

1 **ACCELERATING REINFORCEMENT LEARNING WITH A**
2 **DIRECTIONAL-GAUSSIAN-SMOOTHING EVOLUTION**
3 **STRATEGY**

JIAXIN ZHANG

Computer Science and Mathematics Division
Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

HOANG TRAN

Computer Science and Mathematics Division
Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

GUANNAN ZHANG*

Computer Science and Mathematics Division
Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

ABSTRACT. The objective of reinforcement learning (RL) is to find an optimal strategy for solving a dynamical control problem. Evolution strategy (ES) has been shown great promise in many challenging reinforcement learning (RL) tasks, where the underlying dynamical system is only accessible as a black box such that adjoint methods cannot be used. However, existing ES methods have two limitations that hinder its applicability in RL. First, most existing methods rely on Monte Carlo based gradient estimators to generate search directions. Due to low accuracy of Monte Carlo estimators, the RL training suffers from slow convergence and requires more iterations to reach the optimal solution. Second, the landscape of the reward function can be deceptive and may contain many local maxima, causing ES algorithms to prematurely converge and be unable to explore other parts of the parameter space with potentially greater rewards. In this work, we employ a Directional Gaussian Smoothing Evolutionary Strategy (DGS-ES) to accelerate RL training, which is well-suited to address these two challenges with its ability to (i) provide gradient estimates with high accuracy, and (ii) find nonlocal search direction which lays stress on large-scale variation of the reward function and disregards local fluctuation. Through several benchmark RL tasks demonstrated herein, we show that the DGS-ES method is highly scalable, possesses superior wall-clock time, and achieves competitive reward scores to other popular policy gradient and ES approaches.

2020 *Mathematics Subject Classification.* Primary: 37M05, 93E35; Secondary: 60J25.

Key words and phrases. Reinforcement learning, Gaussian smoothing, non-convex optimization, stochastic control, high-dimensional optimization, Gauss-Hermite quadrature.

* Corresponding author: Guannan Zhang (zhangg@ornl.gov).

Notice: This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05-00OR22725 with the US Department of Energy (DOE). The US government retains and the publisher, by accepting the article for publication, acknowledges that the US government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this manuscript, or allow others to do so, for US government purposes. DOE will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).

1 **1. Introduction.** Reinforcement learning is a class of problems which aim to find,
 2 through trial and error, a feedback policy that prescribes how an agent should
 3 act in an uncertain, complex environment to maximize some notion of cumulative
 4 reward [30]. RL can be viewed as a generalization of stochastic optimal control.
 5 Traditionally, RL algorithms have mainly been employed for small input and action
 6 spaces, and suffered difficulties when scaling to high-dimensional problems. With the
 7 recent emergence of deep learning, powerful non-linear function approximators such
 8 as deep neural networks (DNN) can be integrated into RL and extend the capability
 9 of RL in a variety of challenging tasks which would otherwise be infeasible, ranging
 10 from playing Atari from pixels [18], playing expert-level Go [28] to robotic control
 11 [2]. Among the most popular current deep RL algorithms are Q-learning methods,
 12 policy gradient methods, and evolution strategies. Deep Q-learning algorithms [18]
 13 use a DNN to approximate the optimal Q function, yielding policies that, for a given
 14 state, choose the action that maximizes the Q-value. Policy gradient methods [26]
 15 improve the policies with a gradient estimator obtained from sample trajectories in
 16 action space, examples of which are A3C [17], TRPO [25] and PPO [24].

17 This work concerns RL techniques based on evolution strategies. ES refers to a
 18 family of blackbox optimization algorithms inspired by ideas of natural evolution,
 19 often used to optimize functions when gradient information is inaccessible. This is
 20 exactly the prominent challenge in a typical RL problem, where the environment and
 21 policy are usually nonsmooth or can only be accessed via noisy sampling. It is not
 22 totally surprising ES has become a convincing competitor to Q-learning and policy
 23 gradient in deep RL in recent years. Unlike policy gradient, ES perturbs and performs
 24 policy search directly in the parameter space to find an effective policy, which is
 25 now generally considered to be superior to action perturbation [27]. The policy
 26 search can be guided by a surrogate gradient [23], completely population-based and
 27 gradient-free [29], and hybridized with other exploration strategies such as novelty
 28 search and quality diversity [7]. It has been shown in those works that ES is easy to
 29 parallelize, and requires low communication overhead in a distributed setting. More
 30 importantly, being able to achieve competitive performance on many RL tasks, these
 31 methods are advantageous over other RL approaches in their high scalability and
 32 substantially lower wall-clock time. Given wide availability of distributed computing
 33 resources, all environment simulations at each iteration of training can be conducted
 34 totally in parallel. Thus, a more reasonable metric for the performance of a training
 35 algorithm is the number of non-parallelizable iterations, as opposed to the sample
 36 complexity. In this metric, ES is truly an appealing choice.

37 Nevertheless, there are several challenges that need to be addressed in order to
 38 further improve the performance of ES in training complex policies. First, most
 39 ES methods cope with the non-smoothness of the objective function by considering
 40 a Gaussian-smoothed version of the expected total reward. The gradient of this
 41 function is intractable and must be estimated to provide the policy parameter
 42 update. In the pioneering work [23], a gradient estimator is proposed based on
 43 random parameter sampling. Developing efficient sampling strategies for gradient
 44 estimates has become an interest in ES research since then, and several improvements
 45 have been proposed, based on imposing structures on parameter perturbation [5], or
 46 reusing past evaluations, [6, 15, 16]. Yet, most of these gradient estimators are of
 47 Monte Carlo type, therefore arguably affected by the low accuracy of Monte Carlo
 48 methods. For faster convergence of training (i.e., reducing the number of iterations),
 49 more accurate gradient estimators are desired, particularly in RL tasks where the

1 policy has a large number of parameters to learn. Another prominent challenge
 2 is that the landscape of the objective function is complex and possesses plentiful
 3 local maxima. There is a risk for any optimization algorithm to get trapped in
 4 some of those points and unable to explore the parameter space effectively. The
 5 Gaussian smoothing, with its ability to smooth out function and damps out small,
 6 insignificant fluctuations, is a strong candidate in this very challenge. Specifically,
 7 with a moderately large smoothing parameter (i.e., strong smoothing effect), we
 8 can expect the gradient of the smoothed objective function will be able to look
 9 outside unimportant variation in the adjacent area and detect the general trend of
 10 the function from a distance, therefore an efficient *nonlocal* search direction. This
 11 potential of Gaussian smoothing, however, has not been explored in reinforcement
 12 learning.

13 In this paper, we propose a new strategy to accelerate the time-to-solution of
 14 reinforcement learning by exploiting the Directional Gaussian Smoothing Evolution
 15 Strategy (DGS-ES) method, developed in our recently work [33]. The DGS-ES
 16 method introduced a new directional Gaussian smoothing (DGS) gradient operator,
 17 that smooths the original objective function only along d -orthogonal directions in
 18 the parameter space. In other words, the DGS gradient requires d one-dimensional
 19 Gaussian convolutions, instead of one d -dimensional convolution in the existing
 20 ES methods. There are several advantages of using the DGS gradient operator in
 21 reinforcement learning. First, each component of the DGS gradient, represented
 22 as a one-dimensional integral, can be accurately approximated with various classic
 23 numerical integration techniques. When having Gaussian kernels, we use Gauss-
 24 Hermite quadrature rule which can provide spectral accuracy (see [1]) in the DGS
 25 gradient approximation. Second, the use of Gauss-Hermite quadrature also features
 26 embarrassing parallelism as the random sampling used in existing ES methods. Since
 27 the communication cost between computing processors/cores is neglectable, the total
 28 computing time for each iteration of training does not increase with the number of
 29 environment simulations given sufficient computing resources. Third, the directional
 30 smoothing approach enables nonlocal exploration which takes into account large
 31 variation of the objective function and disregards local fluctuations. This property
 32 will greatly help skipping local optima or saddle points during the training. It is
 33 demonstrated in §4 that the proposed strategy can significantly reduce the number of
 34 iterations in training several benchmark reinforcement learning problems, compared
 35 to the state-of-the-art baselines.

36 The rest of the paper is organized as follows: the RL problem under consideration
 37 and a brief review about the classic Gaussian smoothing technique is given in Section
 38 2, the DGS gradient and the corresponding DGS-ES algorithm is introduced in
 39 Section 3, extensive numerical experiments including tests on benchmark optimization
 40 problems and on several benchmark RL problems are provided in Section 4, some
 41 concluding remarks are given in Section 5.

42 **2. Problem setting.** We study the continuous-time, continuous-state stochastic
 43 control problem via RL. The evolution of the state $\mathbf{s}_t \in \mathcal{S}$ (the state space) depends
 44 on the control (i.e., the action in RL) $\mathbf{a}_t \in \mathcal{A}$ (the action space), by a stochastic
 45 differential equation, i.e., the state dynamics,

$$d\mathbf{s}_t = b(\mathbf{s}_t, \mathbf{a}_t)dt + \sigma(\mathbf{x}_t, \mathbf{a}_t)dw_t, \quad (1)$$

46 where b is the local drift, σ is the local diffusion, and w_t is the standard Brownian
 47 motion. Let $\pi(\mathbf{a}_t|\mathbf{s}_t)$ denote a control strategy (i.e., the policy in RL) that is a

- 1 conditional probability distribution of the action \mathbf{a}_t given the current state \mathbf{s}_t . Then
 2 the objective of RL is to maximize a functional \mathcal{J} of the policy π , i.e.,

$$\mathcal{J}(\pi) := \mathbb{E}_\pi \left[\int_0^T \gamma_t r(\mathbf{s}_t, \mathbf{a}_t) dt \right],$$

- 3 where T is the terminal time, $r(\mathbf{s}_t, \mathbf{a}_t)$ is the reward and $0 \leq \gamma \leq 1$ is a discount
 4 rate. After applying a temporal discretization scheme to the state process in (1),
 5 the continuous objective functional $\mathcal{J}(\pi)$ can be approximated by a discretized
 6 functional,

$$J(\pi) := \sum_{n=0}^N \mathbb{E}_{(\mathbf{s}_{t_n}, \mathbf{a}_{t_n})} [\gamma^{t_n} r(\mathbf{s}_{t_n}, \mathbf{a}_{t_n})], \quad (2)$$

- 7 where $t_n = n\Delta t$ for $n = 0, \dots, N$.

- 8 In policy-based RL, the policy $\pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta})$ is parameterized by $\boldsymbol{\theta} \in \mathbb{R}^d$, where the
 9 vector $\boldsymbol{\theta} := (\theta_1, \dots, \theta_d)^\top$ represents the parameters of the policy, e.g., weights of
 10 a neural network. Then, the task of learning a good policy π becomes iterative
 11 updating the parameter $\boldsymbol{\theta}$ to solve the following optimization problem

$$\max_{\boldsymbol{\theta} \in \mathbb{R}^d} J(\boldsymbol{\theta}), \quad (3)$$

- 12 where we denote $J(\boldsymbol{\theta}) := J(\pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}))$ by an abuse of notation. The main challenges
 13 for the RL problem under consideration lies in three aspects. First, it is usually
 14 true that the local gradient of the environment \mathcal{S} is inaccessible, so automatic
 15 differentiation or adjoint methods cannot be used to obtain the local gradient of
 16 $J(\boldsymbol{\theta})$. Thus, much of the innovation in reinforcement learning algorithms is focused on
 17 addressing the lack of access to or the existence of gradients of the environment/policy.
 18 Second, the landscape of $J(\boldsymbol{\theta})$ are usually highly non-convex and multi-modal, such
 19 that the local gradient (if available) may lead to an unsatisfactory local maximum of
 20 $J(\boldsymbol{\theta})$. Third, the dimension d of the parameters are usually on the order of hundreds
 21 or thousands (e.g., when using a neural network to define the policy), which makes
 22 many existing optimization methods inapplicable. These challenges motivated us to
 23 develop a nonlocal gradient operator and the corresponding approximation for the
 24 RL problem under consideration.

2.1. The standard Gaussian smoothing. We briefly recall the evolution strategy methods, e.g., [13, 23], which use a multivariate Gaussian distribution to generate the population around the current parameter value $\boldsymbol{\theta}_t$ at the t -iteration. When the Gaussian distribution can be factorized to d independent one-dimensional Gaussian distributions, the standard ES method can be mathematically interpreted based on the Gaussian smoothing technique [10, 20]. Specifically, a smoothed version of $J(\boldsymbol{\theta})$ in Eq. (2), denoted by $J_\sigma(\boldsymbol{\theta})$, is defined by

$$\begin{aligned} J_\sigma(\boldsymbol{\theta}) &:= \frac{1}{(2\pi)^{\frac{d}{2}}} \int_{\mathbb{R}^d} J(\boldsymbol{\theta} + \boldsymbol{\sigma} \circ \mathbf{u}) e^{-\frac{1}{2}\|\mathbf{u}\|_2^2} d\mathbf{u} \\ &= \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(0, \mathbf{I}_d)} [J(\boldsymbol{\theta} + \boldsymbol{\sigma} \circ \mathbf{u})], \end{aligned} \quad (4)$$

- 25 where $\mathcal{N}(0, \mathbf{I}_d)$ is a d -dimensional standard Gaussian distribution, the notation
 26 “ \circ ” represents the element-wise product, and the vector $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_d)$ controls
 27 the smoothing effect. It is well known that J_σ is always differentiable even if J
 28 is not. In addition, most of the characteristics of the original objective function
 29 $J(\boldsymbol{\theta})$, e.g., convexity, the Lipschitz constant, are inherited by $J_\sigma(\boldsymbol{\theta})$. When $J(\boldsymbol{\theta})$ is
 30 differentiable, the difference $\nabla J - \nabla J_\sigma$ can be bounded by its Lipschitz constant

1 (see [20], Lemma 3 for details). Thus, the original optimization problem in Eq. (3)
 2 can be replaced by a smoothed version, i.e.,

$$\max_{\boldsymbol{\theta} \in \mathbb{R}^d} J_{\boldsymbol{\sigma}}(\boldsymbol{\theta}), \quad (5)$$

3 where the gradient of $J_{\boldsymbol{\sigma}}(\boldsymbol{\theta})$ is given by

$$\nabla J_{\boldsymbol{\sigma}}(\boldsymbol{\theta}) = \frac{1}{\|\boldsymbol{\sigma}\|_2^{d/2}} \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(0, \mathbf{I}_d)} [J(\boldsymbol{\theta} + \boldsymbol{\sigma} \circ \mathbf{u}) \mathbf{u}]. \quad (6)$$

4 The standard ES method [23] uses Monte Carlo sampling to estimate the gradient
 5 $\nabla J_{\boldsymbol{\sigma}}(\boldsymbol{\theta})$ and update the state $\boldsymbol{\theta}$ from iteration n to $n + 1$ by

$$\boldsymbol{\theta}_{n+1} = \boldsymbol{\theta}_n - \frac{\lambda}{M\sigma} \sum_{m=1}^M J(\boldsymbol{\theta}_n + \boldsymbol{\sigma} \circ \mathbf{u}_m) \mathbf{u}_m, \quad (7)$$

6 where λ is the learning rate, \mathbf{u}_m are sampled from the Gaussian distribution $\mathcal{N}(0, \mathbf{I}_d)$.

7 One drawback of the ES method and its variants is the slow convergence of the
 8 training process, due to the low accuracy of the MC-based gradient estimator (see [3]),
 9 also [33], for extended discussions on the accuracy of gradient approximations using
 10 Eq. (7) and related methods). On the other hand, the evaluations of $J(\boldsymbol{\theta}_n + \boldsymbol{\sigma} \circ \mathbf{u}_m)$
 11 for $m = 1, \dots, M$ at the n -th iteration can be generated totally in parallel, which
 12 makes it well suited to be scaled up to a large number of parallel workers on
 13 modern supercomputers. Therefore, *the motivation of this work is to develop a new*
 14 *gradient operator to replace the one in Eq. (6), such that the new gradient can be*
 15 *approximated in a much more accurate way and the embarrassing parallelism feature*
 16 *can be retained.*

17 **3. The DGS-ES method for RL.** This section introduces our main framework.
 18 We start by introducing in Section 3.1 the DGS gradient operator and its approxi-
 19 mation using the Gauss-Hermite quadrature rule. In Section 3.2, we describe how
 20 to incorporate the DGS gradient operator into the ES for reinforcement learning.

21 **3.1. The DGS gradient operator.** For a given direction $\boldsymbol{\xi} \in \mathbb{R}^d$, the restriction
 22 of the objective function $J(\boldsymbol{\theta})$ along $\boldsymbol{\xi}$ can be represented by

$$G(y | \boldsymbol{\theta}, \boldsymbol{\xi}) = J(\boldsymbol{\theta} + y \boldsymbol{\xi}), \quad y \in \mathbb{R}, \quad (8)$$

23 where $\boldsymbol{\theta}$ is the current state of the agent's parameters. Then, we can define the
 24 one-dimensional Gaussian smoothing of $G(y)$, denoted by $G_{\sigma}(y)$, by

$$G_{\sigma}(y | \boldsymbol{\theta}, \boldsymbol{\xi}) := \mathbb{E}_{v \sim \mathcal{N}(0, 1)} [G(y + \sigma v | \boldsymbol{\theta}, \boldsymbol{\xi})], \quad (9)$$

25 which is also the Gaussian smoothing of $J(\boldsymbol{\theta})$ along $\boldsymbol{\xi}$ in the neighbourhood of $\boldsymbol{\theta}$.
 26 The derivative of $G_{\sigma}(y | \boldsymbol{\theta}, \boldsymbol{\xi})$ at $y = 0$ is given by

$$\mathcal{D}[G_{\sigma}(0 | \boldsymbol{\theta}, \boldsymbol{\xi})] = \frac{1}{\sigma} \mathbb{E}_{v \sim \mathcal{N}(0, 1)} [G(\sigma v | \boldsymbol{\theta}, \boldsymbol{\xi}) v], \quad (10)$$

27 where \mathcal{D} denotes the differential operator. It is easy to see that $\mathcal{D}[G_{\sigma}(0 | \mathbf{x}, \boldsymbol{\xi})]$ only
 28 involves the directionally smoothed objective function given in Eq. (9).

29 We can assemble a new gradient, i.e., the DGS gradient, by putting together the
 30 derivatives in Eq. (10) along d orthogonal directions, i.e.,

$$\nabla_{\boldsymbol{\sigma}, \boldsymbol{\Xi}}[J](\boldsymbol{\theta}) = \boldsymbol{\Xi}^{\top} \begin{bmatrix} \mathcal{D}[G_{\sigma_1}(0 | \boldsymbol{\theta}, \boldsymbol{\xi}_1)] \\ \vdots \\ \mathcal{D}[G_{\sigma_d}(0 | \boldsymbol{\theta}, \boldsymbol{\xi}_d)] \end{bmatrix}, \quad (11)$$

- 1 where $\Xi := (\xi_1, \dots, \xi_d)^\top$ represents the matrix consisting of d orthonormal vectors.
 2 It is important to notice that

$$\nabla J_\sigma(\theta) \neq \nabla_{\sigma, \Xi}[J](\theta)$$

- 3 for any $\sigma > 0$, because of the directional Gaussian smoothing used in Eq. (11).
 4 However, there is consistency between the two quantities as $\sigma \rightarrow 0$, i.e.,

$$\lim_{\sigma \rightarrow 0} |\nabla J_\sigma(\theta) - \nabla_{\sigma, \Xi}[J](\theta)| = 0, \quad (12)$$

- 5 for fixed θ and Ξ . If $\nabla J(\theta)$ exists, then $\nabla_{\sigma, \Xi}[J](\theta)$ will also converge to $\nabla J(\theta)$
 6 as $\sigma \rightarrow 0$. Such consistency naturally led to the idea of replacing $\nabla J_\sigma(\theta)$ with
 7 $\nabla_{\sigma, \Xi}[J](\theta)$ in the ES framework.

8 **3.2. The DGS-ES algorithm.** Since each component of $\nabla_{\sigma, \Xi}[J](\theta)$ in Eq. (11)
 9 only involves a one-dimensional integral, we can use Gaussian quadrature rules [22]
 10 to obtain spectral convergence. In the case of Gaussian smoothing, a natural choice
 11 is the Gauss-Hermite (GH) rule, which is used to approximate integrals of the form
 12 $\int_{\mathbb{R}} g(x) e^{-x^2} dx$. By doing a simple change of variable in Eq. (10), the GH rule can
 13 be directly used to obtain the following estimator:

$$\tilde{\mathcal{D}}^M[G_\sigma(0 | \theta, \xi)] := \frac{1}{\sqrt{\pi}\sigma} \sum_{m=1}^M w_m G(\sqrt{2}\sigma v_m | \theta, \xi) \sqrt{2}v_m \quad (13)$$

- 14 where w_m are the GH quadrature weights defined by

$$w_m = \frac{2^{M+1} M! \sqrt{\pi}}{[H'_M(v_m)]^2}, \quad m = 1, \dots, M, \quad (14)$$

- 15 v_m are the roots of the Hermite polynomial of degree M

$$H_M(v) = (-1)^M e^{v^2} \frac{d^M}{dv^M} (e^{-v^2}), \quad (15)$$

and M is the number of function evaluations, i.e., environment simulations, needed to
 compute the quadrature in Eq. (13). The weights $\{w_m\}_{m=1}^M$ and the roots $\{v_m\}_{m=1}^M$
 can be found in [1]. The approximation error of the GH formula can be bounded by
 $\tilde{\mathcal{D}}^M[G_\sigma]$ is

$$\left| \tilde{\mathcal{D}}^M[G_\sigma] - \mathcal{D}[G_\sigma] \right| \leq C_0 \frac{M! \sqrt{\pi}}{2^M (2M)!} \sigma^{2M-1}, \quad (16)$$

- 16 where $M!$ is the factorial of M , and the constant $C_0 > 0$ is independent of M and σ .

- 17 Applying the GH quadrature rule $\tilde{\mathcal{D}}^M$ to each component of $\nabla_{\sigma, \Xi}[J](\theta)$ in
 18 Eq. (11), we define the following estimator:

$$\tilde{\nabla}_{\sigma, \Xi}^M[J](\theta) := \Xi^\top \begin{bmatrix} \tilde{\mathcal{D}}^M[G_{\sigma_1}(0 | \theta, \xi_1)] \\ \vdots \\ \tilde{\mathcal{D}}^M[G_{\sigma_d}(0 | \theta, \xi_d)] \end{bmatrix}, \quad (17)$$

- 19 which requires a total of $M \times d$ *parallelizable* environment simulations at each
 20 iteration of training. The error bound in Eq. (16) indicates that $\tilde{\nabla}_{\sigma, \Xi}^M[J](\theta)$ is an
 21 accurate estimator of the DGS gradient for a small M , regardless of the value of σ .
 22 This enables the use of relatively big values of σ in Eq. (17) to exploit the *nonlocal*
 23 features of the landscape of $J(\theta)$ in the training process. The nonlocal exploitation
 24 ability of $\tilde{\nabla}_{\sigma, \Xi}^M[J]$ is demonstrated in §4 to be effective in reducing the necessary
 25 number of iterations to achieve a prescribed reward score. An illustration of the

- 1 nonlocal exploitation is given in Figure 1 in optimizing the Ackley function (its
 2 definition is given in Section 4.1).

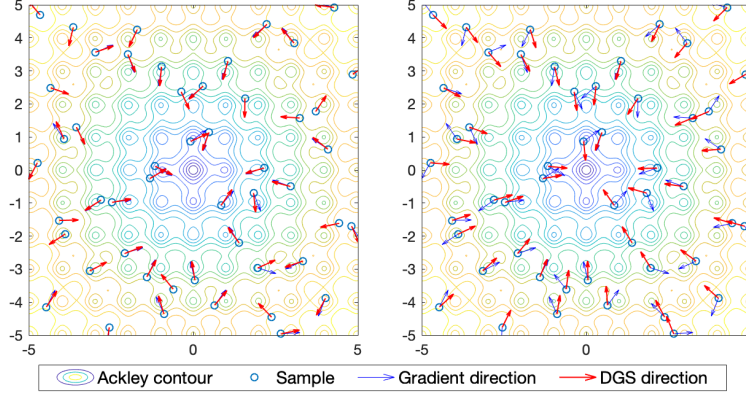


FIGURE 1. Illustration of DGS gradient direction and local gradient direction at 50 random locations on the surface of the Ackley in 2D. (Left) The standard deviation σ is set to 0.01, such that the DGS gradient align with the local gradient at most locations. (Right) The standard deviation σ is set to 1.0, such that most DGS gradient points to the global minimum at $(0, 0)$.

3 On the other hand, as the quadrature weights w_m and v_m defined in Eq. (14)
 4 and Eq. (15), are deterministic values, the DGS estimator $\hat{V}_{\sigma, \Xi}^M[J](\mathbf{x})$ in Eq. (17) is
 5 also a *deterministic* for fixed Ξ and σ . To introduce random exploration ability to
 6 our approach, we add random perturbations to both Ξ and σ . For the orthonor-
 7 mal matrix Ξ , we add a small random rotation, denoted by $\Delta\Xi$, to the current
 8 matrix Ξ . The matrix $\Delta\Xi$ is generated as a random skew-symmetric matrix, of
 9 which the magnitude of each entry is smaller than a prescribed threshold $\alpha > 0$.
 10 The perturbation of σ is conducted by drawing random samples from a uniform
 11 distribution $U(r - \beta, r + \beta)$ with two hyperparameters r and β with $r - \beta > 0$. The
 12 random perturbation of Ξ and σ can be triggered by various types of indicators
 13 e.g., the magnitude of the DGS gradient, the number of iteration done since last
 14 perturbation.

15 We suggest to set the hyperparameters in Algorithm 1 as follows: $M = 7, \alpha =$
 16 $2.0, r = 1.0, \beta = 0.2, \gamma = 0.01$. Note that both input and output variables need to
 17 be properly normalized before running the DGS-ES algorithm. In practice one can
 18 tune the critical hyperparameters (M, r and ℓ_r) given the following suggested range:
 19 $M \in [7, 9], r \in [0.5, 1.0]$ and $\ell_r \in [0.01, 0.1]$. Since our method does not have many
 20 hyperparameters, we suggest to use a simple grid search for hyperparameter tuning,
 21 and special attention should be spent to tuning the smoothing radius σ .

22 **4. Numerical experiments.** We evaluate the DGS-ES method using two sets of
 23 problems: (a) classic high-dimensional benchmark functions to test the performance
 24 of DGS-ES in solving high-dimensional multi-modal optimization problems, and (b)
 25 several reinforcement learning benchmark problems.

Algorithm 1: The DGS-ES for reinforcement learning

```

1: hyperparameters:
    $M$ : the order of GH quadrature rule
    $\alpha$ : the scaling factor for controlling the norm of  $\Delta\Xi$ 
    $r, \beta$ : the mean and perturbation scale for sampling  $\sigma$ 
    $\gamma$ : the tolerance for triggering random perturbation
2: Input:
    $\theta_0$ : the initial parameter value,
    $L$ : the number of parallel workers
3: Output: the final parameter value  $\theta_N$ 
4: Initialize the policy  $\pi$  with  $\theta_0$ 
5: Set  $\Xi = \mathbf{I}_d$ , and  $\sigma_i = r$  for  $i = 1, \dots, d$ 
6: Broadcast  $L$  copies of  $\pi(\mathbf{a}|\mathbf{s}; \theta_0)$  to the  $L$  workers
7: Divide the total GH quadrature points into  $L$  subsets, and send each subset to
   a worker.
8: for  $n = 0, \dots, N - 1$  do
9:   Each worker runs  $Md/L$  environment simulations at their assigned quadrature
   points
10:  Each worker sends  $Md/L$  scores to the master
11:  for  $i = 1, \dots, d$  do
12:    Compute  $\tilde{\mathcal{J}}^M[G_{\sigma_i}(0 | \theta_n, \xi_i)]$  in Eq. (13)
13:  end for
14:  Assemble  $\tilde{\nabla}_{\sigma, \Xi}^M[J](\theta_n)$  in Eq. (17)
15:  Update  $\theta_n$  to  $\theta_{n+1}$  using Adam
16:  if  $\|\tilde{\nabla}_{\sigma, \Xi}^M[J](\theta_n)\|_2 < \gamma$  then
17:    Generate  $\Delta\Xi$  and update  $\Xi = \mathbf{I}_d + \Delta\Xi$ 
18:    Generate  $\sigma$  from  $U(r - \beta, r + \beta)$ 
19:  end if
20:  Broadcast  $\theta_{n+1}$  to the  $L$  workers
21:  Each worker updates the policy to  $\pi(\mathbf{a}|\mathbf{s}; \theta_{n+1})$ 
22: end for

```

1 4.1. **Tests on high-dimensional benchmark functions.** Here we investigate
2 the DGS-ES performance on the high-dimensional functions: Sphere, Ackley, Lévy
3 and Rastrigin functions, whose definitions are given below.

- The Sphere function $F_1(\mathbf{x})$ is defined by

$$F_1(\mathbf{x}) = \sum_{i=1}^d x_i^2,$$

4 where d is the dimension and $\mathbf{x} \in [-5.12, 5.12]^d$ is the input domain. The global
5 minimum is $f(\mathbf{x}^*) = 0$ at $\mathbf{x}^* = (0, \dots, 0)$. It represents *convex and isotropic*
6 landscapes.

- The Ackley function $F_2(\mathbf{x})$ is defined by

$$F_2(\mathbf{x}) = -a \exp \left(-b \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2} \right) - \exp \left(\frac{1}{d} \sum_{i=1}^d \cos(cx_i) \right) + a + \exp(1),$$

where d is the dimension and $a = 20, b = 0.2, c = 2\pi$ are used in our experiments. The input domain $\mathbf{x} \in [-32.768, 32.768]^d$. The global minimum is $f(\mathbf{x}^*) = 0$, at $\mathbf{x}^* = (0, \dots, 0)$. The Ackley function represents *non-convex* landscapes with *nearly flat outer region*. The function poses a risk for optimization algorithms, particularly hill-climbing algorithms, to be trapped in one of its many local minima.

- The Rastrigin function $F_3(\mathbf{x})$ is defined by

$$F_3(\mathbf{x}) = 10d + \sum_{i=1}^d [x_i^2 - 10 \cos(2\pi x_i)], \quad (18)$$

where d is the dimension and $\mathbf{x} \in [-5.12, 5.12]^d$ is the input domain. The global minimum is $f(\mathbf{x}^*) = 0$ at $\mathbf{x}^* = (0, \dots, 0)$. This function represents *multimodal* and *separable* landscapes.

- The Lévy function $F_4(\mathbf{x})$ is defined by

$$F_4(\mathbf{x}) = \sin^2(\pi w_1) + \sum_{i=1}^{d-1} (w_i - 1)^2 [1 + 10 \sin^2(\pi w_i + 1)] + (w_d - 1)^2 [1 + \sin^2(2\pi w_d)],$$

where $w_i = 1 + (x_i - 1)/4$ for $i = 1, \dots, d$ and $\mathbf{x} \in [-10, 10]^d$ is the input domain. The global minimum is $f(\mathbf{x}^*) = 0$ at $\mathbf{x}^* = (0, \dots, 0)$. This represents *multi-modal* and *anisotropic* landscapes.

We compare the DGS-ES method with the following blackbox optimization methods: (a) Evolution Strategy (ES) in [23]; (b) Covariance matrix adaptation evolution strategy (CMA-ES) in [12] (we used the implementation of CMA-ES in the pycma open-source code from <https://github.com/CMA-ES/pycma>) and (c) the BFGS method in the Scipy library.

Figure 2 shows the scalability of the DGS-ES algorithm with the dimension of the objective function. We perform 20 repeated independent trials for each method to show the statistical performance. For the sphere function, the convergence rate does not change when we increase the dimension from 10 to 1000. Such property empirically carries over to the non-convex Lévy and Rastrigin functions. The reason is that both Lévy and Rastrigin have a globally near-convex structure when smoothing out their local minima. In contrast, the Ackley function is highly concave, as a large part of its surface is almost flat regardless of the small local minima. In this case, it takes many iterations for DGS-ES to search for the global minimum hidden in the middle of the flat surface.

Figure 3 shows the comparison between the DGS-ES and the baselines for optimizing the four benchmark functions in Figure 2 in 2000-dimensional space. For the sphere function, BFGS is a clear winner due to its optimal performance in handling convex functions. Among the three ES-type methods, the DGS-ES has much better performance than its competitors. For the other three functions, the DGS-ES method shows significant advantages over other methods.

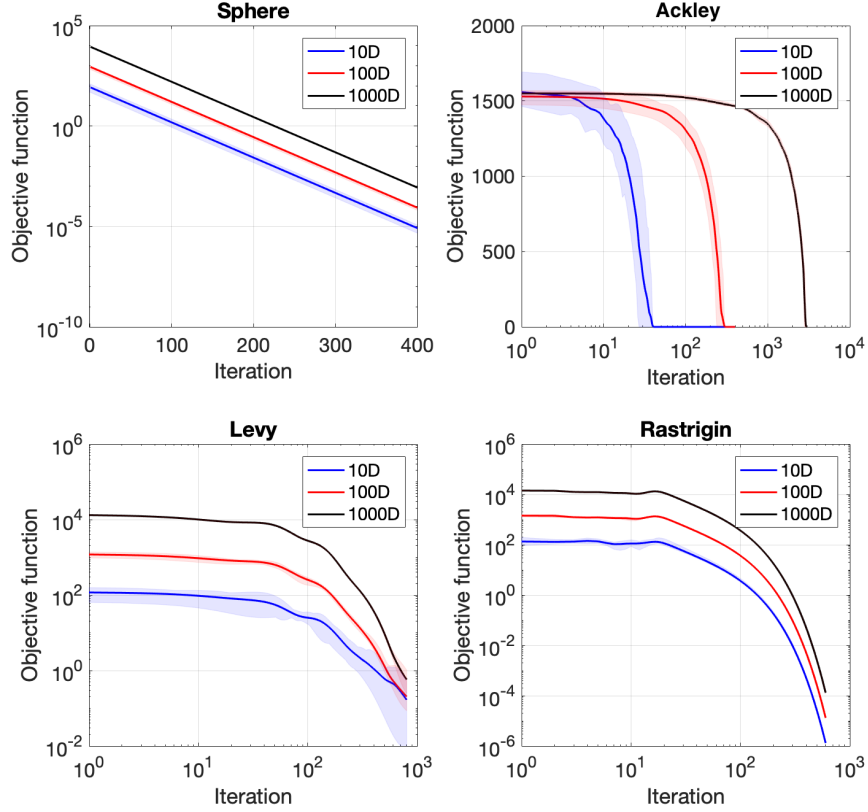


FIGURE 2. Illustration of the dimension dependence of the convergence rate of the DGS-ES method. The convergence rate, i.e., the number of iterations to converge, is *independent* of the dimension for convex functions, e.g., Sphere function, and such property empirically carries over to the non-convex Lévy and Rastrigin functions.

1 4.2. **Tests on RL benchmark problems.** To evaluate the DGS-ES algorithm,
2 we test its performance on two classes of reinforcement learning environments:
3 three classical control theory problems from OpenAI Gym ([https://github.com/](https://github.com/openai/gym)
4 [openai/gym](https://github.com/openai/gym)) [4] and three continuous control tasks simulated using PyBullet (2.6.5)
5 (<https://pybullet.org/>) [8] which is an open-source library. Within OpenAI Gym,
6 we demonstrate the proposed approach on three benchmark examples: `CartPole-v0`
7 (discrete), `MountainCarContinuous-v0` (continuous), `Pendulum-v0` (continuous). The
8 maximum time steps for these examples are 200, 999 and 200, respectively. More
9 details about the environment and reward settings can be found in [4]. We also ex-
10 amine the DGS-ES algorithm on the challenging continuous robotic control problems
11 in PyBullet library, namely `HopperBulletEnv-v0`, `InvertedPendulumBulletEnv-v0`
12 and `ReacherBulletEnv-v0`. In these three tasks, the maximum time steps are 1000,
13 1000 and 150, respectively. For the purpose of reproducible comparison, we employ
14 the original environment settings from the OpenAI Gym and the PyBullet library
15 without modifying the rewards or the environments.

16 For our implementation of DGS-ES, we define our policies as a two-layer feed-
17 forward neural network with 16 hidden nodes and tanh activation functions. For

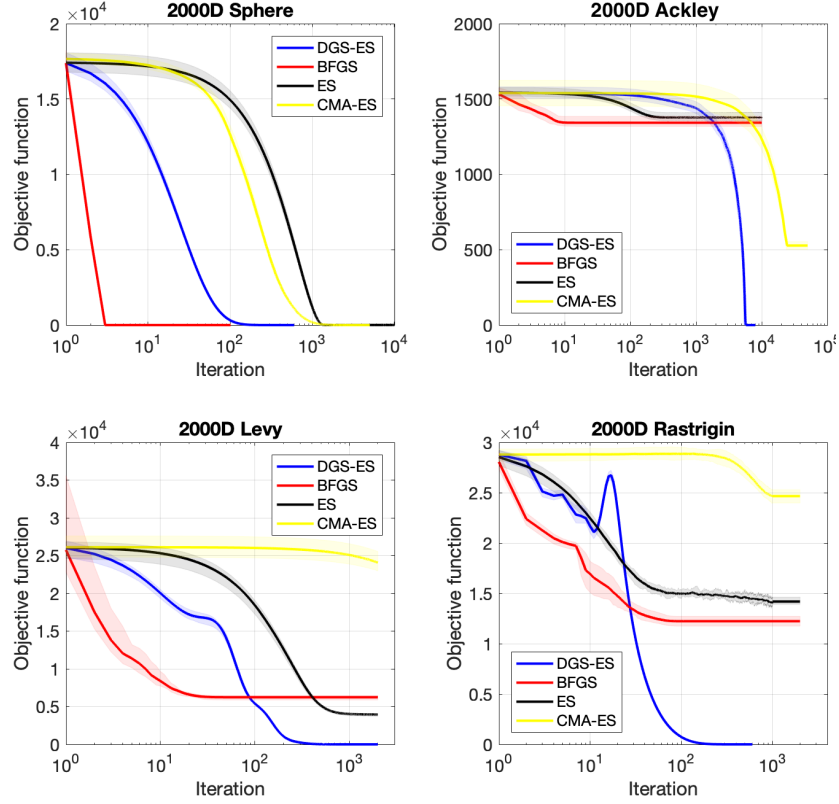


FIGURE 3. Comparison of different blackbox optimization methods on four 2000-dimensional benchmark functions.

1 gradient-based optimization, we use Adam to adaptively update the network pa-
 2 rameters with a learning rate of $\ell_r = 0.1$. We choose the hyperparameters used in
 3 Algorithm 1 as follows: $M = 7, \alpha = 2.0, r = 1.0, \beta = 0.2, \gamma = 0.01$. In practice one
 4 can tune the critical hyperparameters (M, r and ℓ_r) given the following suggested
 5 range: $M \in [7, 9]$, $r \in [0.5, 1.0]$ and $\ell_r \in [0.01, 0.1]$. For each task, our results
 6 are performed over 5 repeated independent trials (different random seeds) of the
 7 Gym/PyBullet simulators and the network policy initialization.

8 The DGS-ES algorithm is specifically amenable to parallelization since it only
 9 needs to communicate scalars, allowing it to scale to over a large number of parallel
 10 workers. We implement a distributed version of Algorithm 1 to the reinforcement
 11 learning tasks. The distributed DGS-ES is implemented using PyTorch [21] combined
 12 with Ray [19] (<https://github.com/ray-project/ray>), which does not rely on
 13 special networking setup and is tested on large-scale high performance computing
 14 facilities with thousands of computing nodes/workers.

15 *Comparison metric.* As the motivation of this work is to accelerate time-to-
 16 solution of reinforcement training under the assumption that sufficient distributed
 17 computing resource is available, we use a different metric to evaluate the performance
 18 of DGS-ES and the baselines. Specifically, we are interested in the average return
 19 $\mathbb{E}[J]$ versus the number of iterations, i.e., N in Algorithm 1, because those iterations
 20 cannot be parallelized.

1 *Baseline methods.* We compare Algorithm 1 against several RL baselines, including
 2 ES, PPO and TRPO, as well as the state-of-the-art algorithms such as ASEBO,
 3 DDPG, and TD3. Below is the information of the packages used.

- 4 • ES: The Evolution Strategy proposed in [23]. We used the implementation of
 5 ES from the open-source code <https://github.com/hardmaru/estool>.
- 6 • ASEBO: Adaptive ES-Active Subspaces for Blackbox Optimization, which was
 7 recently developed by [5]. We used the implementation released by the authors
 8 at <https://github.com/jparkerholder/ASEBO>.
- 9 • PPO: Proximal Policy Optimization in [24], which is available in OpenAI’s
 10 baselines repository at <https://github.com/openai/baselines> [9].
- 11 • TRPO: Trust Region Policy Optimization, developed by [25]. We also used the
 12 OpenAI’s baselines implementation [9].
- 13 • DDPG: Deep Deterministic Policy Gradient, proposed by [14]. We used the
 14 implementation from <https://github.com/georgesung/TD3> where the bench-
 15 mark DDPG in PyBullet is provided.
- 16 • TD3: Twin Delayed Deep Deterministic Policy Gradient [11], which was built
 17 upon the DDPG. The original results were reported for the MuJoCo version
 18 environments using the implementation from [https://github.com/sfujim/](https://github.com/sfujim/TD3)
 19 [TD3](https://github.com/sfujim/TD3), but we used the PyBullet implementation from [https://github.com/](https://github.com/georgesung/TD3)
 20 [georgesung/TD3](https://github.com/georgesung/TD3).

21 The hyperparameters for all algorithms above are set to match the original papers
 22 without further tuning to improve performance on the testing benchmark examples.

23 *Comparative evaluation.* Figure 4 shows the comparison results of CartPole,
 24 Pendulum and MountainCar problems from the OpenAI Gym. We compared the
 25 DGS-ES with classical ES and the improved ASEBO method. In general, the
 26 DGS-ES method features faster convergence than the baselines. For the simplest
 27 CartPole problem, the three methods perform equally well. Discrepancy appears
 28 in the Pendulum test, where the DGS-ES method not only converges faster than
 29 the baselines, but also achieves a higher average return. There is a much bigger
 30 discrepancy between the DGS-ES and the baselines appear in the MountainCar test.
 31 According to the guideline provided in the OpenAI Gym, the success threshold is to
 32 achieve an average return of 90. The DGS-ES method achieves the threshold within
 33 500 iterations, while the average returns of the ES and ASEBO methods are around
 34 zero. It is well known that the challenge of this problem is that the surface of the
 35 objective function $J(\theta)$ is very flat at most locations in the parameter space, which
 36 makes it difficult to capture the peak of $J(\theta)$. This test is a good demonstration
 37 of the nonlocal exploration ability of the DGS-ES method. Since the mean of the
 38 smoothing factor σ is set to 1.0, DGS-ES can capture the peak of $J(\theta)$ much faster
 39 than the baselines. In fact, we can see that it took around 50 iterations for the
 40 DGS-ES to find the peak region, while ES and ASEBO needed more than 3000
 41 iterations (not plotted) to move out of the flat region.

42 Figure 5 shows the comparison results of Hopper-v0, InvertedPendulum-v0 and
 43 Reacher-v0 problems from the PyBullet library. We compare the DGS-ES with six
 44 baselines including ES, ASEBO, PPO, TRPO, DDPG and TD3. As expected, the
 45 DGS-ES method shows better performance in terms of the convergence speed. For the
 46 Hopper-v0 problem, the DGS-ES achieves the highest return of all the methods within

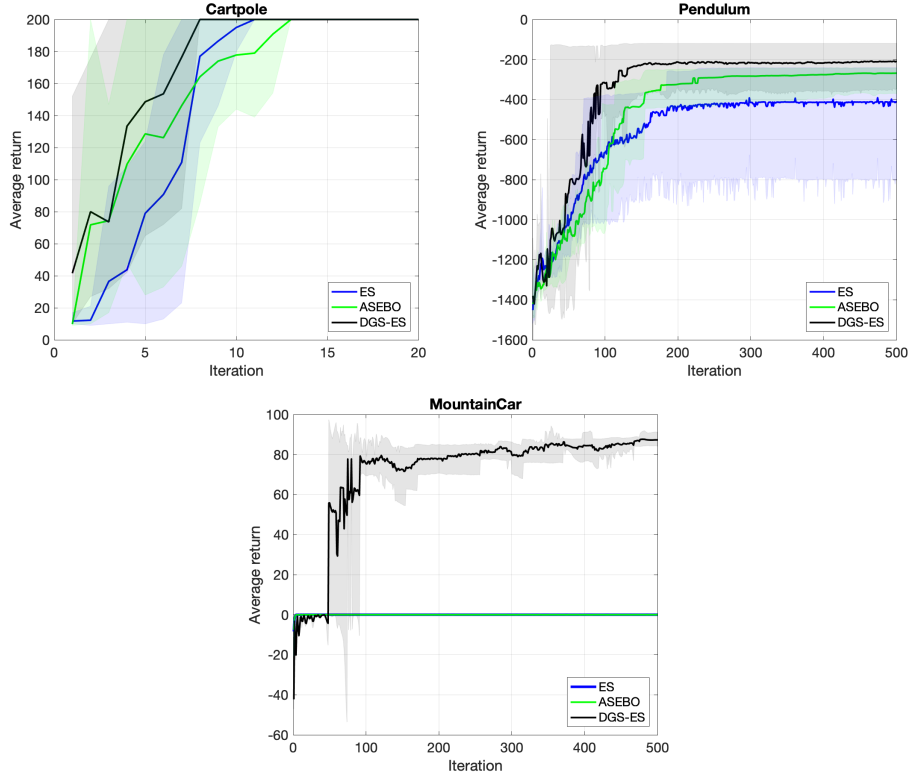


FIGURE 4. Comparison between the DGS-ES and two baselines, i.e., ES and ASEBO, for solving the three problems from OpenAI Gym. The colored curves are the average return over 5 repeated runs with different random seeds, and the corresponding shade represents the interval between the maximum and minimum return.

400 iterations. It is worth pointing out that some of the baselines will eventually catch up with the DGS-ES given a sufficiently large number of iterations. For example, DDPG and TD3 do not provide much improvement within 400 iterations, but DDPG could reach 1650 average return with over 3000 iterations, and TD3 could reach even higher, around 2200 average return, with over 6000 iterations, according to the baselines provided by OpenAI [9] and [11]. This phenomenon illustrates the fast convergence feature of DGS-ES. For the InvertedPendulum-v0 problem, DGS-ES can achieve the maximum return 1000 (default value in the PyBullet library) around 30 iterations. In comparison, ES and ASEBO can reach the maximum return but with more iterations than DGS-ES. According to the benchmark results for PyBullet environments in [11], DDPG and TD3 cannot converge to the maximum return even with a large number of iterations. For the Reacher-v0 problem, DGS-ES and ASEBO are still the top performers, and the advantage of DGS-ES is, again, faster convergence.

Figure 6 illustrates the effect of the radius σ of $\tilde{\nabla}_{\sigma, \Xi}^M[J](\theta)$ on the performance of the DGS-ES method in solving the Reacher-v0 problem. All the simulations are done using the same initialization. We set the mean of θ , i.e., the hyperparameter r

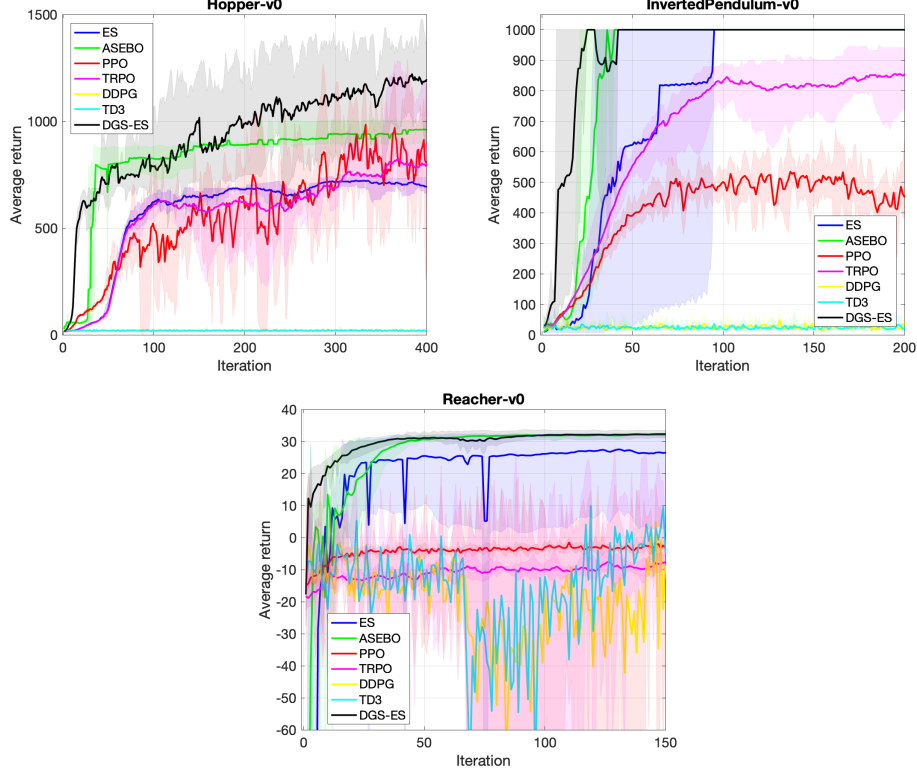


FIGURE 5. Comparison between the DGS-ES method and the baselines, i.e., ES, ASEBO, PPO, TRPO, DDPG and TD3, for solving the three problems from the PyBullet library. The colored curves are the average return over 20 runs with different random seeds, and the corresponding shade represents the interval between the maximal and minimal rewards.

1 in Algorithm 1, to 0.5, 0.05, 0.01, and 0.005. It is easy to see that the performance
2 of DGS-ES deteriorates with the decrease of σ . Since it is evident that the surface
3 of $J(\theta)$ is not convex and may have many local maxima, a relatively big radius σ
4 is necessary to help DGS-ES skip the local maxima. As σ becomes smaller, the
5 DGS gradient $\hat{\nabla}_{\sigma, \Xi}^M[J](\theta)$ converges to the local gradient $\nabla J(\theta)$, which may get the
6 optimizer trapped in a local maximum.

7 **4.3. Aerospace hierarchical stiffened shell design.** Now we demonstrate the
8 DGS-ES method on a real-world stiffened shell design problem. Due to its high
9 strength and stiffness, the hierarchical stiffened shell has been widely used in
10 aerospace engineering [31, 32]. However, it is challenging to fully explore its optimal
11 buckling load-carrying capacity. The goal of this work is to improve the load-carrying
12 capacity by optimizing the representative unit cell of hierarchical stiffened shell where
13 the inputs are 9 size variables (widths, heights and thickness) for major and minor
14 stiffeners, as shown in Figure 7(a), and the output is the carrying-load capability
15 factor which is calculated by high-fidelity numerical simulation, e.g., finite element
16 method (FEM), see Figure 7 (b). Although the high-fidelity FEM simulation is

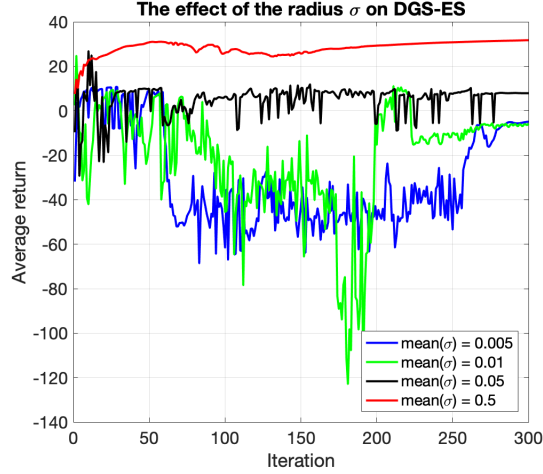


FIGURE 6. Illustration of the effect of the radius σ of $\tilde{\nabla}_{\sigma, \Xi}^M[J](\theta)$ on the performance of the DGS-ES method in solving the Reacher-v0 problem.

1 time-consuming, many open-source codes and commercial software have improved
 2 the computational efficiency via scalable parallelism and GPU acceleration. It is thus
 3 feasible to apply the DGS-ES method by implementing parallel FEM simulations in
 4 supercomputers. Figure 7 presents the comparison between DGS-ES and the other
 5 blackbox optimization methods. We observe that DGS-ES outperforms all other
 6 algorithms and achieves faster convergence using only 100 iterations. The robustness
 7 of DGS-ES also outperforms the alternatives.

8 **5. Conclusion.** Despite the successful demonstration shown in §4, there are several
 9 limitations with the DGS-ES method for reinforcement learning. First, it requires
 10 a powerful enough cluster or distributed computing resources to show superior
 11 performance. Even though more and more parallel environments for complex RL
 12 tasks, those parallel codes might not compatible with a cluster/supercomputer with
 13 a specific architecture. This will need some extra effort to modify environment codes,
 14 in order to exploit the advantage of the DGS-ES approach. Second, asynchronization
 15 between different environment simulations may drag down the total performance.
 16 In a distributed computing system, all the parallel workers receive the same number
 17 of environment simulations. However, as different parameter values may lead to
 18 different termination times of the environment simulations, there will be a potential
 19 waste of computing resources due to such asynchronization. Thus, a better scheduling
 20 algorithm is needed to further improve the performance of DGS-ES in RL. Third, even
 21 though the performance of the DGS-ES is not very sensitive to the hyperparameters,
 22 especially the radius σ in the experiments conducted in this work, optimal or even
 23 viable hyperparameters of the DGS-ES method are still problem dependent, which
 24 means hyperparameter tuning may be needed when applying the method to another
 25 RL problem.

26 **Acknowledgement.** This material was based upon work supported by the U.S.
 27 Department of Energy, Office of Science, Office of Advanced Scientific Computing Re-
 28 search, Applied Mathematics program under contract and award numbers ERKJ352,

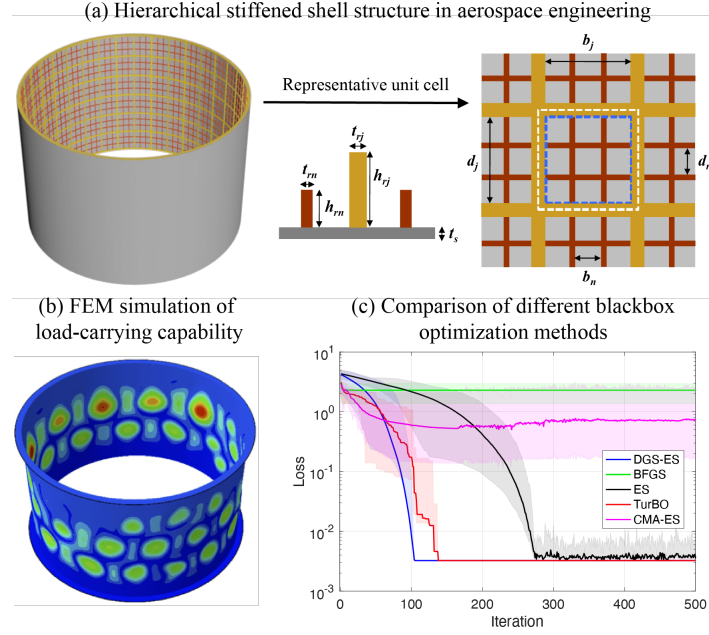


FIGURE 7. Illustration and design optimization of hierarchical stiffened shell structures in aerospace engineering, e.g. rocket.

1 by the DOE SciDac FastMath project, and by the Artificial Intelligence Initiative at
 2 the Oak Ridge National Laboratory (ORNL). ORNL is operated by UT-Battelle,
 3 LLC., for the U.S. Department of Energy under Contract DE-AC05-00OR22725.

4 REFERENCES.

- 5 [1] M. Abramowitz and I. Stegun, editors. *Handbook of Mathematical Functions*.
 6 Dover, New York, 1972.
- 7 [2] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong,
 8 Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech
 9 Zaremba. Hindsight experience replay. In *Advances in Neural Information*
 10 *Processing Systems*, pages 5048–5058, 2017.
- 11 [3] A. S. Berahas, L. Cao, K. Choromanskiv, and K. Scheinberg. A theoretical and
 12 empirical comparison of gradient approximations in derivative-free optimization.
 13 *arXiv:1905.01332*, 2019.
- 14 [4] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John
 15 Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint*
 16 *arXiv:1606.01540*, 2016.
- 17 [5] Krzysztof Choromanski, Aldo Pacchiano, Jack Parker-Holder, , and Yunhao
 18 Tang. From complexity to simplicity: Adaptive es-active subspaces for blackbox
 19 optimization. *NeurIPS*, 2019.
- 20 [6] Krzysztof Choromanski, Aldo Pacchiano, Jack Parker-Holder, , and Yun-
 21 hao Tang. Provably robust blackbox optimization for reinforcement learning.
 22 *arXiv:1903.02993*, 2019.
- 23 [7] Edoardo Conti, Vashisht Madhavan, Felipe Petroski Such, Joel Lehman, Ken-
 24 neth O. Stanley, and Jeff Clune. Improving exploration in evolution strategies

- 1 for deep reinforcement learning via a population of novelty-seeking agents. *NIPS*,
2 2018.
- 3 [8] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simula-
4 tion for games, robotics and machine learning. *GitHub repository*, 2016.
- 5 [9] Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias
6 Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter
7 Zhokhov. Openai baselines. <https://github.com/openai/baselines>, 2017.
- 8 [10] Abraham D. Flaxman, Adam Tauman Kalai, Adam Tauman Kalai, and H. Bren-
9 dan McMahan. Online convex optimization in the bandit setting: gradient
10 descent without a gradient. *Proceedings of the 16th Annual ACM-SIAM sympo-*
11 *sium on Discrete Algorithms*, pages 385–394, 2005.
- 12 [11] Scott Fujimoto, Herke Van Hoof, and David Meger. Addressing function
13 approximation error in actor-critic methods. *arXiv preprint arXiv:1802.09477*,
14 2018.
- 15 [12] Nikolaus Hansen. The cma evolution strategy: a comparing review. In *Towards*
16 *a new evolutionary computation*, pages 75–102. Springer, 2006.
- 17 [13] Nikolaus Hansen and Andreas Ostermeier. Completely derandomized self-
18 adaptation in evolution strategies. *Evolutionary computation*, 9(2):159–195,
19 2001.
- 20 [14] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver,
21 and D. Wierstra. Continuous control with deep reinforcement learning. *ICLR*,
22 2016.
- 23 [15] Niru Maheswaranathan, Luke Metz, George Tucker, Dami Choi, and Jascha
24 Sohl-Dickstein. Guided evolutionary strategies: Augmenting random search
25 with surrogate gradients. *Proceedings of the 36th International Conference on*
26 *Machine Learning*, 2019.
- 27 [16] Florian Meier, Asier Mujika, Marcelo Matheus Gaury, and Angelika Steger.
28 Improving gradient estimation in evolutionary strategies with past descent
29 directions. *Optimization Foundations for Reinforcement Learning Workshop at*
30 *NeurIPS*, 2019.
- 31 [17] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timo-
32 thy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous
33 methods for deep reinforcement learning. *ICML*, pages 1928–1937, 2016.
- 34 [18] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness,
35 Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg
36 Ostrovski, and et al. Human-level control through deep reinforcement learning.
37 *Nature*, 518(7540):529–533, 2015.
- 38 [19] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard
39 Liaw, Eric Liang, Melih Elibol, Zongheng Yang, William Paul, Michael I
40 Jordan, et al. Ray: A distributed framework for emerging {AI} applications. In
41 *13th {USENIX} Symposium on Operating Systems Design and Implementation*
42 (*{OSDI} 18*), pages 561–577, 2018.
- 43 [20] Y. Nesterov and V. Spokoiny. Random gradient-free minimization of convex
44 functions. *Foundations of Computational Mathematics*, pages 1–40, 2015.
- 45 [21] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang,
46 Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer.
47 Automatic differentiation in pytorch. 2017.
- 48 [22] Alfio Quarteroni, Riccardo Sacco, and Fausto Saleri. *Numerical Mathematics*,
49 volume 332. Springer Science Business Media &, 2007.

- 1 [23] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. From
2 complexity to simplicity as a scalable alternative to reinforcement learning.
3 *arXiv preprint arXiv:1703.03864*, 2017.
- 4 [24] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal
5 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 6 [25] John Schulman, Sergey Levine, Pieter Abbeel, Michael I Jordan, and Philipp
7 Moritz. Trust region policy optimization. *ICML*, pages 1889–1897, 2015.
- 8 [26] F. Sehnke, C. Osendorfer, T. Ruckstieb, A. Graves, J. Peters, and J. Schmidhuber.
9 Parameter-exploring policy gradients. *Neural Networks*, 2010.
- 10 [27] Olivier Sigaud and Freek Stulp. Robot skill learning: From reinforcement
11 learning to evolution strategies. *Paladyn Journal of Behavioral Robotics*, 4(1):49–
12 61, 2013.
- 13 [28] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche,
14 J. Schrittwieser, I. Antonoglou, V. Panneershelvam, and et al. M. Lanctot.
15 Mastering the game of go with deep neural networks and tree search. *Nature*,
16 529(7587):484–489, 2016.
- 17 [29] F. P. Such, V. Madhavan, E. Conti, J. Lehman, K. O. Stanley, and J. Clune.
18 Deep neuroevolution: Genetic algorithms are a competitive alternative for
19 training deep neural networks for reinforcement learning. *arXiv preprint*
20 *arXiv:1712.06567*, 2017.
- 21 [30] R.S. Sutton and A. G. Barto, editors. *Reinforcement learning: An introduction*,
22 volume 1. MIT press Cambridge, New York, 1998.
- 23 [31] Kuo Tian, Lei Huang, Yu Sun, Kaifan Du, Peng Hao, and Bo Wang. Fast
24 buckling load numerical prediction method for imperfect shells under axial
25 compression based on pod and vibration correlation technique. *Composite*
26 *Structures*, 252:112721, 2020.
- 27 [32] Kuo Tian, Li Zengcong, Lei Huang, K. Du, Liangliang Jiang, and Bo Wang.
28 Enhanced variable-fidelity surrogate-based optimization framework by gaussian
29 process regression and fuzzy clustering. *Computer Methods in Applied Mechanics*
30 *and Engineering*, 366:113045, 2020.
- 31 [33] Jiaxin Zhang, Hoang Tran, Dan Lu, and Guannan Zhang. Enabling long-
32 range exploration in minimization of multimodal functions. *Proceedings of 37th*
33 *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2021.
- 34 E-mail address: zhangj@ornl.gov
- 35 E-mail address: tranha@ornl.gov
- 36 E-mail address: zhangg@ornl.gov