. Two approaches will be discussed that give potential solutions to address this issue by ensuring model accuracy. The first approach uses an EA to optimize a surrogate model, and the second uses a surrogate model to assist EA.

### 2.2.1 Gaussian Process.

A GP is a probabilistic model where the observations are in a continuous domain  $\mathbb{R}^N$ . It is completely determined by its mean function

$$m(x) = E[f(x)], \quad (16)$$

commonly assumed to be zero, and a covariance function (a positive definite kernel) follows

$$\kappa(x, x') = E[(f(x) - m(x))(f(x') - m(x'))]. \quad (17)$$

$\forall x_i \in \mathbb{R}^N, i = 1, 2, \dots, n$ , the distribution of the function value  $f(x_i), i = 1, 2, \dots, n$  is jointly Gaussian with mean  $\mu = (m(x_1), m(x_2), \dots, m(x_n))$ , and covariance matrix  $\Sigma$  with entry  $\Sigma_{ij} = \kappa(x_i, x_j)$ .

Squared exponential kernel, a commonly used covariance function is applied in the context

$$\kappa(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\theta^2}\right), \quad (18)$$

where  $\theta$  is the length scale factor of the GP.

The GP is defined using the notation and content from [18] and [15]: let  $f(x)$  be an unknown scalar function and  $x \in \mathbb{R}^N$  is a point in an  $N$ -dimensional space. Evaluating  $f$  at  $n$  data points  $X = (x_1, x_2, \dots, x_n)$  yields function values  $f = (f(x_1), f(x_2), \dots, f(x_n))$ . We want to predict new function values  $f(X_*)$  of a test set  $X_*$  with size  $n_*$ .

The vector of known function values and the predicted value  $(f, f_*)$  are jointly normally distributed with mean  $(\mu, \mu_*)$  and covariance matrix

$$\begin{bmatrix} K & K_* \\ K_*^T & K_{**} \end{bmatrix}, \quad (19)$$

where  $K = \kappa(X, X)$ ,  $K_* = \kappa(X, X_*)$ , and  $K_{**} = \kappa(X_*, X_*)$ .

By Bayes' rule,  $f_*$  is normally distributed with mean and covariance

$$\mu_* = K_*^T K^{-1} f \quad (20)$$

$$\Sigma_* = K_{**} - K_*^T K^{-1} K_* \quad (21)$$

### 2.2.2 Surrogate model assisted ES.

*Surrogate model assisted by ES.* Gaussian Process Optimization Procedure (GPOP) proposed by Buche et al. [7] uses an EA optimized surrogate model where a fraction of individuals in each generation is evaluated depending on the surrogate model error.

Covariance Matrix Adaptation Evolution Strategy (CMA-ES) [8] is an ES that is effective in handling ill-conditioned problems. The CMA-ES is effective in handling these problems because a covariance matrix, the inverse of the Hessian is estimated within an iterative procedure. The covariance matrix is used to generate mutations of the CMA-ES and is updated in each generation.

A GP is constructed using a set of evaluated points from random sampling in previous iterations, a CMA-ES is used for parameter searching in order to find the minimum of the GP prediction. To avoid false optima, a merit function is used to help explore new regions of the decision space, regarded as the fitness function for the CMA-ES. Predicted



standard deviation is used as the merit function that follows

$$f_M(x) = \hat{f}(x) - \alpha \sigma_{GP}(x), \quad (22)$$

where  $\hat{f}(x)$  and  $\sigma_{GP}(x) = \sqrt{\Sigma_*}$  (considering a single input  $x'$  in Eq. (21)) are the GP estimate and standard deviation for data  $x$ ,  $\alpha \geq 0$  balances the two terms via scaling the density measure (the density of the points sampled within a region) where a larger  $\alpha$  pushes the search harder into unexplored region. Four merit functions with  $\alpha = 0, 1, 2, 4$  are used and optimized in the context. Finally, the resulting minimum is added to the training set.

To approximate the true objective function with arbitrary precision, a local model is preferred because it uses a smaller number of training points and provides a more precise local approximation of the objective function, compared with a global model. A local model approximates the objective function within a limited region, as opposed to a global approximation that approximates the whole objective function using all evaluated points. In the local model, the new points sampled moves around the current best solution and is restricted to a neighbourhood of the evaluated current best point for model accuracy, which means only points within the region are considered reliable using surrogate estimates. In each iteration,  $N_C$  points are sampled. If more than half of the sampled points are not evaluated successfully (falls out the well-approximated region using GP), another  $N_C/2$  points are generated using the (2, 10) – CMAES and evaluated using the true objective function followed by a model update that chooses the best evaluated point  $x_{\text{best}}$  and re-builds the training set using the  $N_C$  closest points to  $x_{\text{best}}$  and  $N_R$  most recent successfully evaluated points. The GP parameters, specifically the length scale, offset of the function value predicted from 0, and a scaling term for added white noise are then optimized by ES, each predicted optimum is found using the corresponding merit function where the unevaluated optimum is evaluated using the true objective function and the process repeats.

Experimental results show GPOP can be effective in solving unimodal functions and can achieve an average speed-up (number of objective function evaluations to solve the optimization problem within required precision using CMA divided by that of using GPOP) of 3 to 6 for quadratic sphere, Schwefel problem and Rosenbrock’s function. It shows that the GPOP cannot determine the global minima for  $N > 2$  given the large number of local minima and the high oscillation of the function value between minima. The paper does not indicate a generalized training set size for GPOP where the training set sizes used are problem-specific. The main problem with GPOP, compared with CMA, is its higher computational cost of the optimization procedure that scales  $O(N^3)$  and  $O(N)$  with training size and problem dimension respectively, therefore GPOP should be considered mainly for computational expensive problems.

*ES assisted by surrogate model.* The other approach, Local Metamodel Covariance Matrix Adaptation ES (Imm-CMA) proposed by Kern et al. [12] thoroughly evaluates the surrogate model (metamodel) to find a relative precise optimum predicted by the model where the model is updated whenever the ranking of the best  $\mu$  solutions is inconsistent in two consecutive metamodel iterations. The surrogate model in this approach is more thoroughly evaluated compared to the GPOP because the ES proceeds to the next iteration iff. the GP is regarded accurate or otherwise the GP is continuously updated; whereas the GP error in GPOP only indicates the number of the points to be evaluated in the next iteration.

An approximate ranking procedure can be applied to ensure the model accuracy without knowing the true ranking of the complete population. The proposed procedure acts as a ranking procedure in CMA-ES that ensures an effective ranking by evaluating a batch of individuals in each metamodel iteration where the model is updated once the model is regarded inaccurate. The CMA-ES used can be viewed as a  $(\mu/\mu, \lambda)$ -ES with a covariance matrix see Section 2.1.1. In ES generation  $g$ ,  $\lambda$  offspring are generated where the best  $\mu$  offspring  $y_{1:\mu, \lambda}$  are selected according to the predicted fitness by the model  $\hat{f}$ . The ranking of the  $\mu$  selected offspring is recorded as  $\text{ranking}_{1-\mu; \lambda}^{(i)}$  in model iteration  $i$  and

the best  $n_{\text{init}}$  selected individuals are evaluated using the true objective function for training set update. The selected  $\lambda$  individuals  $y_{1:\mu,\lambda}$  are evaluated again by the model in the next model iteration  $i + 1$ , leading to a new ranking  $\text{ranking}_{1-\mu;\lambda}^{(i+1)}$  followed by a comparison. If  $\text{ranking}_{1-\mu;\lambda}^{(i)} = \text{ranking}_{1-\mu;\lambda}^{(i+1)}$ , the model is regarded as reliable, otherwise the next  $n_b$  best unevaluated points based on  $\hat{f}$  are evaluated (by the true objective function) and added to training set, this iterates until all  $\lambda$  offspring in this generation has been evaluated or the number of model iteration exceeds a certain number.

The performance is evaluated using a local weighted regression (LWR) [4] as the metamodel for Imm-CMA. Imm-CMA achieves a speed-up (number of objective function evaluates for the other strategy divided by Imm-CMA) of 5-8 and 2-3 compared with standard CMA-ES [9] on convex quadratic Schwefel function and Rosenbrock function respectively. The performances of Imm-CMA for the above two convex quadratic function also matches the GPOP described above for dimension  $N \leq 4$  and is even better for larger dimension  $N \geq 8$  given a potential less reliable Gaussian Process Regression. On noisy sphere, Imm-CMA gains a small advantage over CMA-ES but the advantage vanishes with increasing dimensionality. Despite the speed-up achieved compared with CMA-ES, constructing the metamodel for Imm-CMA takes computational complexity up to  $N^6$ , so that this algorithm should also be considered for computational expensive problems.

### 3 ANALYSIS

To understand the potential implications of using surrogate model assisted  $(\mu/\mu, \lambda)$ -ES with varying population sizes, in this section, we use the same simple model proposed by Kayhani and Arnold [11] for the use of a surrogate model. Specifically, we propose an EA that generates a population of  $\lambda$  offspring in each generation where the offspring are evaluated using the surrogate model instead of the true objective function. An intermediate recombination is performed based on the inaccurate fitness of the offspring estimated by the surrogate model; the resulting centroid  $x_{\text{centroid}}$  of the  $\mu$  selected offspring is evaluated using the true objective function. The analysis for  $(\mu/\mu, \lambda)$ -ES using simple model is the same as that of the  $(\mu/\mu, \lambda)$ -ES on noisy sphere — in offspring ranking, the former uses a surrogate model to estimate offspring's fitness (achieved at vanishing cost, by the assumption in Section 2.2), and the latter uses a noisy estimate of offspring's fitness where in both cases, the fitness estimation are achieved the same way. The normalized fitness gain of two consecutive centroids  $x_{\text{centroid}}^{(g)}$  and  $x_{\text{centroid}}^{(g+1)}$  over the progress vector  $z_{\text{step}}^{(g)}$  is

$$\begin{aligned}\delta_{GP}^* &= \frac{N}{2R^2} \left( f(x_{\text{centroid}}^{(g+1)}) - f(x_{\text{centroid}}^{(g)}) \right) \\ &= \delta^*(z_{\text{step}}^{(g+1)}) \\ &\stackrel{N \rightarrow \infty}{=} \sigma^* z_{\text{step},1}^{(g+1)} - \frac{(\sigma^*)^2}{2\mu}.\end{aligned}\tag{23}$$

The analysis in Section 2.1.3 still holds which gives the same result as Eq. (11). In this context, the noise-to-signal ratio  $\vartheta$  can be interpreted as the measure of the surrogate model quality relative to the step size of the algorithm.

Since the fitness of  $\lambda$  offspring generated are evaluated by the surrogate model at vanishing cost by earlier assumption. The number of objective function evaluation per iterations is one instead of  $\lambda$  (for  $(\mu/\mu, \lambda)$ -ES without model assistance), therefore the normalized fitness gain when dimensionality  $N \rightarrow \infty$ , by substituting  $\lambda$  with 1 in Eq. (11) is

$$\eta = \frac{1}{1} E[\delta_{GP}^*] \approx \frac{\sigma^* c_{\mu/\mu,\lambda}}{\sqrt{1 + \vartheta^2}} - \frac{(\sigma^*)^2}{2\mu},\tag{24}$$

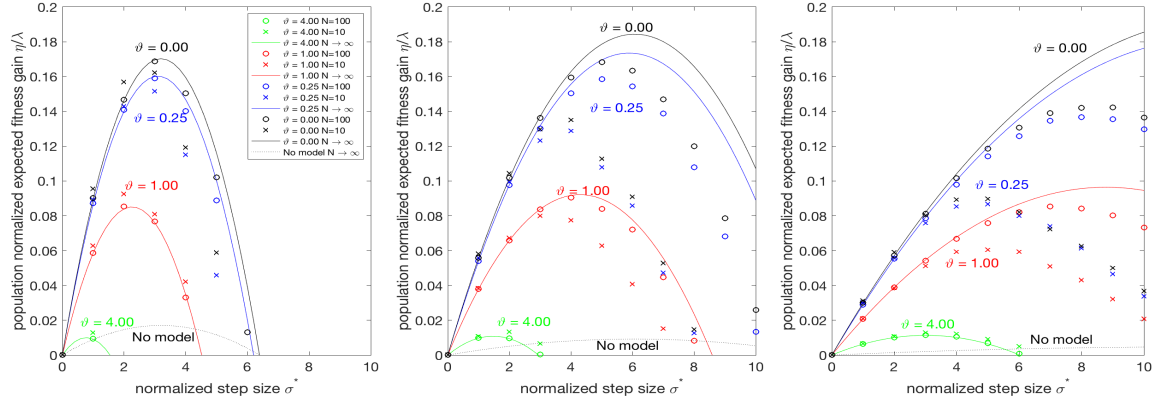


Fig. 1. The figures from left to right show the expected single step behaviour of the surrogate model assisted  $(\mu/\mu, \lambda)$ -ES with unbiased Gaussian distributed surrogate error for  $\lambda = 10, 20, 40$  ( $\mu = \lceil \lambda/4 \rceil$ ) respectively. The solid lines are the results obtained analytically when  $N \rightarrow \infty$ ; the dotted lines in the bottom of the figures show the relationship for corresponding  $(\mu/\mu, \lambda)$ -ES without surrogate model assistance (when  $N \rightarrow \infty$ ); the dots represent the experimental results for  $N = 10$  (crosses) and  $N = 100$  (circles).

To obtain the opt. expected fitness gain  $\eta_{opt}$  and its corresponding opt. normalized step size  $\sigma_{opt}^*$  over a fixed noise-to-signal ratio  $\vartheta$ , by assuming independence of  $\vartheta$  and  $\sigma^*$ , we take the derivative of equation (24) over  $\sigma^*$  and obtain

$$\sigma_{opt}^* = \frac{\mu c_{\mu/\mu, \lambda}}{\sqrt{1 + \vartheta^2}} \quad (25)$$

$$\eta_{opt} = \frac{\sigma_{opt}^* c_{\mu/\mu, \lambda}}{\sqrt{1 + \vartheta^2}} - \frac{(\sigma_{opt}^*)^2}{2\mu} \quad (26)$$

For easy comparison, the expected fitness gain is normalized regarding the population size  $\lambda$ ; the normalized fitness gains for all strategies are divided by their corresponding population size  $\lambda$  (i.e.,  $\eta/\lambda$  is plotted against  $\sigma^*$ ) in Fig. 1. The population-normalized fitness gains against the normalized step size for the  $(\mu/\mu, \lambda)$ -ES with population size  $\lambda = 10, 20, 40$  and corresponding  $\mu = 3, 5, 10$  are plotted from left to right in Fig. 1. The lines show the population-normalized results obtained (after normalization regarding population size  $\lambda$ ) from Eqs. (14), (15), (24), and Eqs. (4), (6),  $\eta = E[\Delta]/p_{eval}$  from the surrogate model assisted (1+1)-ES [11]. The dots represent the experimental results of unbiased Gaussian surrogate error for  $N \in \{10, 100\}$  obtained by averaging 100 runs.

It can be inferred from Fig. 1 that, for a fixed population size, the expected fitness gain decreases as the noise-to-signal-ratio  $\vartheta$  increases. When  $\vartheta \rightarrow \infty$ , the surrogate model becomes useless, and the strategy becomes a random search. For moderate noise-to-signal ratio  $\vartheta$ , the surrogate model assisted  $(\mu/\mu, \lambda)$ -ES can achieve much larger opt. expected fitness gain at a larger normalized step size compared with the surrogate model assisted (1+1)-ES [11]. When  $\vartheta = 1$ , the maximal expected fitness gain achievable for (3/3, 10)-ES, (5/5, 20)-ES, and (10/10, 40)-ES are 0.8426, 1.841, and 3.856 with corresponding  $\sigma^* = 2.258, 4.294$ , and 8.769 respectively; for surrogate model assisted (1+1)-ES, the maximal fitness gain is 0.548 achieved at  $\sigma^* = 1.905$ . As for  $(\mu/\mu, \lambda)$ -ES and (1+1)-ES both without model assistance, the maximal fitness gains are 0.202, 0.170, 0.184, and 0.193 achieved at  $\sigma^* = 1.224, 3.200, 6.080$ , and 12.420 for (1+1)-ES, (3/3, 10)-ES, (5/5, 20)-ES, and (10/10, 40)-ES respectively. The value of maximal fitness achievable for  $(\mu/\mu, \lambda)$ -ES without surrogate

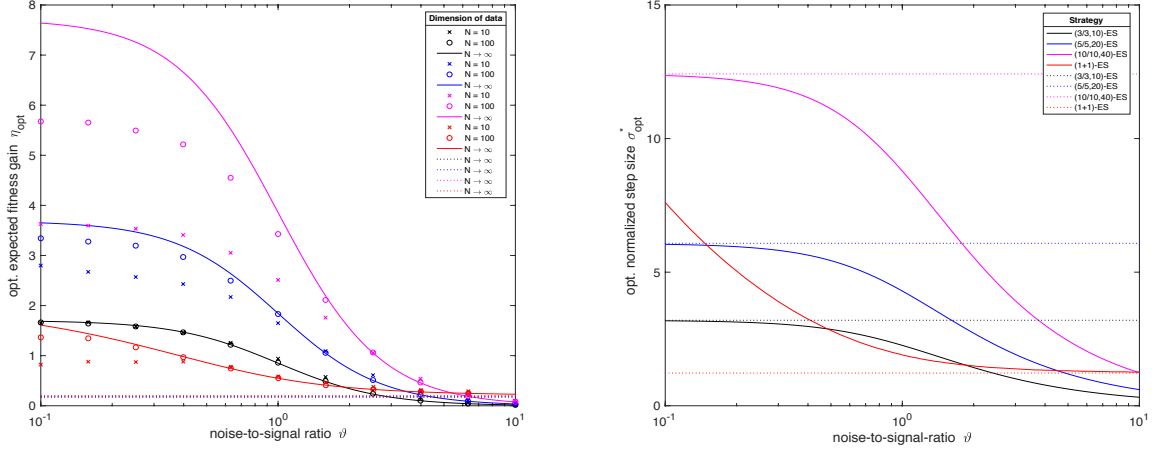


Fig. 2. Opt. expected fitness gain and corresponding opt. normalized step size of the surrogate model assisted  $(\mu/\mu, \lambda)$ -ES and  $(1+1)$ -ES plotted against the noise-to-signal ratio. Colour black, blue, magenta and red represent the result for  $(3/3, 10)$ -ES,  $(5/5, 20)$ -ES,  $(10/10, 40)$ -ES and  $(1+1)$ -ES respectively. The solid line represents the results obtained analytically when  $N \rightarrow \infty$ . The dots are the experimental result for  $N = 10$  (crosses) and  $N = 100$ . The dotted lines show the optimal values for the  $(3/3, 10)$ -ES,  $(5/5, 20)$ -ES,  $(10/10, 40)$ -ES and  $(1+1)$ -ES without surrogate model assistance.

model assistance grows as  $\lambda$  increases and gradually approaches 0.202, which asymptotically equals to that of the  $(1+1)$ -ES [5]. From the above analysis, surrogate assisted  $(\mu/\mu, \lambda)$ -ES does benefit from using a larger population, namely, the expected maximal fitness gain increased by a factor of  $\lambda$  after the surrogate model is applied; and an observable improvement in opt. fitness gain achieved at a larger normalized step size compared with the surrogate model assisted  $(1+1)$ -ES. For  $\vartheta = 0$  (the surrogate model models the objective function exactly), from Eqn. (25) we can obtain the maximal expected fitness gain for surrogate assisted  $(\mu/\mu, \lambda)$ -ES achieved at  $\sigma_{opt}^* = \mu c_{\mu/\mu, \lambda}$  with value  $\eta_{opt} = \mu(c_{\mu/\mu, \lambda})^2/2$ . Even if this indicates the potential benefit of using a growing population size, it is important to note the analytical result derived when  $N \rightarrow \infty$  is an approximation for the finite-dimensional case. Fig. 2 shows the relation of optimal expected fitness gain and the corresponding optimal normalized step size over noise-to-signal ratio derived analytically in the limit of  $N \rightarrow \infty$  for  $(\mu/\mu, \lambda)$ -ES with  $\lambda = 10, 20, 40$  and  $(1+1)$ -ES, all including cases with and without surrogate model assistance. The optimal expected fitness gain is also measured experimentally for  $n \in \{10, 100\}$ .

The speed-up is defined as the median number of objective function evaluations used to solve the test problems (when reaching the termination criteria) using one strategy divided by that of using the surrogate model assisted ES, both using the opt. normalized step size. For comparison purposes, we define two speed-ups, namely  $\text{speed-up}_{\text{self}}$  and  $\text{speed-up}_{\text{model}}$ .  $\text{speed-up}_{\text{self}}$  is defined as the number of objective function calls used for a surrogate model assisted ES divided by that of the ES without surrogate model assistance (specifically  $(\mu/\mu, \lambda)$ -ES if the ES is not specified); and  $\text{speed-up}_{\text{model}}$  as the number of objective function calls needed for  $(1+1)$ -ES to solve the problem divided by that of  $(\mu/\mu, \lambda)$ -ES, both with surrogate model assistance. The results obtained for noise-to-signal  $\vartheta = 0.1$  and 1 are shown in Table 1 and Table 2.

Table 1. Speed-ups for a small noise-to-signal ratio ( $\vartheta = 0.1$ ) observed in analysis

speed-up	Different Evolution Strategy used for comparison			
	(1+1)-ES	(3/3, 10)-ES	(5/5, 20)-ES	(10/10, 40)-ES
speed-up <sub>self</sub> ( $N = 10$ )	4.1	9.9	15.2	18.8
speed-up <sub>self</sub> ( $N = 100$ )	6.8	9.9	18.1	29.4
speed-up <sub>self</sub> ( $N \rightarrow \infty$ )	8.0	10.0	20.0	40.0
speed-up <sub>model</sub> ( $N = 10$ )		1.6	3.0	6.2
speed-up <sub>model</sub> ( $N = 100$ )		1.7	3.3	6.6

Table 2. Speed-ups for moderate noise-to-signal ratio ( $\vartheta = 1$ ) observed in analysis

speed-up	Different Evolution Strategy used for comparison			
	(1+1)-ES	(3/3, 10)-ES	(5/5, 20)-ES	(10/10, 40)-ES
speed-up <sub>self</sub> ( $N = 10$ )	2.7	5.1	9.0	13.0
speed-up <sub>self</sub> ( $N = 100$ )	2.7	5.5	9.9	17.8
speed-up <sub>self</sub> ( $N \rightarrow \infty$ )	2.9	5.9	10.0	19.9
speed-up <sub>model</sub> ( $N = 10$ )		2.4	3.7	4.6
speed-up <sub>model</sub> ( $N = 100$ )		1.5	2.7	4.4

For a finite-dimension e.g.  $N = 10$ , the speed-up<sub>self</sub> achieved for small noise-to-signal ratio ( $\vartheta = 0.1$ ) appears to be around nine and ten for a population size  $\lambda = 10$ , fifteen and sixteen for  $\lambda = 20$ , eighteen and nineteen for  $\lambda = 40$ ; whereas for (1+1)-ES, the speed-up<sub>self</sub> tops out between four and five. When  $N = 100$ , the speed-up<sub>self</sub> for small noise-to-signal ratio are between nine and ten, eighteen and nineteen, and twenty-nine and thirty for (3/3, 10)-ES, (5/5, 20)-ES, and (5/5, 40)-ES respectively; as for (1+1)-ES, the speed-up<sub>self</sub> with a small noise-to-signal ratio is between six and seven when  $N = 100$ . The speed-up<sub>self</sub> (for  $(\mu/\mu, \lambda)$ -ES) mounts with a growing population size  $\lambda$ , but it is worth noting that speed-up<sub>self</sub>  $\xrightarrow{N \rightarrow \infty} \lambda$ , and the gap in speed-up<sub>self</sub> between a finite dimension (fixed) and infinite dimension increases as  $\lambda$  grows (e.g., compared with  $N \rightarrow \infty$ , the speed-up<sub>self</sub> achieved when  $N = 10$  decreases 1%, 24%, and 53% for  $\lambda = 10, 20$ , and 40 respectively).

When comparing (1+1)-ES and  $(\mu/\mu, \lambda)$ -ES both with surrogate model assistance, speed-up<sub>model</sub> for (3/3, 10)-ES, (5/5, 20)-ES, and (5/5, 40)-ES when  $N = 10$  are between two and three, three and four, four and five respectively; whereas  $N = 100$ , the speed-up<sub>model</sub> decrease to a value between one and two, two and three, and four to five respective for  $\lambda = 10, 20$ , and 40. From the observed result, the speed-up for surrogate model assisted  $(\mu/\mu, \lambda)$ -ES is significant over both the surrogate model assisted (1+1)-ES and  $(\mu/\mu, \lambda)$ -ES without model assistance. It seems that the expected fitness gain for the surrogate assisted  $(\mu/\mu, \lambda)$ -ES will increase in line with population size  $\lambda$ , especially with a large dimensionality. Besides, the shrinking gap in speed-up<sub>model</sub> between different finite dimension over a growing population (i.e., the speed-up<sub>model</sub> for  $N = 100$  approaches that of  $N = 10$  as  $\lambda$  grows) also suggests using a larger  $\lambda$  for potential stable speed-up<sub>model</sub> compared with the surrogate assisted (1+1)-ES. Although the analytical results obtained for  $N \rightarrow \infty$  are approximations in the finite-dimension and can be inaccurate but it gives an implication for using larger  $\lambda$  for potential larger fitness gain.

## 4 STEP SIZE ADAPTATION

### 4.1 Cumulative step size adaptation

---

**Algorithm 2** Surrogate Model Assisted  $(\mu/\mu, \lambda)$ -ES (GP- $(\mu/\mu, \lambda)$ -ES)

---

```

1:  $c \leftarrow \frac{\mu+2}{N+\mu+5}$ 
2:  $d \leftarrow 1 + 2\max(0, \sqrt{\frac{\mu-1}{N+1}} - 1) + c$ 
3:  $E\|\mathcal{N}(0, I)\| \approx \sqrt{N}(1 - \frac{1}{4N} + \frac{1}{21N^2})$ 
4:  $s^{(0)} \leftarrow 0$ 
5: while not terminate() do
6:   for  $i = 1, 2, \dots, \lambda$  do
7:     Generate standard normally distributed  $z_i^{(g)} \in \mathbb{R}^N$ 
8:     Evaluate  $x_{\text{centroid}}^{(g)} + \sigma^{(g)}z_i^{(g)}$  using the surrogate model, yielding  $f_\epsilon(x_{\text{centroid}}^{(g)} + \sigma^{(g)}z_i)$ 
9:   end for
10:   $z_{\text{step}}^{(g)} \leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} z_{i;\lambda}^{(g)}$ 
11:   $x_{\text{centroid}}^{(g+1)} \leftarrow x_{\text{centroid}}^{(g)} + \sigma^{(g)}z_{\text{step}}^{(g)}$ 
12:  Evaluate  $x_{\text{centroid}}^{(g+1)}$  using true objective function, yielding  $f(x_{\text{centroid}}^{(g+1)})$ 
13:  Update surrogate model by adding  $(x_{\text{centroid}}^{(g+1)}, f(x_{\text{centroid}}^{(g+1)}))$ 
14:   $s^{(g)} \leftarrow (1-c)s + \sqrt{c(2-c)}\mu z_{\text{step}}^{(g)}$ 
15:   $\sigma^{(g+1)} \leftarrow \sigma^{(g)} \exp\left(\frac{c}{d} \frac{\|s^{(g)}\|}{E\|\mathcal{N}(0, I)\|} - 1\right)$ 
16:   $g \leftarrow g + 1$ 
17: end while

```

---

Even though the analysis in Section 3 suggests a potential better performance for the surrogate-assisted  $(\mu/\mu, \lambda)$ -ES. There is no guarantee that the step size of the strategy can be properly adapted, and further the analysis is very inaccurate in terms of finite dimension. In this section we experiment the surrogate model assisted  $(\mu/\mu, \lambda)$ -ES using the CSA described in Section 2.1.1 and 2.1.2, and exploit the potential insight it may offer. The strategy is evaluated using a Gaussian Process surrogate model in place of the simple model described in Section 3; one single iteration of the surrogate model assisted  $(\mu/\mu, \lambda)$ -ES using CSA (referred to as GP- $(\mu/\mu, \lambda)$ -ES) is shown in Alg. 2 where  $z_{i;\lambda}^{(g)}$  denotes the  $i$ th best mutation ranked according to the fitness of the corresponding offspring estimated by the surrogate model i.e.,  $x_{\text{centroid}}^{(g)} + \sigma^{(g)}z_{i;\lambda}^{(g)}$  is the offspring that has  $i$ th largest fitness ( $i$ th smallest estimated objective function value) according to the surrogate model **fixed?**.

Five ten-dimensional test problems are used to test if the step size of the strategy has been appropriately adapted, namely sphere functions  $f(x) = (x^T x)^{\alpha/2}$  for  $\alpha = \{1, 2, 3\}$ , referred to as linear, quadratic, and cubic spheres;  $f(x) = \sum_{i=1}^N (\sum_{j=1}^i x_j)^2$  (i.e., a convex quadratic function with condition number of the Hessian approximately equal to 175.1) referred to as Schwefel's Problem 1.2 [21]; and quartic function [11] defined as  $f(x) = \sum_{i=1}^{N-1} [\beta(x_{i+1} - x_i^2)^2 + (1 - x_i)^2]$  where  $\beta = 1$ . For  $\beta = 100$ , the quartic function becomes the Rosenbrock function with the condition number of the Hessian at the optimizer exceeds 3,500, making it very hard to find the global optimum without adapting the shape of the mutation distribution. So we use the quartic function with  $\beta = 1$  in the context where the corresponding condition number of the Hessian at the optimizer equals to 49.0. The values of global optima for all test function are zero. For

Table 3. Median test results for  $(\mu/\mu, \lambda)$ -ES without surrogate model assistance.

Test functions	Median number of objective function calls		
	(3/3, 10)-ES	(5/5, 20)-ES	(10/10, 40)-ES
linear sphere	3300	4809	8405
quadratic sphere	1694	2436	4182
cubic sphere	1166	1659	2788
Schwefel' s function	6259	8064	13325
quartic function	6600	8442	14637

Table 4. Median test results and speed-ups of GP- $(\mu/\mu, \lambda)$ -ES.

Test functions	Median number of objective function calls (speed-up <sub>self</sub> , speed-up <sub>model</sub> )			
	(1 + 1)-ES	(3/3, 10)-ES	(5/5, 20)-ES	(10/10, 40)-ES
linear sphere	502	756(4.4, 0.66)	696(6.9, 0.72)	761(11.0, 0.66)
quadratic sphere	214	309(5.5, 0.69)	245(9.9, 0.87)	231(18.1, 0.93)
cubic sphere	205	273(4.3, 0.75)	249(6.7, 0.82)	253(6.7, 0.81)
Schwefel' s function	1503	2278(2.7, 0.66)	$+\infty(/)$	$+\infty(/)$
quartic function	1265	1000(6.6, 1.3)	746(11.3, 1.7)	667(21.9, 1.9)

Note: the value shown for (1+1)-ES is its corresponding speed-up<sub>self</sub>.

each test problem, 100 runs are conducted for (1+1)-ES, (3/3,10)-ES, (5/5,20)-ES, (10/10,40)-ES both with and without surrogate model assistance.

For surrogate model, we use Gaussian process with squared exponential kernel and the length scale parameter in the kernel is set proportional to both the step size of the ES and the square root of data dimension  $N$ . For simplicity, the length scale is set to  $8\sigma\sqrt{N}$  as is used in the surrogate assisted (1+1)-ES [11]. In our experience, using a smaller or larger value by a factor of two i.e., 2, 4 and 16, 32 does not give a significant improvement, so we stick to 8 in the context.

The Gaussian process kernel is constructed using a training size of 40. The training set consists of the 40 most recent candidate solutions evaluated, so that the surrogate model approximates the local landscape of the objective function. All runs are initialized with starting point sampled from a Gaussian distribution with zero mean and unit covariance matrix and initial step size  $\sigma_0 = 1$ . The termination criteria is defined as one solution achieves objective function value below  $10^{-8}$ .

Histograms showing the number of objective function calls needed to solve the test problems within the required accuracy are represented in the first row of Fig. 3, the median number of objective function calls required to solve each test problem for  $(\mu/\mu, \lambda)$ -ES without and with surrogate model assistance (the step size of both are adapted using CSA) is shown in Table 3 and Table 4 respectively; the result for surrogate assisted (1+1)-ES [11] is also included in Table 4 for comparison.

The two speed-ups speed-up<sub>self</sub> and speed-up<sub>model</sub> defined in Section 3 are used. The speed-up<sub>self</sub> on quadratic sphere for GP- $(\mu/\mu, \lambda)$ -ES with  $\lambda = 10, 20, 40$  ( $\mu = \lceil \lambda/4 \rceil$ ) are between five and six, nine and ten, eighteen and nineteen respectively; whereas the corresponding expected results obtained in Analysis are between nine and ten, fifteen and sixteen, eighteen and nineteen respectively. There is a gap between the expected performance (in Section 3) and the

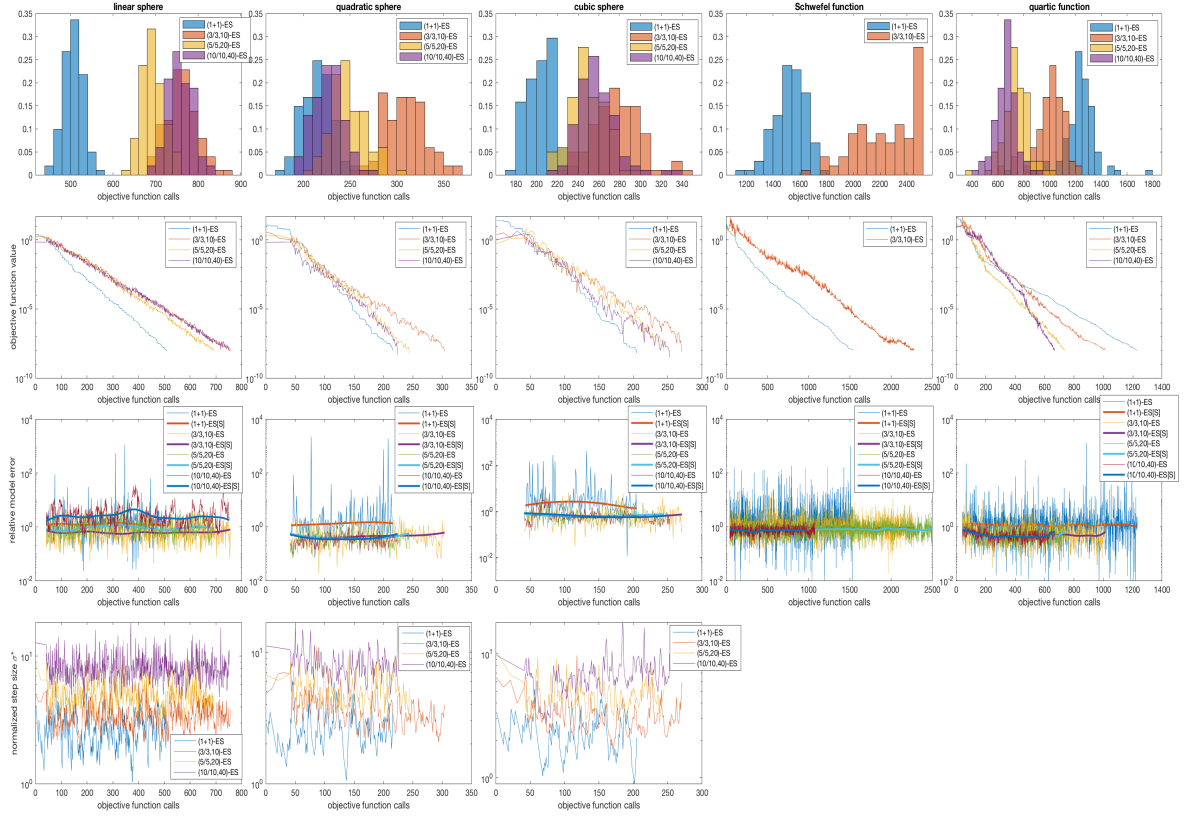


Fig. 3. Results for GP- $(\mu/\mu, \lambda)$ -ES (with GP and CSA applied). Top row: Histograms showing the number of objective function calls needed to solve the five test problems. Second row: Convergence graphs observed in median runs. Third row: Relative model error obtained in median runs ([S] denotes the smoothed relative model error). Last row: Normalized step size measured in median runs for sphere functions.

experimental results (after the GP surrogate model and CSA is in place) for GP- $(\mu/\mu, \lambda)$ -ES. But the difference between expected values and observed values obtained through experiments seems to shrink as the population size  $\lambda$  grows. In terms of  $\text{speed-up}_{\text{model}}$ , the achieved values for GP- $(\mu/\mu, \lambda)$ -ES on quadratic sphere are about 0.7, 0.9, and 0.9 for  $\lambda = 10, 20, 40$  respectively, which is far from the expected maximal value of roughly 2.4, 3.7, and 4.6 obtained in Table 2 when noise-to-signal ratio  $\eta = 1$ . There is about a factor of four in terms of the expected maximal  $\text{speed-up}_{\text{model}}$  and the results obtained in table 4. The  $\text{speed-up}_{\text{self}}$  observed for sphere functions are between four and five for GP- $(\mu/\mu, \lambda)$ -ES with  $\lambda = 10$ , six and ten for GP- $(\mu/\mu, \lambda)$ -ES with  $\lambda = 20$ , and six to nineteen for GP- $(\mu/\mu, \lambda)$ -ES with  $\lambda = 40$ . It is interesting that only GP- $(\mu/\mu, \lambda)$ -ES using  $\lambda = 10$  solves the 10-dimensional Schwefel's function, the increasing population does not give an advantage for solving the problem compared with the (1+1)-ES with surrogate model assistance, but rather a trend for divergence. Overall, the performance of GP- $(\mu/\mu, \lambda)$ -ES does not match that of the surrogate model assisted (1+1)-ES, specifically,  $\text{speed-up}_{\text{model}} < 1$  for all test functions except quartic function where the observed  $\text{speed-up}_{\text{model}}$  are about 1.3, 1.7, and 1.9 for population size  $\lambda = 10, 20, 40$  respectively. It seems a growing population can possibly bring larger  $\text{speed-up}_{\text{model}}$  for quadratic sphere and quartic function, meanwhile



it appears contrary to the other three test functions, especially the Schwefel's function. The speed-up<sub>self</sub> observed in linear and cubic sphere is lower than that in quadratic sphere, indicating a potential less accurate Gaussian process based surrogate model in the two sphere functions mentioned compared with the latter.

The second row of Fig. 3 shows the convergence graphs observed in median runs. Linear convergence are achieved for all but the Schwefel's function, specifically the using GP-( $\mu/\mu, \lambda$ )-ES with  $\lambda = 20, 40$ . Interestingly, using a larger population on Schwefel's function does not help achieve faster convergence compared to (1+1)-ES, but instead, slows the rate of convergence, and finally makes the strategy diverge. Relative model error for the median runs is shown in the third row of the figure. The relative model error in generation ( $g$ ) is defined as

$$\begin{cases} \frac{\|f(y^{(g)}) - f_\epsilon(y^{(g)})\|}{\|f(y) - f(x)\|}, & \text{for surrogate model assisted (1+1)-ES} \\ \sqrt{\frac{\text{var}(fY^{(g)} - fY_\epsilon^{(g)})}{\text{var}(fY^{(g)} - fX^{(g)})}}, & \text{for surrogate model assisted } (\mu/\mu, \lambda)\text{-ES,} \end{cases} \quad (27)$$

where  $x^{(g)}$  and  $y^{(g)}$  are the parent and offspring in generation ( $g$ ) for surrogate model assisted (1+1)-ES from Kayhani and Arnold [11];  $fY^{(g)} - fX^{(g)} = \{f(x_{\text{centroid}}^{(g)} + \sigma^{(g)}z_{\text{step}}^{(g)}) - f(x_{\text{centroid}}^{(g)}) : 1 \leq i \leq \lambda\}$ ,  $fY^{(g)} - fY_\epsilon^{(g)} = \{f(x_{\text{centroid}}^{(g)} + \sigma^{(g)}z_{\text{step}}^{(g)}) - f_\epsilon(x_{\text{centroid}}^{(g)} + \sigma^{(g)}z_{\text{step}}^{(g)}) : 1 \leq i \leq \lambda\}$ , and  $\text{var}(fY^{(g)} - fY_\epsilon^{(g)})$  returns the variance of the  $\lambda$  differences between the true and estimated objective function values for the offspring generated in generation ( $g$ ), and by taking square root over the division of two variances we can obtain the standard deviation similar to the case of surrogate model assisted (1+1)-ES. For easy comparison, the relative model error is smoothed logarithmically by convolution with a Gaussian kernel with window size 40 that is represented as the bold line in the centre of the plots (denoted [S] in the Fig.), interpreted as the a relative constant noise-to-signal ratio. The relative model error for all surrogate assisted ES in observed is approximately 1; but we notice the relative model error for GP-( $\mu/\mu, \lambda$ )-ES is much smaller than that of the surrogate assisted (1+1)-ES, so is the variance of the relative model error. This is possibly the benefit of using a large population size. In linear sphere, the relative model error obtained by using GP-( $\mu/\mu, \lambda$ )-ES with  $\lambda = 40$  is significantly higher compared with the rest strategies, which may explain the decline in speed-up<sub>model</sub> from 0.72 to 0.66 after  $\lambda$  increased from 20 to 40. In the other problems, the relative model error for GP-( $\mu/\mu, \lambda$ )-ES with different  $\lambda$  seem to stay at the same level. In the quadratic sphere, the relative constant noise-to-signal ratio is approximately 0.7 for GP-( $\mu/\mu, \lambda$ )-ES with different  $\lambda$ ; and according to Fig. 2, when  $\vartheta = 0.7$ , the opt. step size are between two and three, four and five, ten and eleven for  $\lambda = 10, 20, 40$  respectively. The corresponding expected speed-up<sub>model</sub> for  $\lambda = 10, 20, 40$  obtained by using the opt. expected fitness gain from Fig. 2 when  $\vartheta = 0.7$  are about 1.7, 3.1, and 4.3 for  $\lambda = 10, 20, 40$  respectively. The bottom row in Fig. 3 shows the normalized step size over the number of objective function calls observed in the median runs in the three sphere functions tested. The GP-( $\mu/\mu, \lambda$ )-ES achieves a larger normalized step size  $\sigma^*$  compared to the surrogate model assisted (1+1)-ES with  $\sigma^*$  growing in line with  $\lambda$ , which coincides with our knowledge that using a larger population size in ( $\mu/\mu, \lambda$ )-ES can give a larger step size i.e., in the extreme case when  $\lambda \rightarrow \infty$ , we obtain the gradient at the current point (parent  $x$ ). Given the relative constant noise-to-signal is 0.7 (i.e.,  $\vartheta = 0.7$ ) for quadratic sphere, we can obtain the opt. expected normalized step size  $\sigma^*$  from Fig. 2 and compare that with the experimental result obtained in the bottom row of Fig. 3. The obtained  $\sigma^*$  in experiment for each GP-( $\mu/\mu, \lambda$ )-ES with different  $\lambda$  is very close to that of the analytical result shown in Fig. 2, where the former ranges from two to four, four to six, nine to eleven for  $\lambda = 10, 20, 40$  respectively, and the latter has discussed above. The results indicates that the step size of the GP-( $\mu/\mu, \lambda$ )-ES is appropriately adapted by CSA.

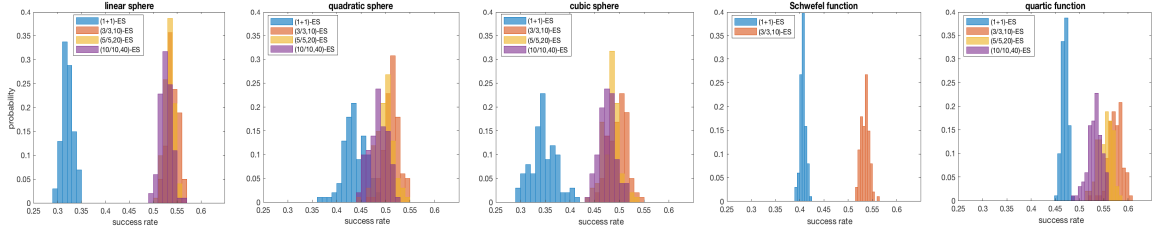


Fig. 4. Histograms of the success rate (proportion of good steps in a single run) for each test function using  $(1+1)$ -ES and  $GP-(\mu/\mu, \lambda)$ -ES with  $\lambda = 10, 20, 40$ .

Despite the potentially well-adapted step size, there is a big gap between the expected  $speed-up_{model}$  discussed above and the experimental result illustrated in Table 4. The reasons behind the resulting experimental results are not yet clear. Assuming the step size is well-adapted, the natural question to ask is what is the quality of the direction ( $z_{step}$ ) used. If the resulting direction is poor (i.e., taking the direction cannot improve the fitness), the individual cannot gain any fitness even with an optimally adapted step size. On the other hand, taking a good step with a properly adapted step size can help generate better offspring (offspring with larger fitness gain compared with its parent). We further plot the histogram of success rate (the proportion of a good step in each run) for all test problems in Fig. 4 by replicating 100 runs for each test problems. Overall, the success rates for  $GP-(\mu/\mu, \lambda)$ -ES are higher compared with the surrogate model assisted  $(1+1)$ -ES, where the latter achieves about 0.32, 0.42, 0.35, and 0.4 for linear sphere, quadratic sphere, cubic sphere, Schwefel’s function, and quartic function respectively. Whereas the success rates for  $GP-(\mu/\mu, \lambda)$ -ES with  $\lambda = 10, 20, 40$  are at about the same level in all five test problems, specifically, ranging from 0.45 to 0.55 in three sphere functions, about 0.53 in Schwefel’s function, and between 0.5 and 0.6 in quartic function. Recombination and selection over a larger population size in  $(\mu/\mu, \lambda)$ -ES do help improve the chance of making a good step compared with the  $(1+1)$ -ES. It is interesting that the success rate of an increasing  $\lambda$  first increases and then decreases. The benefit of using a larger step size in terms of sampling a good direction in the context of a GP surrogate model is not yet clear. The overall success rate of  $GP-(\mu/\mu, \lambda)$ -ES approximately 0.5, meaning the strategy makes a bad step every other step.

## 4.2 Safeguard of plus-selection

The  $GP-(\mu/\mu, \lambda)$ -ES gives a success rate approximately equals to 0.5 for all population sizes used in all test functions. It comes natural to ask, how much we are to benefit if we can avoid or simply reject those bad steps. A recent paper in surrogate model assisted ES considers  $(1+1)$ -ES [11] where the step size of the strategy is successfully adapted based on the success rate of a good step size. The mentioned step size adaptation mechanism considers decreasing the step size in two cases. In the first case, the estimated fitness of the offspring (using the surrogate model) is inferior to the true fitness of its parent; as for the second case, the true fitness of the offspring evaluated is inferior to that of its parent. Using a similar idea, we propose a surrogate model assisted ES that is a cross between  $(1+1)$ -ES and  $(\mu/\mu, \lambda)$ -ES (referred to as GP-cross-ES) and a new step size adaptation mechanism using CSA and the safeguard of plus-selection.

The safeguard of plus-selection is used to reject the poor steps (resulted from a bad direction that cannot improve the fitness of the current parent). If the offspring is inferior to its parent, meaning the step generated in this generation is bad, the safeguard of plus-selection simply rejects the inferior offspring. The proposed GP-cross-ES and the step size adaptation are shown in Alg. 3: in each generation, one offspring  $y = x_{centroid} + \sigma^{(g)} z_{step}^{(g)}$  is evaluated using the true

---

**Algorithm 3** Surrogate model assisted ES with plus-selection (GP-cross-ES)

---

```

1:  $c \leftarrow \frac{\mu+2}{N+\mu+5}$ 
2:  $d \leftarrow 1 + 2\max(0, \sqrt{\frac{\mu-1}{N+1}} - 1) + c$ 
3:  $E \|\mathcal{N}(0, I)\| \approx \sqrt{N}(1 - \frac{1}{4N} + \frac{1}{21N^2})$ 
4:  $s^{(0)} \leftarrow 0$ 
5:  $D \leftarrow 0.72$ 
6: while not terminate() do
7:   for  $i = 1, 2, \dots, \lambda$  do
8:     Generate standard normally distributed  $z_i^{(g)} \in \mathbb{R}^N$ 
9:     Evaluate  $x_{\text{centroid}} + \sigma^{(g)} z_i^{(g)}$  using the GP surrogate model, yielding  $f_\epsilon(x_{\text{centroid}} + \sigma^{(g)} z_i^{(g)})$ 
10:   end for
11:    $z_{\text{step}}^{(g)} \leftarrow \frac{1}{\mu} \sum_{i=1}^{\mu} z_{i;\lambda}^{(g)}$ 
12:    $y \leftarrow x_{\text{centroid}} + \sigma^{(g)} z_{\text{step}}^{(g)}$ 
13:   Evaluate  $y$  using true objective function, yielding  $f(y)$ 
14:   Update surrogate model by adding  $(y, f(y))$ 
15:   if  $f(x_{\text{centroid}}) < f(y)$  (occurrence of a bad step) then
16:      $\sigma \leftarrow \sigma D$ 
17:   else
18:      $x_{\text{centroid}} \leftarrow y$ 
19:      $s^{(g)} \leftarrow (1 - c)s^{(g)} + \sqrt{c(2 - c)}\mu z_{\text{step}}^{(g)}$ 
20:      $\sigma^{(g+1)} \leftarrow \sigma^{(g)} \exp\left(\frac{c}{d} \frac{\|s^{(g)}\|}{E\|\mathcal{N}(0, I)\|} - 1\right)$ 
21:   end if
22:    $g \leftarrow g + 1$ 
23: end while

```

---

objective function and added to the training set to update GP; the true fitness of the offspring  $f(y = x_{\text{centroid}} + \sigma^{(g)} z_{\text{step}}^{(g)})$  is compared to that of its parent  $f(x_{\text{centroid}})$ . If the fitness of the offspring is inferior to its parent, indicating the resulted step  $z_{\text{step}}^{(g)}$  is poor in this generation, the offspring  $y$  is discarded and the step size is decreased by a factor of  $D$  using an idea similar to the step size adaptation of the surrogate model assisted (1+1)-ES [11]; the bad step information is not added to the evolution path  $s^{(g)}$  because we want to build an evolution path that is based on the good step information of previous generations. Otherwise (the step made is good), we update the parent by setting  $x_{\text{centroid}} = y$  and update the step size accordingly using CSA.

To test the proposed GP-cross-ES and the step size adaptation mechanism, we use the same test functions and generate the corresponding Figs in Section 4.1. The number of objective function evaluations used in median runs and the corresponding speed-ups are reported in Table 5. The overall performance of GP-cross-ES improves as the population size increases, and outperforms the surrogate model assisted (1+1)-ES [11] in all test functions after  $\lambda \geq 20$ . The speed-up<sub>model</sub> for all three population sizes tested is around 1.5 for the linear sphere, between 1 and 1.5 for both the quadratic sphere and cubic sphere, about 0.7 to 1.4 for the Schwefel's function, and between 0.8 and 1.6 for the quartic function. The linear sphere does not seem to benefit much from a growing population size and the speed-up<sub>model</sub> seems to plateau after  $\lambda \geq 20$ ; whereas in the quadratic sphere, Schwefel's function, and quartic function, the speed-up<sub>model</sub>

Table 5. Median test results and speed-ups of GP-cross-ES.

Test functions	Median number of objective function calls (speed-up <sub>self</sub> , speed-up <sub>model</sub> )			
	(1 + 1)-ES	(3/3, 10)-ES	(5/5, 20)-ES	(10/10, 40)-ES
linear sphere	502	367(9.0, 1.4)	316(15.1, 1.6)	321(26.3, 1.6)
quadratic sphere	212	211(8.0, 1.0)	164(14.9, 1.3)	146(29.0, 1.5)
cubic sphere	202	213(5.5, 0.96)	176(9.4, 1.3)	178(15.8, 1.4)
Schwefel's function	1503	2070(3.0, 0.73)	1355(6.0, 1.1)	1051(12.7, 1.4)
quartic function	1250	1511(4.4, 0.83)	1016(8.3, 1.3)	796(18.4, 1.6)

Note: the value shown for (1+1)-ES is its corresponding speed-up<sub>self</sub>.

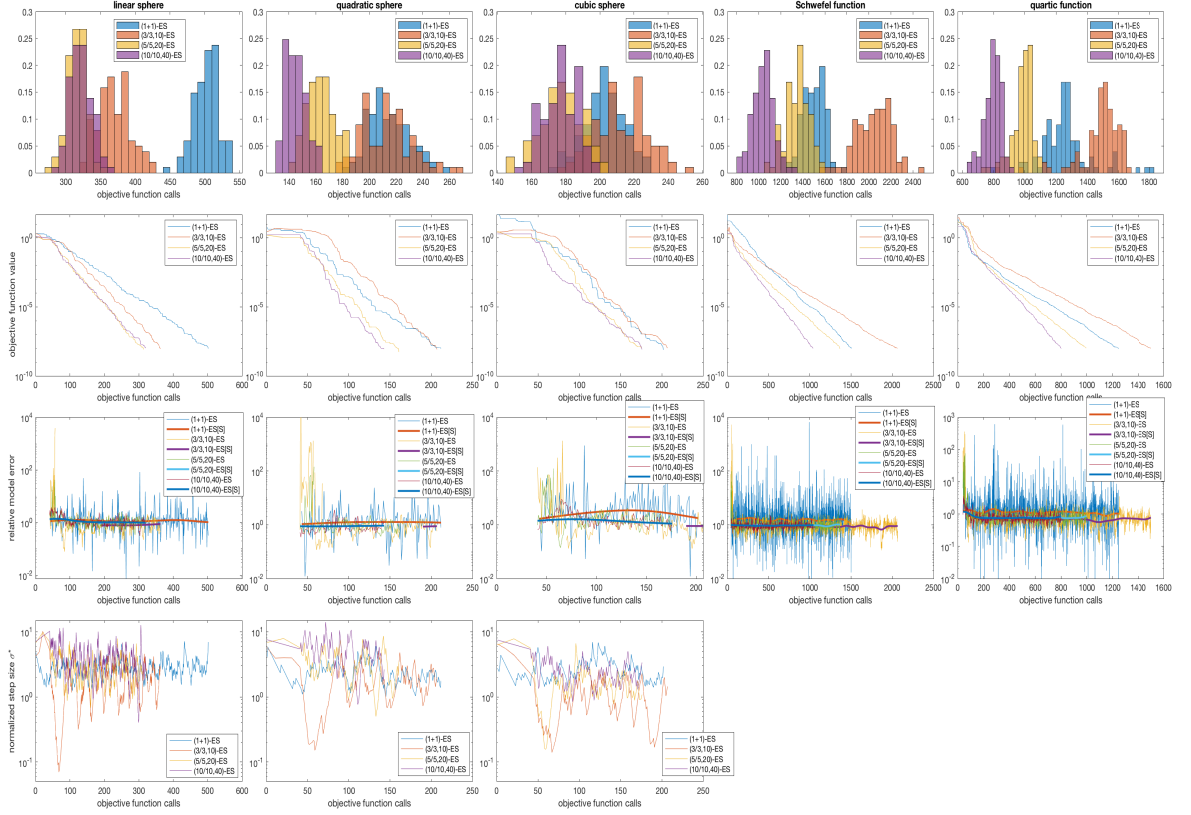


Fig. 5. Results of GP-cross-ES obtained by applying the safeguard of plus-selection. Top row: Histogram showing the number of objective function calls needed to solve the five test problems. Second row: Convergence graphs observed in median runs. Third row: Relative model error obtained in median runs ([S] denotes the smoothed plot). Last row: normalized step size measured in median runs for sphere functions.

does not show a clear sign of plateau. This observation also coincides with the speed-up<sub>self</sub> observed on the quadratic sphere, Schwefel's function, and quartic function where the growth in speed-up is proportional to the growth in population size. For example, in the Schwefel's function, the speed-up<sub>self</sub> doubles from 3.0 to 6.0 then approximately

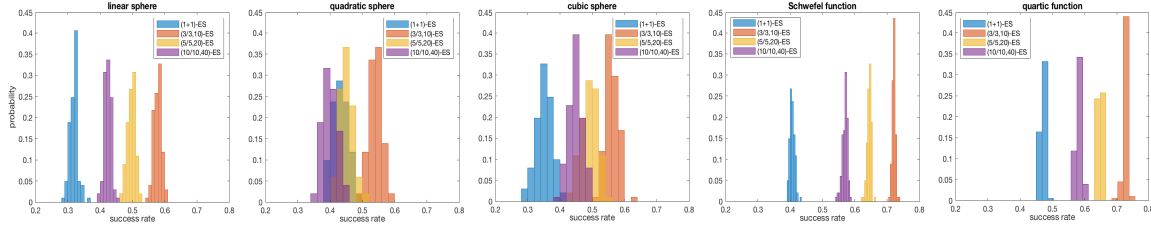


Fig. 6. Result obtained by applying the safeguard of plus-selection. The histograms shows the success rate (proportion of good steps in a single run) for each test function using (1+1)-ES and GP-cross-ES with  $\lambda = 10, 20, 40$ .

doubles to 12.7 after the population size doubles from 10 to 20 and finally 40. The histograms of objective function calls that are needed to solve each test problem (first row of Fig. 5) also illustrate the decrease in objective function calls needed to solve each test problem as the population size grows. Convergence graphs for the medium run are shown in the second row of Fig. 5. Linear convergence appears to be achieved in all runs. The third row reports the relative model error defined in Section 4.1. It is interesting that the relative model error for proposed GP-cross-ES (in any case of  $\lambda$  being used) is still lower than that of the surrogate model assisted (1+1)-ES. One possible explanation for a larger variance in the relative model error observed for the surrogate model assisted (1+1)-ES (compared with the GP- $(\mu/\mu, \lambda)$ -ES and GP-cross-ES) is the small population size (only a single offspring is generated), whereas the GP- $(\mu/\mu, \lambda)$ -ES and GP-cross-ES generate  $\lambda$  offspring in each generation. The bottom row in Fig. 5 shows the normalized step size observed in median runs for the three sphere functions used. The normalized step size  $\sigma^*$  in all runs are approximately between 0.5 and 10, where the strategy using plus-selection loses the benefit of a potential larger step size as opposed to the comma-selection (e.g., the  $\sigma^*$  for the GP- $(\mu/\mu, \lambda)$ -ES is significantly higher than that of the GP-cross-ES). Another observation is the decreasing trend in the variance of relative model error with a growing  $\lambda$  for GP-cross-ES, which indicates the possibility of a potential good step resulting from using a large population size in selection and recombination. An illustration of the continuous occurrence of bad steps that can result in a continuous decrease of  $\sigma^*$  is shown in the first 80 objective function calls in linear sphere for (3/3,10)-ES (GP-cross-ES with  $\lambda = 10, \mu = 3$ ). But the occurrence or the length (the number of objective functions such situations take up) of such situations is reduced as  $\lambda$  increases.

Similarly, the histograms of success rates (the proportion of a good step size in each run) for all test problems are plotted in Fig. 6 by replicating 100 runs for each test problem. The success rate for GP-cross-ES is overall larger than that of the surrogate model assisted (1+1)-ES with the exception of the quadratic sphere in which the success rate of GP-cross-ES with  $\lambda = 10$  and 20 being the same level as surrogate model assisted (1+1)-ES. For all other test functions, the success rate for GP-cross-ES decreases as  $\lambda$  increases. In both the linear sphere and cubic sphere, the success rate for GP-cross-ES ranges from 0.4 to 0.5, 0.45 to 0.55, and 0.55 to 0.65 for  $\lambda = 10, 20$ , and 40 respectively; whereas surrogate model assisted (1+1)-ES being about 0.3 to 0.35. The success rate observed in the Schwefel's function and quartic function are between 0.55 and 0.61; 0.62 and 0.65; and 0.70 and 0.75 for GP-cross-ES with  $\lambda = 10, 20$ , and 40 respectively, and ranges from 0.55 to 0.62 for surrogate model assisted (1+1)-ES. Even if the success rate reduces with a larger  $\lambda$  for GP-cross-ES, from the sharper slope, compared with GP- $(\mu/\mu, \lambda)$ -ES, in the convergence plots, we can infer that the step is regarded as "good" by the GP-cross-ES has a potential better quality than that of GP- $(\mu/\mu, \lambda)$ -ES. That is, the

fitness gain obtained by taking the step is larger than that of the GP- $(\mu/\mu, \lambda)$ -ES. By using a larger  $\lambda$ , the GP-cross-ES with proposed step-size adaption mechanism takes "good" steps that can bring potential larger fitness gain compared with the GP- $(\mu/\mu, \lambda)$ -ES. One possible explanation of the better performance is the quality increase of taking the good steps outweighs the reduce in probability regarding the occurrence of a good step .

## 5 CONCLUSIONS

In this thesis, we applied the simple model of unbiased Gaussian distributed noise for surrogate modelling approach from Kayhani and Arnold [11] to analyze the surrogate model. By using this approach, we analyzed the behaviours of the proposed GP- $(\mu/\mu, \lambda)$ -ES on the quadratic sphere function and observed a very significant speed-up<sub>self</sub> as the surrogate model is more fully exploited compared with the surrogate model assisted (1+1)-ES. If the surrogate model is accurate, we would expect a speed-up<sub>self</sub> equals to the population size  $\lambda$ . In experiments, the speed-up<sub>self</sub> achieved is about half the expected value for the quadratic sphere and quartic function; for linear sphere and cubic sphere, the speed-up<sub>self</sub> achieved is a factor between two and six smaller than that of the expected value. In the Schwefel's function, the GP- $(\mu/\mu, \lambda)$ -ES does not even converge when  $\lambda \geq 20$ . In most cases with the exception of the quartic function, the performance of the GP- $(\mu/\mu, \lambda)$ -ES is inferior to that of the surrogate model assisted (1 + 1)-ES. There is a factor approximately four between the analytical results and the experimental results for speed-up<sub>model</sub> in the quadratic sphere function. One possibility is the assumption that the Gaussian noise are uncorrelated may not hold in the case of generating multiple offspring in each generation as is happened in  $(\mu/\mu, \lambda)$ -ES, but can be valid when generating a single offspring in (1+1)-ES.

We tried to interpret the cause of the gap from the single step behaviour of the strategy. One observation is the increased success rate but a potential decreased quality of a good step compared with the surrogate model assisted (1+1)-ES, indicating many of the steps generated do not give much improvement (fitness gain). Based on the analysis and the observations of the GP- $(\mu/\mu, \lambda)$ -ES, we proposed the GP-cross-ES and a corresponding step size adaptation mechanism taking benefit of the plus-selection where the bad steps are rejected. The strategy is evaluated numerically using a set of test functions. It shows that the step size adaptation mechanism adapted the step size successfully in all runs and the GP-cross-ES achieves an overall speed-up<sub>model</sub> about 1.4 to 1.6 for a population size  $\lambda \geq 20$  in all five test functions.

In future work, we will study the behaviours of surrogate assisted CMA-ES using the same analysis that can potentially handle ill-conditioned problems. Further goals include length scale adaptation mechanism in the Gaussian Process and surrogate model accuracy control that can possibly further reduce the gap between the expected analytical results and the experimental results.

## ACKNOWLEDGMENT

I would like to express my sincere gratitude my supervisor Dr. Dirk V. Arnold for his continuous support during the past two semesters, as well as his patience, motivation, and passion that motivates me all the time both in study and beyond. Besides my supervisor, I would like to thank my parents for their love and unconditional support during the past years without which I would not be where I am now.

## REFERENCES

- [1] D. V. Arnold and H.-G. Beyer. 2000. Efficiency and mutation strength adaptation of the  $(\mu/\mu, \lambda)$ -ES in a noisy environment. In *Proceedings of the 6th International Conference on Parallel Problem Solving from Nature (PPSN VI)*. Springer-Verlag, London, UK, 39–48.

- [2] D. V. Arnold and H.-G. Beyer. 2001. Local performance of the  $(\mu/\mu, \lambda)$ -ES in a noisy environment. In *Foundations of Genetic Algorithms 6*. W. N. Martin and W. M. Spears, (Eds.) Morgan Kaufmann, San Francisco, 127–141.
- [3] D. V. Arnold and H.-G. Beyer. 2004. Performance analysis of evolutionary optimization with cumulative step length adaptation. *IEEE Transactions on Automatic Control*, 49, 4, (Apr. 2004), 617–622.
- [4] C. G. Atkeson, A. W. Moore, and S. Schaal. 1997. Locally weighted learning for control. In *Lazy learning*. Springer, 75–113.
- [5] H. G. Beyer. 1995. Toward a theory of evolution strategies: on the benefits of sex—the  $(\mu/\mu, \lambda)$  theory. *Evolutionary Computation*, 3, 1, (Mar. 1995), 81–111.
- [6] H.-G. Beyer. 2013. *The theory of evolution strategies*. Springer Science & Business Media.
- [7] D. Buche, N. N. Schraudolph, and P. Koumoutsakos. 2005. Accelerating evolutionary algorithms with gaussian process fitness function models. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 35, 2, (May 2005), 183–194.
- [8] N. Hansen. 2016. The CMA evolution strategy: a tutorial. *arXiv preprint arXiv:1604.00772*.
- [9] N. Hansen and S. Kern. 2004. Evaluating the CMA evolution strategy on multimodal test functions. In *International Conference on Parallel Problem Solving from Nature*. Springer, 282–291.
- [10] Y. Jin. 2011. Surrogate-assisted evolutionary computation: recent advances and future challenges. *Swarm and Evolutionary Computation*, 1, 2, 61–70.
- [11] A. Kayhani and D. V. Arnold. 2018. Design of a surrogate model assisted  $(1 + 1)$ -ES. In *Parallel Problem Solving from Nature - PPSN XV - 15th International Conference, Coimbra, Portugal, September 8-12, 2018, Proceedings, Part I*, 16–28.
- [12] S. Kern, N. Hansen, and P. Koumoutsakos. 2006. Local meta-models for optimization using evolution strategies. In *Parallel Problem Solving from Nature - PPSN IX*. T. P. Runarsson et al., (Ed.) Springer Berlin Heidelberg, 939–948.
- [13] S. Kern, S. D. Müller, N. Hansen, D. Büche, J. Ocenasek, and Petros Koumoutsakos. 2004. Learning probability distributions in continuous evolutionary algorithms—a comparative review. *Natural Computing*, 3, 1, 77–112.
- [14] I. Loshchilov. 2016. LM-CMA: an Alternative to L-BFGS for Large Scale Black-box Optimization. *Evolutionary Computation*.
- [15] K. P. Murphy. 2012. *Machine Learning: A Probabilistic Perspective*. The MIT Press.
- [16] D. V. Arnold N. Hansen and A. Auger. 2015. Evolution strategies. In *Springer Handbook of Computational Intelligence*. Springer, 871–898.
- [17] G. Andreas O. Andreas and H. Nikolaus. 1994. A derandomized approach to self-adaptation of evolution strategies. *Evolutionary Computation*, 2, 4, (Dec. 1994), 369–380.
- [18] C. E. Rasmussen. 2004. Gaussian processes in machine learning. In *Advanced Lectures on Machine Learning*. Springer, 63–71.
- [19] A. Ratle. 1998. Accelerating the convergence of evolutionary algorithms by fitness landscape approximation. In *Proceedings of the 5th International Conference on Parallel Problem Solving from Nature (PPSN V)*. Springer-Verlag, London, UK, 87–96.
- [20] I. Rechenberg. 1973. *Evolutionstrategie : Optimierung technischer Systeme nach Prinzipien der biologischen Evolution*. Number 15 in *Problemata*. Frommann-Holzboog, Stuttgart-Bad Cannstatt.
- [21] H. -P. Schwefel. 1981. *Numerical Optimization of Computer Models*. John Wiley & Sons, Inc., New York, NY, USA.