# A Surrogate-Assisted Controller for Expensive Evolutionary Reinforcement Learning

Yuxing Wang[a], Tiantian Zhang[a], Yongzhe Chang[a], Bin Liang[b], Xueqian Wang[a] and Bo Yuan[a,*]

[a]*Shenzhen International Graduate School, Tsinghua University, Shenzhen, 518055, China*
[b]*Department of Automation, Tsinghua University, Beijing, 100084, China*

## ARTICLE INFO

## ABSTRACT

The integration of Reinforcement Learning (RL) and Evolutionary Algorithms (EAs) aims at simultaneously exploiting the sample efficiency as well as the diversity and robustness of the two paradigms. Recently, hybrid learning frameworks based on this principle have achieved great success in various challenging tasks. However, in these methods, policies from the genetic population are evaluated via interactions with the real environments, limiting their applicability when such interactions are prohibitively costly. In this work, we propose Surrogate-assisted Controller (SC), a novel and efficient module that can be integrated into existing frameworks to alleviate the burden of EAs by partially replacing the expensive fitness evaluation. The key challenge in applying this module is to prevent the optimization process from being misled by the possible false minima introduced by the surrogate. To address this issue, we present two strategies for SC to control the workflow of hybrid frameworks. Experiments on six continuous control tasks from the OpenAI Gym platform show that SC can not only significantly reduce the cost of interacting with the environment, but also boost the performance of the original hybrid frameworks with collaborative learning and evolutionary processes.

## 1. Introduction

Reinforcement Learning (RL) has demonstrated promising achievements in various domains, ranging from Atari games [26], GO [36], to robot control tasks [20]. Among these successes, Deep Learning (DL) techniques [11] such as Deep Neural Networks (DNNs) have been widely used for decision-making [40, 47]. The combination of RL with DL is generally called Deep Reinforcement Learning (DRL). However, recent studies show that DRL suffers from premature convergence to local optima in the training process, and the RL agents are highly sensitive to hyperparameter settings, implementation details and uncertainties of the environmental dynamics [14], preventing agents from learning stable policies.

In the meantime, black-box optimization techniques such as Evolutionary Algorithms (EAs) [50, 8] have shown competitive results compared to DRL algorithms [33, 48]. Firstly, the population-based mechanism makes EAs explore the parameter space better than DRL. Secondly, since EAs only consider the total returns across the entire episode, they are indifferent to the issue of sparse reward and robust to environmental noise [33]. However, EAs suffer from low sample efficiency [18], due to their inherent black-box properties, and they often do not make full use of the feedback signals and historical data from the environment.

An emergent research direction is dedicated to exploiting the benefits of both solutions following the theory of evolution [45], where the learning of individuals can increase the evolutionary advantage of species, which can subsequently make the population learn faster. Recently proposed



**Figure 1:** An overview of Surrogate-assisted Controller (SC). The SC switches between the real fitness function and the approximated fitness function constructed by the surrogate model to efficiently evaluate the genetic population in hybrid frameworks, while boosting the performance of the original hybrid frameworks. As shown at the bottom, in some cases, inaccurate surrogates may bring extra benefits: the predicted fitness function (dashed curve) has large deviations from the real one (solid curve) but nevertheless features the same global minimum and constructs a new fitness landscape that is more friendly to evolutionary search.

Evolutionary Reinforcement Learning (ERL) [18], Proximal Distilled ERL (PDERL) [2] and other hybrid frameworks based on EA and off-policy DRL [29, 23, 49, 22] have demonstrated encouraging progresses on single-objective hard-exploration tasks with large continuous state and action spaces [3], outperforming pure DRL and EA. Particularly, one common feature of these approaches is that a diverse genetic population of policies is generated by EA to drive the

exploration, while an off-policy actor-critic algorithm thoroughly exploits the population's environmental experience to search for high-performing policies.

However, there is a notable issue with these hybrid frameworks: the evaluation of the genetic population requires all individuals to be presented to the environment during the training process. Although population-based interactions can provide diverse historical experiences, they usually require a large amount of computational time and can be prohibitively expensive [34, 16]. For instance, in robotic scenarios, frequent evaluations may lead to massive resources consumption or even equipment damage. These unfavorable characteristics greatly limit the applicability of existing hybrid frameworks in complex simulated circumstances or non-trivial real-world scenarios.

A possible solution is to introduce the mechanism of surrogate. In typical Surrogate-assisted Evolutionary Optimization (SEO), surrogates, also known as meta-models, are trained to partially replace the expensive fitness functions [42, 28]. Surrogate models such as Kriging [7], polynomials [17] and neural networks have been successfully applied to domains including constrained optimization [1] and multi-objective optimization [21, 44]. Nevertheless, in the standard RL context, due to the large uncertainties of genotype-phenotype-fitness mapping [38], conventional methods face great challenges of high computational cost and modeling complexity (Section 2.2) with only very limited studies conducted on simple discrete control tasks.

In this paper, we design a generic and effective module, called Surrogate-assisted Controller (SC), that can be conveniently combined with existing hybrid RL frameworks to improve their practicability. As shown in Figure 1, SC employs an approximated fitness function together with the real fitness function to help evaluate the genetic population. Its sample-efficient surrogate model can be naturally implemented in the existing hybrid RL frameworks, making full use of the diverse experiences to evaluate the fitness of individuals without environmental interactions. Note that the primary criterion for applying the surrogate model is the introduced prediction error [32]. SC mitigates this risk using two management strategies with an elite protection mechanism (Section 4.2), which can strategically schedule the real and the surrogate-assisted evaluation as well as effectively prevent the dramatic fluctuation of performance and the spread of detrimental information from individuals with inaccurate fitness values. Our empirical studies in Section 5 show that SC can not only successfully transfers the computational burden of real fitness evaluations to the efficient surrogate model and boost the performance of original hybrid approaches, but also stabilize the interactions between the RL agent and the genetic population.

The major contributions of our work are summarized as follows:

- A novel module named Surrogate-assisted Controller (SC) is proposed that can be easily integrated into existing hybrid RL frameworks to significantly relieve

the computational cost of interacting with the real environment during the optimization process.

- Two effective management strategies are presented for SC to control the workflow of the hybrid frameworks, which can strategically schedule the surrogate-assisted evaluation and the real fitness evaluation.

- We combine SC with ERL and PDERL, creating two new frameworks named SERL and SPDERL, respectively, to highlight its principle and flexibility. Comprehensive experimental studies on Mujoco benchmarks show that SC can not only effectively bring better sample-efficiency to the original hybrid frameworks, but also stabilize the learning and evolution processes with superior performance.

The remainder of this paper is organized as follows. Section 2 introduces the related work on hybrid frameworks and surrogate-assisted methods for solving RL problems. The problem definition is specified in Section 3 and the details of SC and its components are presented in Section 4. In Section 5, the numerical validation and in-depth analyses are conducted and Section 6 concludes our work with some discussions on future research directions.

## 2. Related work

### 2.1. Combining EA and off-policy DRL

As an alternative solution for RL problems [33], EAs have also been combined with DRL to leverage the benefits of both solutions. The first hybrid framework ERL [18] combines an off-policy DRL agent based on DDPG [20] with a population evolved by the Genetic Algorithm (GA) [46]. In each generation, individuals are evaluated over a few episodes and the fitness is given by averaging the total rewards. Based on this information, policy optimization is then conducted over the parameter space by genetic operators. Meanwhile, the RL agent is trained on the diverse experiences produced by the population. In the periodical synchronization step, it is injected into the population and this bi-directional interaction controls the information flow between the RL agent and the genetic population. To further extend the ERL framework, PDERL [2] uses the distillation crossover operator to alleviate the catastrophic forgetting caused by the recombination of neural networks. In addition, various evolutionary methods [29, 23, 39] have also been combined with other DRL frameworks such as TD3 [10] and SAC [13] to further leverage the advantages of both gradient-based and gradient-free methods. In our paper, we mainly focus on the paradigm of ERL and PDERL to show the effectiveness of our proposed methods.

### 2.2. Surrogate-assisted methods for RL problems

Surrogate-assisted methods have been widely investigated to reduce unnecessary expensive evaluations [15]. Recently, a few studies on using the surrogate model to enhance

RL algorithms or EAs in the context of sequential decision-making tasks have been conducted and can be generally divided into two categories.

The first category focuses on learning a model of the environment. Many studies [4, 12, 6] try to construct a transition model, which is trained in a standard supervised manner using a large amount of historical data, to imitate at least some aspects of the system's physical dynamics. Then, the fixed surrogate model is embedded into RL algorithms to help learn a control strategy from the experiences obtained by continuous interactions with the surrogate model. However, for domains with high noise or when the availability of the historical data is limited, this kind of methods may no longer be effective.

The other category does not require modeling the environment and mainly focuses on evaluation. For example, Kriging can be employed to directly map the relationships between neural networks (policies) and their fitness. Typically, genotypic distances between neural networks are needed by the Kriging model [37]. However, the large-scale settings and complex problems may make the computation of these distances practically impossible. Although approximate distances such as the phenotypic distance [38] can offer some help, the high dimensionality of the inputs can result in an increase of computational costs, and the parameter settings for Kriging remain a challenge for input vectors of different sizes [38].

For instance, Evolutionary Surrogate-assisted Prescription (ESP) [9] incorporates a surrogate model into the EA. Given a set of input states $\mathcal{S}$, the policy neural network takes each $s \in \mathcal{S}$ as the input and outputs the action $a$. The surrogate model of ESP, represented by a random forest or a deep neural network, is used to predict the outcome of each state-action pair $(s, a)$, and the fitness of each policy is given by averaging the outcomes over $\mathcal{S}$. In each generation, the surrogate is first trained on the historical data by minimizing the Mean Square Error (MSE) loss between the real and the predicted outcomes. Subsequently, individuals (called *prescriptors*) in the population are evolved with the trained surrogate. Finally, selected elites are presented to the real environment to generate new training data for the surrogate. Unfortunately, for challenging environments with large continuous state and action spaces or with rich feedback signals, the lack of gradient information makes EAs suffer from brittle convergence [33, 18]. Furthermore, the update of the policy neural networks in ESP is purely based on predicted fitness, which increases the risk of misleading the evolutionary optimization in the wrong direction under complex circumstances. Consequently, the frequency of applying the surrogate also needs to be managed properly to ensure the stability of the training process.

In general, apart from limited successes on simple discrete control tasks such as Cart-Pole [3], conventional surrogate-assisted EAs [38, 37, 9] still face significant challenges. To explore the potential of surrogate-assisted methods in continuous and complex RL contexts, in this work, we extend the surrogate model to hybrid RL frameworks with the objective to reduce the cost of evaluations while making full use of the efficiency of the RL agent and the exploration capability of the EA.

## 3. Preliminaries

In DRL, each problem is modeled as a Markov Decision Process (MDP), which can be specified by a 5-tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma \rangle$. With the state space $\mathcal{S}$ and the action space $\mathcal{A}, \mathcal{P} : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the transition function of the environment; $r(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, and the discount factor $\gamma \in (0, 1]$ specifies the degree to which rewards are discounted over time.

$$Q(s, a|\theta^Q) = \mathbb{E}\left[ \sum_{i=0}^{\infty} \gamma^i r_{t+i+1} \Big| s_t = s, a_t = a \right] \quad (1)$$

The actor-critic architecture based on the policy gradient approach is widely used in DRL [40]. The critic network $Q(s, a|\theta^Q)$, as shown in Eq.(1), is used to estimate the expectation of the discounted accumulative rewards of the state-action pair $(s, a)$. According to the Bellman equation, its recursive expression is:

$$Q(s, a|\theta^Q) = \mathbb{E}\left[ r(s, a) + \gamma Q(s', a'|\theta^Q) \right] \quad (2)$$

where $s'$ and $a'$ represent the next state and action, respectively. In each iteration, transitions with the batch size of $N$ are sampled randomly from the experience replay buffer to update the parameters of $Q(s, a|\theta^Q)$ by minimizing the Temporal Difference (TD) loss based on the standard back-propagation:

$$\mathcal{L}_{Q(s,a|\theta^Q)} = \frac{1}{N} \sum_i \left( y_i - Q(s_i, a_i|\theta^Q) \right)^2 \quad (3)$$

$$y_i = r(s_i, a_i) + \gamma Q\left( s', \mu(s'|\theta^\mu)|\theta^Q \right) \quad (4)$$

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i Q\left( s_i, \mu(s_i|\theta^\mu) \right) \nabla_{\theta^\mu} \mu(s_i|\theta^\mu) \quad (5)$$

Then, the critic is used to assist the training of the actor network $\mu(s|\theta^\mu)$ via policy gradient, according to Eq.(5). With the back-propagation method, the parameter $\theta^\mu$ of the actor network is updated in the direction towards maximizing the $Q$ value.

## 4. Methodology

In this section, we present the Surrogate-assisted Controller (SC) with two management strategies and introduce how to incorporate it into the hybrid framework. Figure 2 shows an example of combining SC with ERL, where the

**Figure 2:** Integration of SC and ERL. The ERL framework is shown in the light gray box and SC is shown in the light green box. To evaluate the population without environment interactions, the recent state information is sampled from the replay buffer as the evaluation memory. After that, the RL-Critic plays the role of a surrogate model to evaluate actors based on the historical data. Finally, the predicted population fitness is consumed by evolutionary methods.

---

**Algorithm 1** Surrogate-assisted Evaluation

---

**Input**: Policy set of population $\mu_{pop} = \{\mu_1, \mu_2, ..., \mu_n\}$;
      The surrogate model $Q(s, a|\theta^Q)$;
      Replay buffer $\mathcal{R}$ for RL agent training.
**Output**: Population fitness $F$
1: Sample $k$ latest states from $\mathcal{R}$ to the evaluation memory $\mathcal{R}_{eva}$ as the evaluation samples;
2: **for** $i = 1$ to $n$ **do**
3:    Initialize the fitness of $\mu_i$: $f_i = 0$;
4:    **for** $j = 1$ to $k$ **do**
5:       $f_i = f_i + Q\big(s_j, \mu_i(s_j|\theta^{\mu_i})|\theta^Q\big)/k$;
6:    **end for**
7: **end for**
8: **return** Population fitness $F = \{f_1, f_2, ..., f_n\}$.

---

RL-Critic exploits the diverse experiences collected by the genetic population to simultaneously train the RL-Actor and assist the evaluation, referred to as Surrogate-assisted Evolutionary Reinforcement Learning (SERL). In practice, SC includes three components: critic-based surrogate model, management strategies and evaluation memory.



**Figure 3:** The preliminary experiment on the surrogate's approximation accuracy over two environments: Hopper (left, trained for 1 million steps) and Walker (right, trained for 3 million steps). The curves of approximation accuracy clearly show that the accuracy of the surrogate tends to increase during the training process.

## 4.1. Critic-based surrogate model

Using the real fitness function can be time-consuming or even dangerous for expensive problems. A surrogate, which is a predictive model, can be used to assist the evaluation. Note that, in the domain of Evolutionary Computation (EC), the fitness of an individual can be determined in many forms, even as a non-markovian definition. In our work, we consider the real fitness value as the sum of the reward over the episode, following the typical setting of RL paradigm.

However, the genotype-phenotype-fitness mapping of sequential decision-making problems is often difficult to learn, but with the definition of fitness mentioned above, it is straightforward for the surrogate model to estimate the outcome of a state-action pair $(s, a)$, as shown in Eq.(1). In ESP [9], individuals are evaluated by an extra dedicated surrogate model. However, in hybrid frameworks, the critic module $Q(s, a|\theta^Q)$ of the off-policy RL agent can be naturally used as a surrogate. The actor module $\mu(s|\theta^\mu)$ represented by a fully connected policy neural network takes a state vector $s$ as the input and decides what action vector $a$ to perform. Then, the concatenated vector $(s, a)$ is evaluated by the critic-based surrogate model. Thus, with the evaluation memory (described in Section 4.3) that contains the information of the $k$ latest state vectors drawn from the replay buffer, the fitness value of an actor in the genetic population can be obtained by averaging its predicted $Q$ values over all states. Then we have:

$$f = \frac{1}{k} \sum_{j=1}^{k} Q\big(s_j, \mu(s_j|\theta^\mu)|\theta^Q\big) \qquad (6)$$

The complete process of the surrogate-assisted evaluation is presented in Algorithm 1, and its intrinsic motivation is to evaluate a policy whether it is powerful and robust enough to perform as much as high-quality actions when facing diverse states. Furthermore, a prominent feature is that there is no extra cost involved in training the surrogate model, as it is part of the standard training procedure of the RL. Different from conventional actor-critic methods such as Asynchronous Advantage Actor-Critic (A3C) [24], where the critic is only used to guide the improvement of policies, the critic-based surrogate model is able to simultaneously

---

**Figure 4:** Two kinds of management strategies: the Generation-based Control (GC) and the Individual-based Control (IC). For GC, the parameter update direction is partially based on real fitness information, while the update direction is totally provided by the real fitness of the pre-selected candidates in IC.

train the RL-Actor via policy gradient and provide relatively accurate fitness estimations for individuals with various improvements in the estimation of $Q$ values [10, 43].

The approximation quality of the surrogate is measured by the Spearman's rank correlation coefficient $r_s \in [-1, 1]$ between the real and the predicted fitness values, where $n$ is the number of data points (population size) and $d$ represents the difference between the two ranks:

$$r_s = 1 - \frac{6\sum_{i=1}^{n} d_i^2}{n(n^2 - 1)} \quad (7)$$

A preliminary experiment on two standard continuous control tasks from MuJoCo [41] is conducted to show how the approximation accuracy of the critic-based surrogate changes during ERL's training process. In each generation, we apply the surrogate-assisted evaluation to the newly generated population with $k = 50,000$ in the first place. After evaluating the population in the real environment, $r_s$ is calculated to report the current approximation accuracy of the surrogate. The result in Figure 3 is reported over 6 runs, it suggests that the approximation accuracy increase during the training process. Although at the beginning of training, the surrogate has relatively low approximation accuracy, previous studies [15, 9] suggest that this kind of uncertainty may not cause negative effects on evolutionary search. Instead, it may further push the population to explore the fitness landscape. Moreover, the training data generated by these candidates can subsequently help improve the approximation accuracy of the surrogate.

## 4.2. Management strategies

Although the surrogate model can provide an approximately accurate fitness estimation, using predicted fitness to assist the evolutionary operations throughout the optimization process may easily introduce false minima, leading to a drop in performance or an increase in the computational cost of evaluating low-performing policies. As a result, surrogates should be used along with real fitness functions and we consider the following two methods: generation-based and individual-based strategies.

---

**Algorithm 2** Generation-based control

**Input**: Policy set of population $\mu_{pop} = \{\mu_1, \mu_2, ..., \mu_n\}$;
      The surrogate model $Q_{rl}$;
      Replay buffer $\mathcal{R}$ for RL agent training;
      Control factor $\omega$;
      Random number generator $G_{rand} \in (0, 1]$.
**Output**: New population $\mu_{pop*}$

 1: **if** $G_{rand} > \omega$ **then**
 2:     $F_{real} = \text{Evaluation}(\mu_{pop})$
 3:     Copy the elite actor
 4:     $\mu_{pop*} = \text{Evolutionary methods}(F_{real}, \mu_{pop})$
 5: **else**
 6:     $F_{pre} = \text{Surrogate-assisted evaluation}(\mu_{pop}, Q_{rl}, \mathcal{R})$
 7:     $\mu_{pop*} = \text{Evolutionary methods}(F_{pre}, \mu_{pop})$
 8:     Maintain the recorded elite actor in $\mu_{pop*}$
 9: **end if**
10: **return** New population $\mu_{pop*}$

---

### 4.2.1. Generation-based control

In the generation-based control, a fixed hyperparameter $\omega \in [0, 1)$ is used to indicate the probability of using the surrogate-assisted evaluation. In this setting, the evolutionary operators are partially based on predicted fitness values. Note that, at the beginning of evolution, the population is evaluated in the real environment to collect necessary state data. Algorithm 2 shows the pseudo-code of the generation-based control.

Previous studies have provided some theoretical analyses of the convergence of surrogate-assisted EAs [5, 31]. Here,

**Algorithm 3** Individual-based control

---

**Input**: Policy set of population $\mu_{pop} = \{\mu_1, \mu_2, ..., \mu_n\}$;
  The surrogate model $Q_{rl}$; The RL actor $\mu_{rl}$;
  Replay buffer $\mathcal{R}$ for RL agent training;
  Candidate population size $n^*$; Scaling factor $\sigma$.
**Output**: New population $\mu_{pop*}$

1: **for** $i = 1$ to $n^* - n$ **do**
2:   Sample $\epsilon_i \sim \mathcal{N}(0, 1)$
3:   Inject a new candidate $(\mu_{rl} + \sigma\epsilon_i)$ to $\mu_{pop}$
4: **end for**
5: $F_{pre}$ = Surrogate-assisted evaluation($\mu_{pop}, Q_{rl}, \mathcal{R}$)
6: Retain the best $n - 1$ actors in $\mu_{pop}$ according to $F_{pre}$
7: Inject the recorded elite actor to $\mu_{pop}$
8: $F_{real}$ = Evaluation($\mu_{pop}$)
9: Copy the elite actor
10: $\mu_{pop*}$ = Evolutionary methods($F_{real}, \mu_{pop}$)
11: **return** New population $\mu_{pop*}$

---

we demonstrate how the evolution of parameters is affected by different management strategies. As shown in Figure 4, inaccuracy may be introduced by the surrogate model under the generation-based control, misguiding the update direction. However, it can be corrected appropriately based on real fitness evaluations, as the SC switches between the real and the predicted fitness functions.

### 4.2.2. Individual-based control

In our individual-based approach, the evolutionary operators work on real fitness. Assume that the population size is $n$. Before the population is evaluated in the real environment, a candidate population with $n^*$ offspring is generated with $n^* > n$. After being evaluated by the surrogate model, only the best $n$ individuals are presented to the real environment. This method is also referred to as "preselection strategy". In principle, to generate the candidate population, apart from the original population, $n^* - n$ extra individuals are produced by adding Gaussian noise to the RL actor or the best actor found so far (Algorithm 3). By default, we mutate the RL actor to better explore its surrounding landscape.

As shown in Figure 4, the surrogate model preselects those mutated individuals with relatively higher predicted fitness and filters out any solution that are likely to fail, then the preselected genetic population actually forms a "trust region" [30]. The real fitness evaluation is finally conducted for stable optimization, resulting in a more directional and smooth parameter update path.

### 4.2.3. Elite protection

The evolution parts of previous hybrid methods such as ERL and PDERL follow the spirit of elitism mechanism [18], where the selected elites with high fitness are kept in the population. In practice, the elites copy a part of their genes to other individuals and are protected from mutations.

Apart from the elitism mechanism within hybrid frameworks, another issue may arise when the surrogate model is used for fitness evaluation. As discussed in Section 4.1,

**Table 1**
Action and state dimensions in various environments

| Environment | State dimension | Action dimension |
|:---:|:---:|:---:|
| Ant | 111 | 8 |
| Hopper | 11 | 3 |
| Walker | 17 | 6 |
| Swimmer | 8 | 2 |
| Reacher | 11 | 2 |
| HalfCheetah | 17 | 6 |

inevitably, the ranking of individuals based on the predicted fitness may not perfectly align with that based on the real fitness, resulting in instability of evolution and fluctuation of performance.

To handle this issue, in the generation-based control, SC keeps track of the current best actor when the population is evaluated by the real fitness function, and makes sure it stays in the population after performing evolutionary operators based on predicted fitness. In the individual-based control, the top $n - 1$ actors from the candidate population and the current best actor constructs a new population to be applied to the real environment. Our experiment in Section 5.5 shows that the elite protection mechanism effectively prevents the dramatic fluctuation of performance and the spread of detrimental information from individuals with inaccurate fitness while using the fixed evolution control.

### 4.3. Evaluation memory

Off-policy DRL methods such as DQN [25], DDPG [20], and TD3 [10] maintain a constantly updated replay buffer to improve the sample efficiency of the RL agent. In our approach, the most recent part of the historical data (state information only), referred to as evaluation memory, is exploited to evaluate the population. This memory not only contains diverse state samples but also keeps track of the current optimization process. In addition, evaluation memory does not have to be maintained all the time. It is only created when the surrogate model is called for evaluating individuals and has a negligible memory footprint.

## 5. Experiments and Evaluations

To highlight the value of SC, we focus on the implementation of SC with two state-of-the-art hybrid frameworks ERL and PDERL, referred to as SERL and SPDERL, respectively. We aim to answer the following questions: (1) Does SC improve the computational efficiency and the performance of the original hybrid frameworks? (2) How sensitive is the optimization process to the control parameters of SC's management strategies? (3) What impacts does SC bring to the internal dynamics of the original hybrid frameworks?

### 5.1. Environmental settings

We performed experiments on 6 continuous control tasks: HalfCheetah, Ant, Hopper, Swimmer, Reacher and Walker with the MuJoCo[1] physics engine [41], and Table

---

[1]https://mujoco.org/

**Table 2**
Summary of the baselines and the proposed algorithms

| Algorithms | EA Part | RL Part | Surrogate Fitness | Management Strategies |
|:---:|:---:|:---:|:---:|:---:|
| ERL | GA (N-point Crossover & Gaussian Mutation) | DDPG | No | No |
| SERL-I | GA (N-point Crossover & Gaussian Mutation) | DDPG | Yes | Individual-based Control |
| SERL-G | GA (N-point Crossover & Gaussian Mutation) | DDPG | Yes | Generation-based Control |
| PDERL | GA (Distillation Crossover & Proximal Mutation) | DDPG | No | No |
| SPDERL-I | GA (Distillation Crossover & Proximal Mutation) | DDPG | Yes | Individual-based Control |
| SPDERL-G | GA (Distillation Crossover & Proximal Mutation) | DDPG | Yes | Generation-based Control |

**Table 3**
Hyperparameters of SERL and SPDERL

| Hyperparameter | Value |
|:---|:---:|
| Hidden layers of the actor network | $(64, 64)$ |
| Hidden layers of the critic network | $(400, 300)$ |
| Activation function of the actor network | Tanh |
| Activation function of the critic network | ELU |
| Target weight $\tau$ | 0.001 |
| RL actor learning rate | $5e^{-5}$ |
| RL critic learning rate | $5e^{-4}$ |
| Replay buffer size | $1e^{6}$ |
| RL agent batch size | 128 |
| Discount factor | 0.99 |
| Optimizer | Adam |
| Genetic actor learning rate | $1e^{-3}$ |
| Genetic memory size | 8000 |
| Population size | 10 |
| Genetic agent crossover batch size | 128 |
| Genetic agent mutation batch size | 256 |
| Distillation crossover epochs | 12 |
| Mutation probability | 0.9 |
| Mutation strength | 0.1 |

**Table 4**
Hyperparameters used in various environments

| Environment | Algorithm | Elite | Trials | Sync |
|:---:|:---|:---:|:---:|:---:|
| Ant | SERL | 0.3 | 1 | 1 |
| | SPDERL | 0.2 | 1 | 1 |
| Hopper | SERL | 0.3 | 5 | 1 |
| | SPDERL | 0.2 | 3 | 1 |
| Walker | SERL | 0.2 | 3 | 1 |
| | SPDERL | 0.2 | 5 | 1 |
| Swimmer | SERL | 0.1 | 1 | 10 |
| | SPDERL | 0.1 | 1 | 10 |
| Reacher | SERL | 0.1 | 1 | 10 |
| | SPDERL | 0.1 | 1 | 10 |
| HalfCheetah | SERL | 0.1 | 1 | 1 |
| | SPDERL | 0.1 | 1 | 10 |

1 shows the action and state dimensions in all environments. All these tasks are packaged according to the standard OpenAI Gym API[2] and friendly to simulation. At the beginning of each simulation, an initial state vector determined by internal random seeds is provided, and at each subsequent time step, the policy neural network calculates what action to perform according to the latest state vector. The environment simulates this action and returns a new state vector and the corresponding reward. Additionally, the reward function is task-specific and we consider the real fitness value as the sum of the rewards over one episode or the average reward over several runs.

## 5.2. Algorithm settings

We use the official implementations of ERL[3] [18] and PDERL[4] [2] as the major baselines and follow all their hyperparameter settings for EAs, RL agents, neural networks and the population size ($n = 10$). The actor (agent's policy) is represented by a fully connected neural network with two hidden layers, each containing 64 neurons. The number of neurons in the input and output layers is task-specific (Table 1), and the hidden and output layers take the Tanh activation function. Thus, the policy space is encoded by 5702

parameters in the Walker task, 5702 in HalfCheetah, 5058 in Reacher, 11848 in Ant and 5123 in Hopper. Furthermore, the critic network is also fully connected, and the numbers of neurons in the first and second hidden layers are 400 and 300, respectively. We use the ELU activation between hidden layers.

The main difference between ERL and PDERL is in the EA part. In ERL, the GA employs the typical N-point crossover and Gaussian mutation, while in PDERL, the crossover is implemented by distillation, and the Proximal mutation, based on the SM-G-SUM operator [19] is used. Table 2 provides a brief summary of ERL, PDERL and our methods. Moreover, we use a standard Genetic Algorithm [18] and the DDPG implemented by OpenAI Spinningup[5] as extra baselines to compare our methods against pure EA and RL algorithms. Table 3 shows the hyperparameters of SERL and SPDERL in this work, they are common over all environments. "Elite" in Table 4 represents the proportion of elite individuals in the genetic population. "Trials" is referred to as the number of evaluation times of an individual, for Walker and Hopper, where the reward variance is relatively higher than other environments, policies need run more times to obtain their average fitness. Finally, "Sync" is the synchronization period of the RL-Actor.

For SC's default hyperparameters, the maximum evaluation memory size $k$ is limited to $50,000$, and we fix $\omega$ to 0.6 for the generation-based control, which means that, in each generation, the population has a 60% chance of being

**Figure 5:** Learning curves on six MuJoCo environments: HalfCheetah, Ant, Hopper, Swimmer, Reacher and Walker.

evaluated by the surrogate model. For the individual-based control, $\alpha = (n^* - n)/n$ is the control factor for the candidate population size with $\alpha = 1$ in our experiments and we set the scaling factor $\sigma = 0.01$ for generating Gaussian noise. During the training, SC keeps recording the current best actor when the population is evaluated in the real environment for elite protection. The periodical synchronization of the RL actor and the population is performed only after real fitness evaluation. We refer to SERL and SPDERL with generation-based control as SERL-G and SPDERL-G. Similarly, SERL-I and SPDERL-I indicate the combination of SERL and PDERL with individual-based control, respectively.

### 5.3. Overall performance

We train each proposed method with 6 different random seeds on 6 MuJoCo environments following the convention in literature [41]. During the training process, the average of 5 test results of the best actor from the genetic population is reported as the performance of each algorithm. Figure 5 illustrates their learning processes during training. The solid curves represent the mean values and the shaded regions indicate the standard deviations.

In most environments, SC can significantly improve the learning speed and the performance of the original hybrid frameworks, and also make the learning process more stable with lower variance. For example, SERL-G can outperform ERL across all environments. Except for Swimmer, the improvement of SERL-I is more evident in the early training phase (within 1 million steps) compared to SERL-G, and it outperforms other methods in Reacher. When it comes to SPDERL, both SPDERL-I and SPDERL-G outperform PDERL in Hopper and Walker. For Ant and HalfCheetah, SPDERL-I can achieve higher final performance while

SPDERL-G can accelerate the performance improvement in the early training phase. As for pure RL and EA methods, GA struggles in most environments except Swimmer. In HalfCheetah, DDPG performs better than any other methods in the early training phase but the performance tends to converge after 3 million steps.

It is worth noting that DDPG and all methods under the individual-based control fail in Swimmer and, as explained in [18, 29], DRL methods face great challenges in effectively learning the gradient information. To alleviate this issue, we generate the candidate population by mutating the best actor that has been found by genetic operations (Appendix A). In this specific environment, the evolutionary search is more suitable for driving the optimization process.

### 5.4. Parameter analysis

In this part, we investigate the influence of the following parameters: the ratio of surrogate-assisted evaluation $\omega$, the control factor $\alpha$ of the candidate population size, and the capacity $k$ of evaluation memory, as shown in Figure 6.

**Control factors.** We vary $\omega$ from 0.2 to 0.8 in SERL-G and $\alpha$ from 0.5 to 2.0 in SERL-I, respectively. The results indicate that both the final performance and the learning speed are generally improved by increasing $\omega$ under the generation-based control. Although a high value (e.g., $\omega = 0.8$) may lead to a little drop in the final performance, it significantly reduces the number of interactions with the environment and speeds up the learning process. For individual-based control, a relatively small $\alpha$ value (e.g., $\alpha = 0.5$) is more cost-effective. As the surrogate model needs to evaluate additional individuals generated by mutations, a high value of $\alpha$ may result in overhead, especially when the input of the surrogate model is high-dimensional.

**Figure 6:** Parameter analysis of SERL-G with different generation-based control factors and evaluation memory sizes (first two columns) and SERL-I with different individual-based control factors and evaluation memory sizes (last two columns) in Walker and Hopper environments.



**Figure 7:** Comparisons of SERL-G (first two figures) and SERL-I (last two figures) with and without Elite Protection (EP) mechanism in Walker and Hopper environments.

**Memory capacity.** Furthermore, we investigate the impact of the capacity of the evaluation memory. The second row of Figure 6 shows the performance of SERL-G and SERL-I with fixed control factors $\omega = 0.6$, $\alpha = 1.0$ and various $k$ values in two environments. Overall, a large evaluation memory contains more diverse state data and significantly helps improve the quality of the evaluation, but it also increases the computational cost of SC. In general, the generation-based control is better suited with large evaluation memories to counteract the bias in surrogate-assisted evaluation. By contrast, since the evolutionary methods are based on real fitness evaluation with low deviations, a relatively small $k$ is suitable for individual-based control.

## 5.5. Elite protection evaluation

Figure 7 shows the performance of SERL-I and SERL-G without the elite protection mechanism. In this setting, SC only speeds up the policy improvement in the early period of the training process and then encounters a dramatic drop in performance. Although the surrogate model can provide a roughly accurate estimation of population fitness, its estimation of elites is possibly biased, which increases the risk of discarding elites from the population. This experiment underlines the importance of elite protection while using the surrogate for fitness evaluation.

## 5.6. Changes of the internal dynamics

**Interactions between RL and EA.** To gain a deeper insight into the internal transformation of hybrid frameworks in the presence of SC, we highlight the changes in the internal dynamic between the RL agent and the genetic population. We keep a separate record of the accumulative rates of the RL actor being selected, discarded, or chosen as an elite in the population in Hopper. In Figure 8, although ERL, SERL-G, and SERL-I present similar dynamics, the integration of SC significantly stabilizes the internal dynamic of learning and evolution, which is more pronounced in PDERL and the two variants of SPDERL, where the evolutionary search is the major driving force for the training process. It is reasonable to hypothesize that optimization based both on the real and predicted fitness function may benefit the evolutionary search and make the optimization process more stable.

**Intriguing behavioral patterns.** A notable discovery is that the agents trained by our proposed methods can produce highly intriguing behavioral patterns than the original hybrid approaches, as shown in Figure 9. The solutions found by hybrid frameworks with SC are more intriguing and stable. On some control tasks, they can better adapt themselves to the environments and perform more competently.

**Figure 8:** Accumulative rates of the RL agent being selected, discarded, or chosen as the elite during the training process in Hopper. The results indicate that SC can potentially stabilize the internal dynamics of the original frameworks.



(a) HalfCheetach agent



(b) Hopper agent

**Figure 9:** Several intriguing behavioral patterns of the agents trained by SPDERL-I and SPDERL-G. (a) A HalfCheetah agent trained by SPDERL-I with average performance of 14000 points over 50 test seeds. The agent is able to adjust its posture more appropriately and run faster. (b) A Hopper agent trained by SPDERL-G with average performance of 4100 points over 50 test seeds. The agent jumps faster and learns to better stabilize the center of gravity.

## 5.7. Computational efficiency

The number of real evaluations is potentially limited due to time and/or money, it is critical to make full use of available resources to achieve expected performance. In this part, we aim to validate the improvement of computational efficiency of original hybrid frameworks when SC is introduced, mainly from the following aspects: the sample reduction, time consumption and the Floating Point Operations (FLOPs).

**Sample consumption.** We conduct a sample reduction study of different methods in various environments when reaching the target scores. The number of environmental interactions (time steps) is equal to the number of consumed samples, and the results are reported over 6 suns. According to Table 5, SC can significantly reduce the sample consumption across most domains, especially in the environment with relatively higher noise and the reward variance like Walker and Hopper. For instance, ERL needs to perform almost 3 million environmental interactions more than SERL-G and SERL-I in Walker to reach the same score, and PDERL needs to consume 3 hundred thousand interactions more than SPDERL-G and SPDERL-I in Hopper. For different control strategies, the individual-based control requires a relatively small number of samples compared with generation-based

**Table 5**
The average time steps (in million, equal to the number of consumed samples) of different algorithms in various environments when reaching the target score. "−" represents that DDPG or GA can not reach the target.

| Env | Score | DDPG | GA | ERL | SERL-G | SERL-I | PDERL | SPDERL-G | SPDERL-I |
|-----|-------|------|-----|-----|--------|--------|-------|----------|----------|
| Ant | 5,000 | − | − | 4.266 | **1.564** | **1.685** | 1.464 | 1.784 | **1.085** |
| Hopper | 2,500 | − | − | 2.519 | **1.073** | **0.786** | 0.857 | **0.554** | 0.561 |
| Walker | 2,000 | 0.775 | − | 3.572 | **0.767** | **0.505** | 1.312 | **0.796** | 0.508 |
| Swimmer | 300 | − | 2.640 | 0.968 | **0.516** | 1.034 | 0.742 | 1.489 | 3.568 |
| Reacher | −5 | − | − | 1.051 | **0.322** | **0.224** | 0.717 | **0.513** | 0.672 |
| HalfCheetah | 10,000 | 0.986 | − | 4.092 | **2.565** | **1.761** | 2.234 | 2.153 | 2.201 |



**Figure 10:** Comparison of ERL, SERL-I and SERL-G in terms of the training time and performance in Ant and Hopper.



**Figure 11:** Comparison of ERL, SERL-I and SERL-G in terms of the number of FLOPs and performance (red star) achieved at 3M environment steps in Walker.

control, due to its mechanism of pre-selecting potential high-performing individuals, which can be subsequently presented in real environments to generate higher-quality samples, as discussed in Section 4.2.2. In a nutshell, SC can further improve the sample-efficiency of the original hybrid frameworks, as the diverse historical data is not only used for training the DRL algorithm but also employed to help evaluate the genetic population.

**Time consumption.** We focus on the training time and corresponding performance of different methods, especially on the time costs reduced by the surrogate model. As SC is applied on top of the original algorithms, there are no changes in the network structures of actors and critics. The real evaluation of the genetic population is serial, and the results reported are averaged over 6 runs, and the total training time of each run is limited to 8 hours. The parameter settings are: $\omega = 0.6$ for SERL-G and $\alpha = 1.0$ for SERL-I. According to Figure 10, in the Ant environment, it is hard

for ERL to reach 3000 points within 8 hours, while SERL-G takes approximately 8 hours and SERL-I only takes 4.5 hours. In the Hopper environment, it takes 8 hours for ERL to reach above 1500 points, 7 hours for SERL-G to reach 2500 points, while SERL-I is much more time-efficient, achieving the same performance as SERL-G in only 3 hours. When using real fitness evaluation, the actor only needs to perform forward propagation through its policy neural network, such as the calculation of $\mu(s_j|\theta^\mu)$, and the reward is provided by the environment itself, which may be time-consuming and expensive. While using the surrogate-assisted evaluation, additional forward propagation will be introduced to calculate the averaged value of $Q(s_j, \mu(s_j|\theta^\mu)|\theta^Q)$ via diverse historical state information. SC can efficiently transfer the computational burden of real evaluations to additional forward propagation through the surrogate model, and the computational time of these calculations can be easily reduced by parallelization and shared memory.

**Floating Point Operations.** The computational efficiency of SC is finally verified by comparing the FLOPs consumed by ERL, SERL-I, SERL-G and their corresponding performance (Figure 11). The neural networks in our work are all fully connected and the number of updates of the RL agent is equal to the environment steps. Thus, regardless of which hybrid framework that SC is combined with, the computational cost of the entire training process can always be divided into the surrogate-assisted evaluation cost and the original optimization cost, as shown in Appendix B. As mentioned above, while bringing significant performance improvement, SC is computationally efficient in that only a small amount of forward propagation for surrogate-assisted evaluations is introduced.

## 6. Conclusion and Future Work

The application of hybrid RL frameworks to expensive learning problems has largely been limited by the cost of evaluating the population in real environments. In this work, we propose a surrogate-assisted controller with two management strategies, which can be easily integrated into existing hybrid frameworks to simultaneously facilitate the optimization of the RL agent and the evaluation of the population. To the best of our knowledge, this is the first attempt on introducing the surrogate model into hybrid RL frameworks. Empirical evaluations show that the combination of SC

**Figure 12:** Learning curves using different approaches of generating the candidate population in the Swimmer environment.

**Table 6**
Symbolic meanings and values in the Calculation of FLOPs

| Symbol | Value |
|---|---|
| Environment steps $T$ | $3,000,000$ |
| Population size $n$ | 10 |
| Candidate population size $n_c$ | 20 |
| RL Agent batch size $b$ | 128 |
| Evaluation memory size $k$ | $50,000$ |
| Evolutionary generations $G$ | 240 |
| Evolutionary generations using the surrogate $G_s$ | 390 |
| Flops of forward propagation of RL-Actor $A_f$ | $11,136$ |
| Flops of forward propagation of RL-Critic $C_f$ | $249,800$ |
| Flops of backward propagation of RL-Actor $A_b$ | $22,272$ |
| Flops of backward propagation of RL-Critic $C_b$ | $499,600$ |

with two state-of-the-art evolutionary reinforcement learning frameworks ERL and PDERL can effectively reduce the evaluation cost and boost the performance. Furthermore, SC brings beneficial changes to the internal dynamics of learning and evolution, resulting in more collaborative interactions between the RL agent and the EA population.

Learning and evolution are two symbiotic counterparts in nature. It is of great scientific importance to fully appreciate and explore this hybrid paradigm towards implementing truly competent artificial intelligence. As to future work, there are plenty of fascinating directions to advance our proposed techniques, including real-world settings, such as embodied AI, designing self-adaptive evolution control strategies and more effective evaluation criteria for the surrogate. For complex and challenging settings, multi-agent evolutionary reinforcement learning with multiple surrogates is expected to further extend the horizon.

## Appendix A. The Swimmer environment

Results in Figure 5 demonstrate that all individual-based control methods fail in the Swimmer environment because the RL agent is unfortunately misled by the deceptive gradient information, and candidates mutated from the RL actor could not provide useful information for policy improvement. To alleviate this issue, we generate the candidate population by mutating the best actor that has been found by genetic operations (Figure 12). In this specific environment, the evolutionary search is more suitable for driving the optimization process.

## Appendix B. Measurement of FLOPs

In this section, we introduce how the FLOPs values in Figure 11 are calculated. The neural networks in our work are all fully connected and the number of updates of the RL agent is equal to the environment steps. Similar to [35], we consider the FLOPs of each forward pass being half of that

of the backward pass. We also assume that the consumed FLOPs of activation functions and operations that do not require passing through neural networks are negligible. We follow the procedures in [27] to measure the FLOPs consumed by each method within 3M environment steps. The calculation formulas for FLOPs consumed by ERL, SERL-I and SERL-G are shown in Eqs.(8), (9) and (10), respectively.

$$F_{ERL} = T(A_f + b \times (2A_f + 3C_f + C_b + A_b)) \quad (8)$$

$$F_{SERL(I)} = F_{ERL} + G_s \times n \times k \times (A_f + C_f) \quad (9)$$

$$F_{SERL(G)} = F_{ERL} + G \times n_c \times k \times (A_f + C_f) \quad (10)$$

All the symbols and their symbolic meanings that are involved in calculating are shown in Table 6.

## Appendix C. Combining SC with CEM-RL

We conduct an additional experiment on combining SC with another state-of-the-art hybrid framework CEM-RL [29], referred to as Surrogate-assisted CEM-RL (SCEM-RL), to show SC can be applied to a different type of EA. The EA parts of ERL and PDERL are based on GA, while in CEM-RL, it is based on Cross-Entropy Method (CEM), which is a typical Evolutionary Strategy (ES). CEM-RL combines CEM with TD3 learners [10], an off-policy DRL algorithm related to DDPG [20].

We train SCEM-RL-G (SCEM-RL with generation-based control) and SCEM-RL-I (SCEM-RL with individual-based control) with 6 different random seeds. Figure 13 illustrates their learning curves during training on three control tasks from MuJoCo [41]. Apart from the clear advantages of SC in the training performance and sample consumption, similar to the results in Section 5.3, on HalfCheetah and Ant, the improvement over the original hybrid framework is more evident under the individual-based control, while the generation-based control method is more favorable on Walker.

**Figure 13:** Learning curves on three MuJoCo environments: Ant, HalfCheetah and Walker.

# References

[1] Charles Audet, J Denni, Douglas Moore, Andrew Booker, and Paul Frank. A surrogate-model-based method for constrained optimization. In *8th symposium on multidisciplinary analysis and optimization*, page 4891, 2000.

[2] Cristian Bodnar, Ben Day, and Pietro Lió. Proximal distilled evolutionary reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3283–3290, 2020.

[3] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.

[4] Di Cao, Junbo Zhao, Weihao Hu, Fei Ding, Nanpeng Yu, Qi Huang, and Zhe Chen. Model-free voltage control of active distribution system with pvs using surrogate model-based deep reinforcement learning. *Applied Energy*, 306:117982, 2022.

[5] Yu Chen, Weicheng Xie, and Xiufen Zou. How can surrogates influence the convergence of evolutionary algorithms? *Swarm and Evolutionary Computation*, 12:18–23, 2013.

[6] Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pages 465–472. Citeseer, 2011.

[7] Huachao Dong, Peng Wang, Chongbo Fu, and Baowei Song. Kriging-assisted teaching-learning-based optimization (ktlbo) to solve computationally expensive constrained problems. *Information Sciences*, 556:404–435, 2021.

[8] Madalina M Drugan. Reinforcement learning versus evolutionary computation: A survey on hybrid algorithms. *Swarm and evolutionary computation*, 44:228–246, 2019.

[9] Olivier Francon, Santiago Gonzalez, Babak Hodjat, Elliot Meyerson, Risto Miikkulainen, Xin Qiu, and Hormoz Shahrzad. Effective reinforcement learning through evolutionary surrogate-assisted prescription. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference*, pages 814–822, 2020.

[10] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, pages 1587–1596. PMLR, 2018.

[11] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.

[12] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *arXiv preprint arXiv:1809.01999*, 2018.

[13] Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.

[14] Andrew Ilyas, Logan Engstrom, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and Aleksander Madry. A closer look at deep policy gradients. *arXiv preprint arXiv:1811.02553*, 2018.

[15] Yaochu Jin. Surrogate-assisted evolutionary computation: Recent advances and future challenges. *Swarm and Evolutionary Computation*, 1(2):61–70, 2011.

[16] Yaochu Jin, Markus Olhofer, and Bernhard Sendhoff. Managing approximate models in evolutionary aerodynamic design optimization. In *Proceedings of the 2001 Congress on Evolutionary Computation (IEEE Cat. No. 01TH8546)*, volume 1, pages 592–599. IEEE, 2001.

[17] Andy Keane, Alexander Forrester, and Andras Sobester. *Engineering design via surrogate modelling: a practical guide*. American Institute of Aeronautics and Astronautics, Inc., 2008.

[18] Shauharda Khadka and Kagan Tumer. Evolution-guided policy gradient in reinforcement learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 1196–1208, 2018.

[19] Joel Lehman, Jay Chen, Jeff Clune, and Kenneth O Stanley. Safe mutations for deep and recurrent neural networks through output gradients. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 117–124, 2018.

[20] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

[21] Ilya Loshchilov, Marc Schoenauer, and Michèle Sebag. A mono surrogate for multiobjective optimization. In *Proceedings of the 12th annual conference on Genetic and evolutionary computation*, pages 471–478, 2010.

[22] Shuai Lü, Shuai Han, Wenbo Zhou, and Junwei Zhang. Recruitment-imitation mechanism for evolutionary reinforcement learning. *Information Sciences*, 553:172–188, 2021.

[23] Enrico Marchesini, Davide Corsi, and Alessandro Farinelli. Genetic soft updates for policy evolution in deep reinforcement learning. In *International Conference on Learning Representations*, 2020.

[24] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.

[25] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.

[26] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.

[27] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. *arXiv preprint arXiv:1611.06440*, 2016.

[28] Jeng-Shyang Pan, Nengxian Liu, Shu-Chuan Chu, and Taotao Lai. An efficient surrogate-assisted hybrid optimization algorithm for expensive optimization problems. *Information Sciences*, 561:304–325, 2021.

[29] Aloïs Pourchot and Olivier Sigaud. Cem-rl: Combining evolutionary and gradient-based methods for policy search. *arXiv preprint*

*arXiv:1810.01222*, 2018.

[30] MJD Powell. On the convergence of a wide range of trust region methods for unconstrained optimization. *IMA journal of numerical analysis*, 30(1):289–301, 2010.

[31] Nestor V Queipo, Raphael T Haftka, Wei Shyy, Tushar Goel, Rajkumar Vaidyanathan, and P Kevin Tucker. Surrogate-based analysis and optimization. *Progress in aerospace sciences*, 41(1):1–28, 2005.

[32] Alain Ratle. Optimal sampling strategies for learning a fitness model. In *Proceedings of the 1999 Congress on Evolutionary Computation-CEC99 (Cat. No. 99TH8406)*, volume 3, pages 2078–2085. IEEE, 1999.

[33] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.

[34] Gisbert Schneider. Neural networks are useful tools for drug design. *Neural Networks*, 13(1):15–16, 2000.

[35] Younggyo Seo, Lili Chen, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. State entropy maximization with random encoders for efficient exploration. *arXiv preprint arXiv:2102.09430*, 2021.

[36] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.

[37] Jörg Stork, Martin Zaefferer, and Thomas Bartz-Beielstein. Improving neuroevolution efficiency by surrogate model-based optimization with phenotypic distance kernels. In *International Conference on the Applications of Evolutionary Computation (Part of EvoStar)*, pages 504–519. Springer, 2019.

[38] Jörg Stork, Martin Zaefferer, Thomas Bartz-Beielstein, and AE Eiben. Surrogate models for enhancing the efficiency of neuroevolution in reinforcement learning. In *Proceedings of the genetic and evolutionary computation conference*, pages 934–942, 2019.

[39] Karush Suri, Xiao Qi Shi, Konstantinos N Plataniotis, and Yuri A Lawryshyn. Maximum mutation reinforcement learning for scalable control. *arXiv preprint arXiv:2007.13690*, 2020.

[40] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

[41] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012.

[42] Hao Tong, Changwu Huang, Leandro L Minku, and Xin Yao. Surrogate models in evolutionary single-objective optimization: A new taxonomy and experimental study. *Information Sciences*, 562:414–437, 2021.

[43] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.

[44] Xilu Wang, Yaochu Jin, Sebastian Schmitt, and Markus Olhofer. An adaptive bayesian approach to surrogate-assisted evolutionary multi-objective optimization. *Information Sciences*, 519:317–331, 2020.

[45] Bruce H Weber and David J Depew. *Evolution and learning: The Baldwin effect reconsidered*. Mit Press, 2003.

[46] Darrell Whitley. A genetic algorithm tutorial. *Statistics and computing*, 4(2):65–85, 1994.

[47] Dennis G Wilson, Sylvain Cussat-Blanc, Hervé Luga, and Julian F Miller. Evolving simple programs for playing atari games. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 229–236, 2018.

[48] Peng Yang, Qi Yang, Ke Tang, and Xin Yao. Parallel exploration via negatively correlated search. *Frontiers of Computer Science*, 15(5):1–13, 2021.

[49] Peng Yang, Hu Zhang, Yanglong Yu, Mingjia Li, and Ke Tang. Evolutionary reinforcement learning via cooperative coevolutionary negatively correlated search. *Swarm and Evolutionary Computation*, 68:100974, 2022.

[50] Xin Yao. Evolving artificial neural networks. *Proceedings of the IEEE*, 87(9):1423–1447, 1999.