





# Evolutionary Diversity Optimization with Clustering-based Selection for Reinforcement Learning

*Yutong Wang, <u>Ke Xue</u>, Chao Qian* ICLR'22 LAMDA 11 Pre by Ke Xue 22/06/24

### Introduction

Reinforcement Learning (RL)

- Most RL methods aim to obtain a single optimal policy
- Some complex scenarios need **a set of diverse policies** 
  - Exploration
  - Policy ensemble
  - Model/environment generation
  - Few-shot adaption

*How to efficiently obtain a set of high-quality policies with diverse behaviors is a challenging problem in RL* 

## Background

• Evolutionary Algorithm



- Key properties
  - Population-based search
  - Only evaluation is required
- Application: complex optimization problems
  - Black-box optimization problems
  - Multi-objective problems

### Background

• Quality-diversity



### Background - QD



### Background - QD

```
Algorithm 2 QD-Optimization algorithm (I iterations)
  \mathcal{A} \leftarrow \emptyset
                                                                                                                                                    ▷ Creation of an empty container.
  for iter = 1 \rightarrow I do
                                                                                                                                      ▶ The main loop repeats during I iterations.
       if iter == 1 then
                                                                                                                                                                           ▶ Initialization.
            \mathcal{P}_{\text{parents}} \leftarrow \text{random}()
                                                                                                                 ▶ The first 2 batches of individuals are generated randomly.
            \mathcal{P}_{\text{offspring}} \leftarrow \text{random}()
                                                                                  > The next controllers are generated using the container and/or the previous batch.
        else
                                                                                    ▷ Selection of a batch of individuals from the container and/or the previous batch.
            \mathcal{P}_{\text{parents}} \leftarrow \text{selection}(\mathcal{A}, \mathcal{P}_{\text{offspring}})
            \mathcal{P}_{offspring} \leftarrow variation(\mathcal{P}_{parents}) Reproduction
                                                                                      \triangleright Creation of a randomly modified copy of \mathcal{P}_{parents} (mutation and/or crossover).
       for each \theta \in \mathcal{P}_{offspring} do
             \{f_{\boldsymbol{\theta}}, \boldsymbol{b}_{\boldsymbol{\theta}}\} \leftarrow f(\boldsymbol{\theta})
                                                                                       ▶ Evaluation of the individual and recording of its descriptor and performance.
            if ADD_TO_CONTAINER(\theta, \mathcal{A}) then
                                                                               ▶ "ADD_TO_CONTAINER" returns true if the individual has been added to the container.
                                                                                                                                                     ▶ The parent might get a reward.
                 UPDATE_SCORES(parent(\theta), Reward, \mathcal{A})
            else
                 UPDATE scores(parent(\theta), -Penalty, \mathcal{A})
                                                                                                                                                  ▷ Otherwise, it might get a penalty.
                                                                                 ▶ Update of the attributes of all the individuals in the container (e.g. novelty score).
       UPDATE CONTAINER(\mathcal{A})
   return A
```

- Policy gradient assisted MAP-Elites (PGA-ME) [GECCO'21 Best paper]
   Using PG to update the solutions
- Differentiable Quality-diversity (DQD) [NeurIPS'21 Oral, < 2.4%]</li>
   Using differentiable *behavior descriptor b*

### Method - motivation

### The **inefficient** selection

- Uniform selection
- Biased selection
- Pareto-based selection

Method	Selection
Vanilla ES	The only parent solution
NSR-ES	Probabilistic selection
CVT-ES	Uniform selection
ME-ES	<b>Biased selection</b>
DvD-ES	All parent solutions
QD-RL	Pareto-based selection



### Method - motivation

The **inefficient** selection

- Uniform selection
- Biased selection
- Pareto-based selection





### Method - motivation

The **inefficient** selection

- Uniform selection
- Biased selection
- Pareto-based selection



-5

-10

-15

-20

-25

-30

-35

-40

### Method

#### Algorithm 1: EDO-CS

```
Input: number K of selected policies, number T of total iterations, number T' of updating
             iterations, behavior characterization b(\pi_{\theta}), archive size l, archive A, bandit B
   Output: archive A
   // Warm up
 1 for j = 1 : l do
        Randomly generate policy \pi_{\theta_i};
 2
        Get cumulative rewards R and behavior b(\pi_{\theta_j}) by evaluating the policy \pi_{\theta_j};
 3
        Add (\pi_{\theta_i}, R, b(\pi_{\theta_i})) into archive A
 4
 5 end
 6 t = 0;
                                                         Clustering-based selection
 7 while t < T do
        // Selection
        Use K-means to divide the policies in archive A into K clusters \{C_k\}_{k=1}^K;
 8
        Select K policies \{\pi_{\theta_k}\}_{k=1}^K, each one from a cluster;
 9
        // Reproduction
                                                                                                Div(\boldsymbol{\theta}) = \frac{1}{k} \sum_{\pi'_{\boldsymbol{\theta}} \in A'_{\pi_{\boldsymbol{\theta}},k}} \|b(\pi_{\boldsymbol{\theta}}) - b(\pi'_{\boldsymbol{\theta}})\|_{2}.
        for k = 1 : K do
10
             // Update in parallel
             Sample \lambda_k from the bandit B;
11
             for i = 1 : T' do
12
                  Set the objective function J(\boldsymbol{\theta}_k) = (1 - \lambda_k) \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}_k}} [R(\tau)] + \lambda_k Div(\boldsymbol{\theta}_k);
13
                  Use ES to update \theta_k as Eq. (6)
14
             end
15
             Get cumulative rewards R and behavior b(\pi_{\theta'_k}) by evaluating the updated policy \pi_{\theta'_k}
16
        end
17
        Update archive A and bandit B;
18
        t = t + T'
19
20 end
```

### Method



Method	Selection	Reproduction	EAs type	From archive
Vanilla ES	The only parent solution	Quality	(1, 1)	×
NSR-ES	Probabilistic selection	Quality and diversity	(K, 1)	×
CVT-ES	Uniform selection	Quality and diversity	(K+K)	$\checkmark$
ME-ES	<b>Biased selection</b>	Quality or diversity	(K+K)	$\checkmark$
DvD-ES	All parent solutions	Quality and diversity	(K, K)	×
QD-RL	Pareto-based selection	Quality or diversity	(K+K)	$\checkmark$
EDO-CS	Clustering-based selection	Quality and diversity	(K+K)	$\checkmark$

### Experiment



### Experiment

### • Multi-modal

Environment	EDO-CS	QD-RL	ME-ES	DvD-ES	CVT-ES	NSR-ES	Vanilla ES
HalfCheetahFwd	4284	2930	2700	-3419	3219	1346	-5543
HalfCheetahBwd	6548	6013	5953	6353	4672	5366	3911
AntFwd	<b>4617</b>	4291	4316	4507	3856	1737	1911
AntBwd	4697	4164	4123	3498	2958	3961	-851
Performance Ranking	1	3	3.5	3.75	4.75	5.25	6.75

• Single-modal



### Experiment – ablation studies

• Adaptive  $\lambda$ 

Setting of $\lambda$	AntWall-v0	HalfCheetahFwd	HalfCheetahBwd	AntFwd	AntBwd
Adaptive $\lambda$	-529	4284	6548	4617	4697
$\lambda = 0$	-650	3856	5931	4340	4426
$\lambda = 0.5$	-850	3877	6077	2800	3093

### Experiment – ablation studies

- Adaptive  $\lambda$
- Compared with direct optimization

$$f_{pair} = (1 - \omega) \sum_{i=1}^{K} \mathbb{E}_{\tau \sim \pi_{\theta_i}} [R(\tau)] + \omega \sum_{i=1}^{K} \sum_{j \neq i} \|b(\pi_{\theta_i}) - b(\pi_{\theta_j})\|_2 \qquad f_{det} = (1 - \omega) \sum_{i=1}^{K} \mathbb{E}_{\tau \sim \pi_{\theta_i}} [R(\tau)] + \omega \cdot det(\Pi)$$

Method	AntWall-v0	HalfCheetahFwd	HalfCheetahBwd	AntFwd	AntBwd
EDO-CS	-529	4284	6548	4617	4697
$EDO-DOS_{pair}$	-706	4188	5847	5013	4392
$EDO-DOS_{det}$	-536	-5591	6529	-530	2417

Method	AntWall-v0	HalfCheetahFwdBwd	AntFwdBwd
EDO-CS	0.015	0.021	0.023
$EDO-DOS_{pair}$	9.570	9.802	9.807
$EDO-DOS_{det}$	8.900	9.266	9.567

### Experiment – ablation studies

- Adaptive  $\lambda$
- Compared with direct optimization
- Clustering algorithm
- Number *T* ′ of updating iterations
- Population size *K*
- Archive size *l*

More results are shown in the Appendix

### Discussion

Conclusion

- QD is an interesting and attractive research area
- We proposed a *"Simple but efficient"* selection method for QD Future work
- Application of QD
  - -Human-AI coordination
  - Environment generation





# Thanks you!

Our code can be found at: <u>www.lamda.nju.edu.cn/qianc/code-EDOCS.html</u>

