

# Apply Novelty Search And Evolution Strategy To Improve Exploration For Deep Reinforcement Learning

Võ Đức Dương   duong922003@gmail.com

Ngày 16 tháng 4 năm 2025

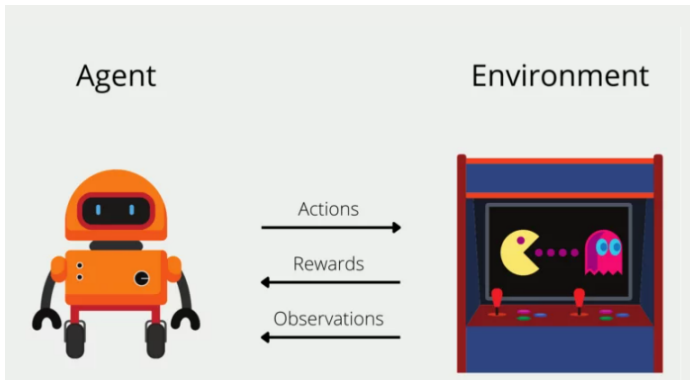
# Mục lục

- 1 Giới thiệu đề tài
- 2 Cơ sở lý thuyết
  - DQN
  - ES
  - NS
  - QD
- 3 Phương pháp tiếp cận
  - NS-ES
  - QD-ES algo: NSR-ES
  - NS-ES customized
  - NSR-ES customized
- 4 Experiment results

# Mục lục

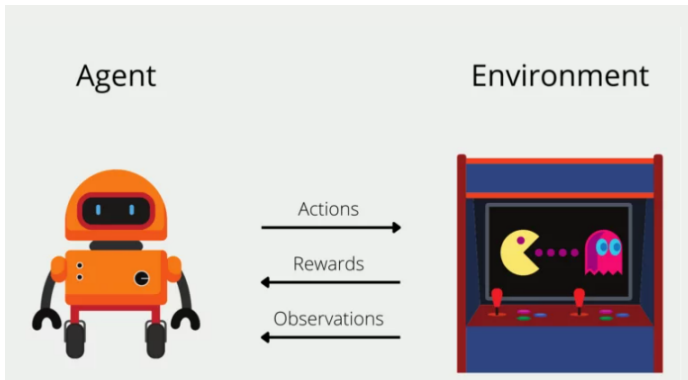
- 1 Giới thiệu đề tài
- 2 Cơ sở lý thuyết
  - DQN
  - ES
  - NS
  - QD
- 3 Phương pháp tiếp cận
  - NS-ES
  - QD-ES algo: NSR-ES
  - NS-ES customized
  - NSR-ES customized
- 4 Experiment results

# Reinforcement Learning



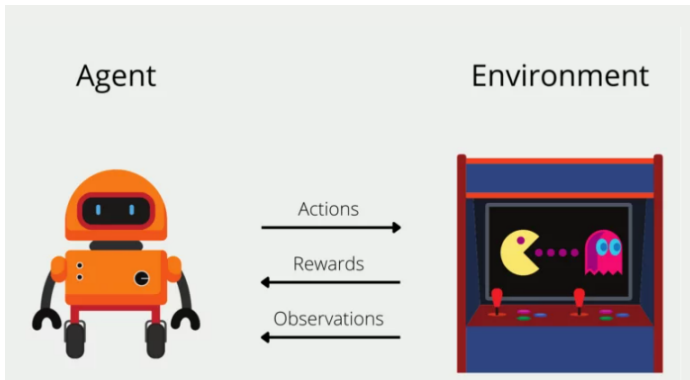
- Agent: tối ưu hóa điểm thưởng

# Reinforcement Learning



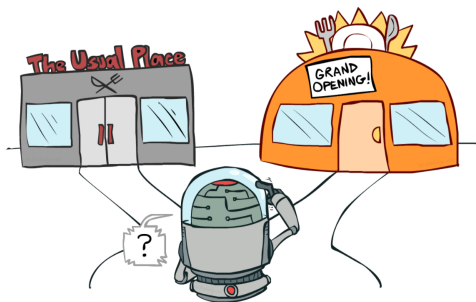
- Agent: tối ưu hóa điểm thưởng
- **Vấn đề:** hàm (tín hiệu) phần thưởng có thể gây nhầm lẫn hoặc khan hiếm. Khiến agent kẹt ở điểm tối ưu cục bộ

# Reinforcement Learning



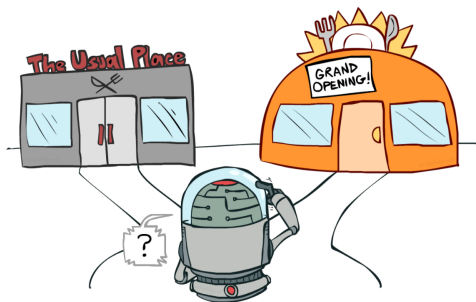
- Agent: tối ưu hóa điểm thưởng
- **Vấn đề:** hàm (tín hiệu) phần thưởng có thể gây nhầm lẫn hoặc khan hiếm. Khiến agent kẹt ở điểm tối ưu cục bộ  
→ Cần chiến lược khám phá thông minh hơn

# Exploration in RL



- Giải pháp: thúc đẩy agent đến các trạng thái hiếm hoặc chưa được khám phá → khuyến khích khám phá (exploration)

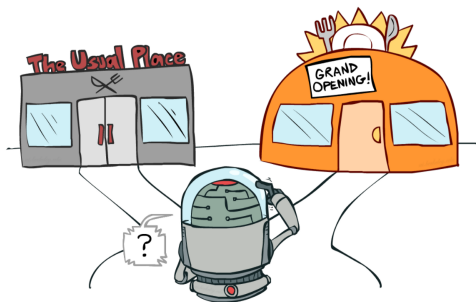
# Exploration in RL



- Giải pháp: thúc đẩy agent đến các trạng thái hiếm hoặc chưa được khám phá → khuyến khích khám phá (exploration)
- **Ví dụ:** Kỹ thuật đếm tần suất thăm trạng thái dựa trên mã hóa trạng thái hoặc mô hình mật độ



# Exploration in RL



- Giải pháp: thúc đẩy agent đến các trạng thái hiếm hoặc chưa được khám phá → khuyến khích khám phá (exploration)
- **Ví dụ:** Kỹ thuật đếm tần suất thăm trạng thái dựa trên mã hóa trạng thái hoặc mô hình mật độ
- **Vấn đề:** Các trạng thái được đếm một cách độc lập với nhau.

# NS & QD for DRL

- **Hướng tiếp cận khác:** Thiết kế hoặc học một mô tả về hành vi trong suốt lifetime của agent.

# NS & QD for DRL

- **Hướng tiếp cận khác:** Thiết kế hoặc học một mô tả về hành vi trong suốt lifetime của agent.
- **Mục tiêu:** Khuyến khích agent thực hiện hành vi khác với những hành vi trước đó. Ví dụ *novelty search (NS)* và *quality diversity (QD)*.
- **Đặc điểm:** Khám phá bằng một quần thể các agent.

# NS & QD for DRL

- **Hướng tiếp cận khác:** Thiết kế hoặc học một mô tả về hành vi trong suốt lifetime của agent.
- **Mục tiêu:** Khuyến khích agent thực hiện hành vi khác với những hành vi trước đó. Ví dụ *novelty search (NS)* và *quality diversity (QD)*.
- **Đặc điểm:** Khám phá bằng một quần thể các agent.

**Nội dung đề tài:** Sử dụng ý tưởng của NS và QD trong các thuật toán tiến hóa để giải quyết các bài toán Deep Reinforcement Learning (DRL).

# Mục lục

- 1 Giới thiệu đề tài
- 2 Cơ sở lý thuyết
  - DQN
  - ES
  - NS
  - QD
- 3 Phương pháp tiếp cận
  - NS-ES
  - QD-ES algo: NSR-ES
  - NS-ES customized
  - NSR-ES customized
- 4 Experiment results

# DQN

- DQN (Deep Q-Network) là thuật toán được mở rộng từ thuật toán Q-Learning, sử dụng một mạng nơ-ron để ước lượng Q-value thay vì sử dụng Q-table như Q-Learning truyền thống.

# Q-Learning

- Hàm  $Q(s, a)$  đại diện cho giá trị kỳ vọng nhận được khi thực hiện hành động  $a$  ở trạng thái  $s$  và tuân theo chính sách  $\pi$ .
- Cập nhật  $Q$ -value thông qua phương trình Bellman:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

# DQN

- DQN thay thế  $Q$ -table bằng một mạng nơ-ron  $Q_{\theta}(s, a)$  có tham số  $\theta$ . Phương trình Bellman được ước lượng thông qua tối ưu hóa hàm mất mát:

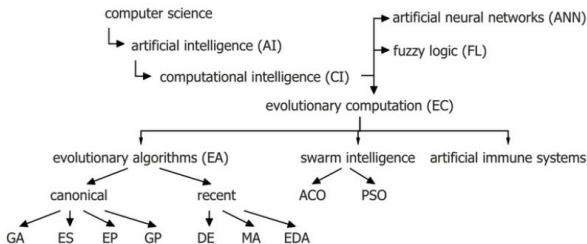
$$\mathcal{L}(\theta) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[ \left( r + \gamma \max_{a'} Q_{\theta^-}(s', a') - Q_{\theta}(s, a) \right)^2 \right]$$

- $Q_{\theta}$ : Mạng chính (online network).
- $Q_{\theta^-}$ : Mạng mục tiêu (target network), cập nhật định kỳ từ mạng chính.
- $\mathcal{D}$ : Bộ nhớ kinh nghiệm (Replay Buffer) lưu trữ các trạng thái, hành động, phần thưởng và trạng thái tiếp theo để huấn luyện.

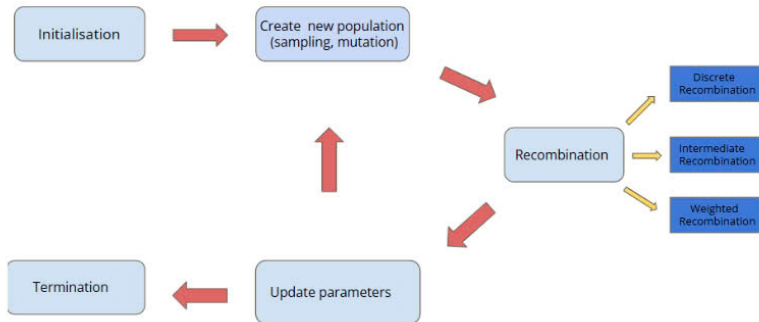


# Evolution Strategies (ES)

- Evolution Strategy (ES) là một thuật toán thuộc nhóm thuật toán tiến hóa (EAs)
- ES thường sử dụng cho các biểu diễn số thực, với phép biến đổi (variation) chủ yếu là đột biến (mutation)



# Evolution Strategies (ES)



Hình: Ý tưởng cơ bản của Evolution Strategies

# Novelty Search (NS)

- NS được lấy cảm hứng từ xu hướng tự nhiên hướng tới sự đa dạng.
- NS trong thuật toán tiến hóa là một phương pháp được sử dụng để tìm kiếm các giải pháp mới mẻ
- NS khuyến khích sự đa dạng trong hành vi, ưu tiên những giải pháp khác biệt so với các giải pháp trước đó thay vì tập trung vào một mục tiêu cố định.

# NS - working ideas

- Behavior Characterization

- Mỗi chính sách  $\pi$  được gán một đặc trưng hành vi (BC) là  $b(\pi)$ , mô tả cách nó hoạt động trong môi trường

# NS - working ideas

- Behavior Characterization
  - Mỗi chính sách  $\pi$  được gán một đặc trưng hành vi (BC) là  $b(\pi)$ , mô tả cách nó hoạt động trong môi trường
- Behavior Archive
  - Trong quá trình huấn luyện, mỗi một BC được thêm vào archive A

# NS - working ideas

- Behavior Characterization

- Mỗi chính sách  $\pi$  được gán một đặc trưng hành vi (BC) là  $b(\pi)$ , mô tả cách nó hoạt động trong môi trường

- Behavior Archive

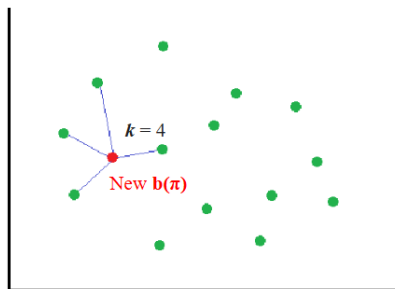
- Trong quá trình huấn luyện, mỗi một BC được thêm vào archive A

- Tính toán mức độ mới lạ

- Novelty score của một chính sách  $\pi_\theta$  được đo bằng khoảng cách giữa  $b(\pi_\theta)$  và các BC trước đó trong archive.
- Chọn k hàng xóm gần nhất của  $b(\pi)$  trong Archive A và tính trung bình khoảng cách L2 giữa  $b(\pi_\theta)$  và các hàng xóm gần nhất

# NS - working ideas

$\pi$   $\longrightarrow$  Compute  $\mathbf{b}(\pi)$   $\longrightarrow$  Add  $\mathbf{b}(\pi)$  into archive



Calculate average  $L2$  distance to  $k$  nearest neighbors

Hình: Minh họa cách thuật toán NS hoạt động

# Quality Diversity - QD

- QD trong thuật toán tiến hóa là một phương pháp để tìm kiếm một tập hợp đa dạng các giải pháp có chất lượng cao thay vì chỉ tập trung vào một giải pháp tối ưu duy nhất.



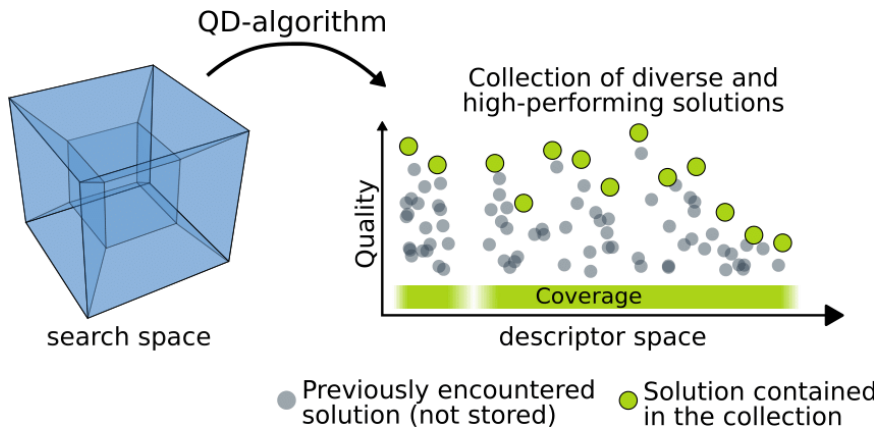
# Quality Diversity - QD

- QD trong thuật toán tiến hóa là một phương pháp để tìm kiếm một tập hợp đa dạng các giải pháp có chất lượng cao thay vì chỉ tập trung vào một giải pháp tối ưu duy nhất.
- Ý tưởng chính
  - **Đa dạng:** Chia không gian tìm kiếm thành các khu vực khác nhau, mỗi khu vực đại diện cho một dạng hành vi hoặc đặc điểm cụ thể của giải pháp
  - **Chất lượng cao:** Mỗi cá thể cần phải đạt được chất lượng cao theo một hàm mục tiêu

# Quality Diversity - QD

- QD trong thuật toán tiến hóa là một phương pháp để tìm kiếm một tập hợp đa dạng các giải pháp có chất lượng cao thay vì chỉ tập trung vào một giải pháp tối ưu duy nhất.
- Ý tưởng chính
  - **Đa dạng:** Chia không gian tìm kiếm thành các khu vực khác nhau, mỗi khu vực đại diện cho một dạng hành vi hoặc đặc điểm cụ thể của giải pháp
  - **Chất lượng cao:** Mỗi cá thể cần phải đạt được chất lượng cao theo một hàm mục tiêu
- *Không triển khai trực tiếp QD:* sử dụng NS kết hợp với thông tin từ fitness value (độ thích nghi)

# Quality Diversity - QD



Hình: Minh họa Archive trong QD

# Mục lục

- 1 Giới thiệu đề tài
- 2 Cơ sở lý thuyết
  - DQN
  - ES
  - NS
  - QD
- 3 Phương pháp tiếp cận
  - NS-ES
  - QD-ES algo: NSR-ES
  - NS-ES customized
  - NSR-ES customized
- 4 Experiment results

# NS-ES

## Algorithm 1 NS-ES

---

```

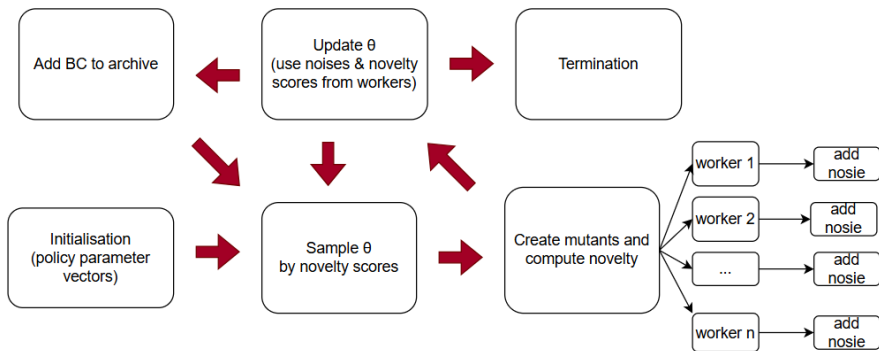
1: Input: learning rate  $\alpha$ , noise standard deviation  $\sigma$ , number of policies to maintain  $M$ , iterations  $T$ , behavior characterization  $b(\pi_\theta)$ 
2: Initialize:  $M$  randomly initialized policy parameter vectors  $\{\theta_0^1, \theta_0^2, \dots, \theta_0^M\}$ , archive  $A$ , number of workers  $n$ 
3: for  $j = 1$  to  $M$  do
4:   Compute  $b(\pi_{\theta_0^j})$ 
5:   Add  $b(\pi_{\theta_0^j})$  to  $A$ 
6: end for
7: for  $t = 0$  to  $T - 1$  do
8:   Sample  $\theta_t^m$  from  $\{\theta_t^1, \theta_t^2, \dots, \theta_t^M\}$  via eq 1
9:   for  $i = 1$  to  $n$  do
10:    Sample  $\epsilon_i \sim \mathcal{N}(0, \sigma^2 I)$ 
11:    Compute  $\theta_t^{i,m} = \theta_t^m + \epsilon_i$ 
12:    Compute  $b(\pi_{\theta_t^{i,m}})$ 
13:    Compute  $N_i = N(\theta_t^{i,m}, A)$ 
14:    Send  $N_i$  from each worker to coordinator
15:   end for
16:   Set  $\theta_{t+1}^m = \theta_t^m + \alpha \frac{1}{n\sigma} \sum_{i=1}^n N_i \epsilon_i$ 
17:   Compute  $b(\pi_{\theta_{t+1}^m})$ 
18:   Add  $b(\pi_{\theta_{t+1}^m})$  to  $A$ 
19: end for

```

---

## Hình: Thuật toán NS-ES

# NS-ES - working ideas



# NS-ES - working ideas

- Khởi tạo

# NS-ES - working ideas

- Khởi tạo

- Tạo M vector tham số ngẫu nhiên  $\theta_0^1, \theta_0^2, \dots, \theta_0^M$  đại diện cho M agent.
- Với mỗi agent, tính giá trị đặc trưng hành vi  $b(\pi_{\theta_0^j})$  và thêm vào tập lưu trữ A.



# NS-ES - working ideas

- Khởi tạo
  - Tạo M vector tham số ngẫu nhiên  $\theta_0^1, \theta_0^2, \dots, \theta_0^M$  đại diện cho M agent.
  - Với mỗi agent, tính giá trị đặc trưng hành vi  $b(\pi_{\theta_0^j})$  và thêm vào tập lưu trữ A.
- Với mỗi vòng lặp t:

# NS-ES - working ideas

- Khởi tạo
  - Tạo M vector tham số ngẫu nhiên  $\theta_0^1, \theta_0^2, \dots, \theta_0^M$  đại diện cho M agent.
  - Với mỗi agent, tính giá trị đặc trưng hành vi  $b(\pi_{\theta_0^j})$  và thêm vào tập lưu trữ A.
- Với mỗi vòng lặp t:
  - Tính xác suất được chọn của mỗi  $\theta_i$  bằng cách lấy độ mới lạ của  $b(\pi_{\theta_i})$  đó chia cho tổng độ mới lạ của tất cả các BC trong archive A.
  - Chọn ngẫu nhiên một  $\theta$  dựa theo các phân phối xác suất đã tính.

# NS-ES - working ideas

- Khởi tạo

- Tạo M vector tham số ngẫu nhiên  $\theta_0^1, \theta_0^2, \dots, \theta_0^M$  đại diện cho M agent.
- Với mỗi agent, tính giá trị đặc trưng hành vi  $b(\pi_{\theta_0^j})$  và thêm vào tập lưu trữ A.

- Với mỗi vòng lặp t:

- Tính xác suất được chọn của mỗi  $\theta_i$  bằng cách lấy độ mới lạ của  $b(\pi_{\theta_i})$  đó chia cho tổng độ mới lạ của tất cả các BC trong archive A.
- Chọn ngẫu nhiên một  $\theta$  dựa theo các phân phối xác suất đã tính.
- Gửi  $\theta$  đến mỗi worker. Thêm nhiều và tính toán  $b(\pi_{\theta})$  cùng với  $N(\theta, A)$

# NS-ES - working ideas

- Khởi tạo

- Tạo M vector tham số ngẫu nhiên  $\theta_0^1, \theta_0^2, \dots, \theta_0^M$  đại diện cho M agent.
- Với mỗi agent, tính giá trị đặc trưng hành vi  $b(\pi_{\theta^j})$  và thêm vào tập lưu trữ A.

- Với mỗi vòng lặp t:

- Tính xác suất được chọn của mỗi  $\theta_i$  bằng cách lấy độ mới lạ của  $b(\pi_{\theta_i})$  đó chia cho tổng độ mới lạ của tất cả các BC trong archive A.
- Chọn ngẫu nhiên một  $\theta$  dựa theo các phân phối xác suất đã tính.
- Gửi  $\theta$  đến mỗi worker. Thêm nhiễu và tính toán  $b(\pi_{\theta})$  cùng với  $N(\theta, A)$
- Lấy các kết quả từ các worker và tiến hành cập nhật  $\theta$

# NS-ES - Experiment configs

- Policy: Mạng nơ-ron (trạng thái  $\rightarrow$  xác suất thực hiện hành động).
- Policy parameters: trọng số của mạng
- Population: Tập hợp các policy ( $M = 3$ )
- n\_workers: tận dụng khả năng tính toán song song ( $n = 5$ )
- Behavior characterization: Tính trung bình các state
- Behavior archive: list các BC với môi trường CartPole-v0 và MountainCar-v0 / deque(maxlen = 500) (BreakoutNoFrameskip-v4)
- Novelty score (N function): Trung bình euclidean distance của BC hiện tại với  $k = 5$  BC gần nhất
- Sample: Theo phân phối xác suất được tính từ novelty score
- Mutation: Gửi policy đến từng worker  $\rightarrow$  Thêm nhiễu, tính toán BC và novelty score (N)  $\rightarrow$  Tổng hợp kết quả  $\rightarrow$  Thực hiện bước thêm nhiễu cuối cùng cho policy

# NS-ES: problem

- NS-ES giúp agent thoát khỏi các điểm tối ưu cục bộ trong hàm phần thưởng.

# NS-ES: problem

- NS-ES giúp agent thoát khỏi các điểm tối ưu cục bộ trong hàm phần thưởng.
- **Vấn đề:** NS-ES bỏ qua hoàn toàn tín hiệu phần thưởng trong khi đây là một thông tin hữu ích.

→ Ảnh hưởng đến hiệu suất của thuật toán

# NS-ES: problem

- NS-ES giúp agent thoát khỏi các điểm tối ưu cục bộ trong hàm phần thưởng.
- **Vấn đề:** NS-ES bỏ qua hoàn toàn tín hiệu phần thưởng trong khi đây là một thông tin hữu ích.

→ Ảnh hưởng đến hiệu suất của thuật toán

Một biến thể của NS-ES sử dụng thông tin phần thưởng đó là NSR-ES.

- NSR-ES về cơ bản giống với NS-ES, điểm khác biệt duy nhất là NSR-ES sử dụng thêm “tín hiệu phần thưởng” để cập nhật trọng số



# NSR-ES

---

## Algorithm 2 NSR-ES

---

```

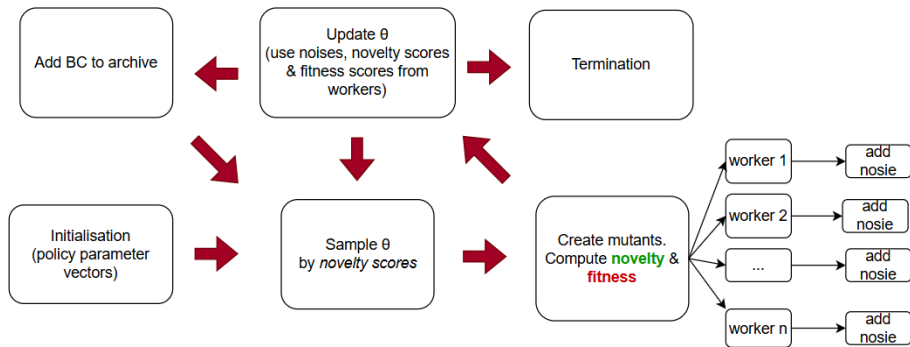
1: Input: learning rate  $\alpha$ , noise standard deviation  $\sigma$ , number of policies to maintain  $M$ , iterations  $T$ , behavior characterization  $b(\pi_\theta)$ 
2: Initialize:  $M$  sets of randomly initialized policy parameters  $\{\theta_0^1, \theta_0^2, \dots, \theta_0^M\}$ , archive  $A$ , number of workers  $n$ 
3: for  $j = 1$  to  $M$  do
4:   Compute  $b(\pi_{\theta_0^j})$ 
5:   Add  $b(\pi_{\theta_0^j})$  to  $A$ 
6: end for
7: for  $t = 0$  to  $T - 1$  do
8:   Sample  $\theta_t^m$  from  $\{\theta_t^0, \theta_t^1, \dots, \theta_t^M\}$  via eq. 1
9:   for  $i = 1$  to  $n$  do
10:    Sample  $\epsilon_i \sim \mathcal{N}(0, \sigma^2 I)$ 
11:    Compute  $\theta_t^{i,m} = \theta_t^m + \epsilon_i$ 
12:    Compute  $b(\pi_{\theta_t^{i,m}})$ 
13:    Compute  $N_i = N(\theta_t^{i,m}, A)$ 
14:    Compute  $F_i = f(\theta_t^{i,m})$ 
15:    Send  $N_i$  and  $F_i$  from each worker to coordinator
16:   end for
17:   Set  $\theta_{t+1}^m = \theta_t^m + \alpha \frac{1}{n\sigma} \sum_{i=1}^n \frac{N_i + F_i}{2} \epsilon_i$ 
18:   Compute  $b(\pi_{\theta_{t+1}^m})$ 
19:   Add  $b(\pi_{\theta_{t+1}^m})$  to  $A$ 
20: end for

```

---

Hình: Thuật toán NSR-ES

# NSR-ES



## NS-ESc

---

**Algorithm 3** NS-ESc

---

**Input:** noise standard deviation  $\sigma$ , population size  $M$ , max generations  $G$ , behavior characterization  $b(\pi_\theta)$ , number of individuals to remain  $k$

**Initialize:** a population of  $M$  random policy parameters  $pop = [\theta_0^1; \theta_0^2; \dots; \theta_0^M]$ , archive  $A$

**for**  $j = 1$  **to**  $M$  **do**

Compute  $b(\pi_{\theta_j^0})$ ; Add  $b(\pi_{\theta_j^0})$  to  $A$

**end for**

**for**  $t = 0$  **to**  $G - 1$  **do**

$next\_pop = []$ ,  $mutant\_pop = []$

**for**  $i = 1$  **to**  $M$  **do**

Compute  $b(\pi_{\theta_t^i})$

Compute  $N_t = N(\theta_t^i, A)$

**end for**

**Ranking** current population by novelty score:  $ranked\_pop = \text{ranking}(pop)$

**Select** top  $k$   $\theta_t^i$  from  $ranked\_pop$ , add to  $next\_pop$  and  $mutant\_pop$

**While**  $mutant\_population\ size < M - k$  **do**

Randomly select  $\theta_t^i$  from the rest of  $ranked\_pop$

Add  $\theta_t^i$  to  $mutant\_pop$

**end while**

**for**  $i = 1$  **to**  $M - k$  **do**

Select  $\theta_t^i$  from  $mutant\_pop$

Sample  $\varepsilon_t \sim N(0, \sigma^2 I)$

Compute  $\theta_t^i = \theta_t^i + \varepsilon_t$ ; Add  $\pi_{\theta_t^i}$  to  $next\_pop$

Compute  $b(\pi_{\theta_t^i})$ ; Add  $b(\pi_{\theta_t^i})$  to  $A$

**end for**

Update  $pop = next\_pop$

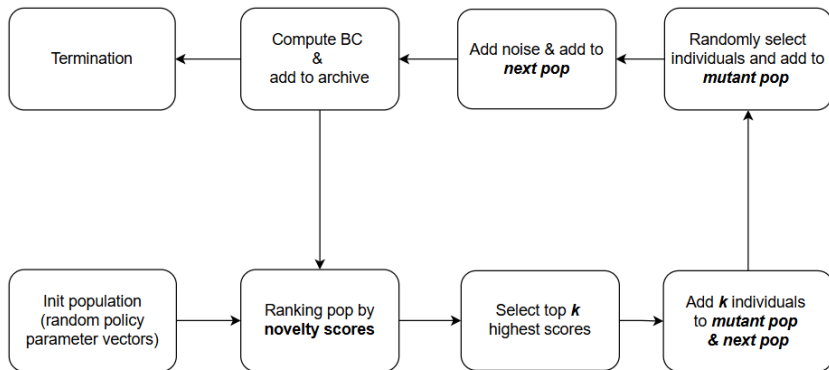
**end for**

# NS-ESc - working ideas

Điểm khác so với NS-ES:

- NS-ES huấn luyện qua từng vòng lặp (iteration), trong khi NS-ESc huấn luyện theo thế hệ.
- Ở mỗi thế hệ,  $k$  cá thể có *novelty score* **cao nhất** sẽ được giữ lại và đưa vào **next\_pop** (quần thể ở thế hệ kế tiếp).
- **mutant\_pop** được xác định bằng cách lấy  $k$  cá thể có novelty score cao nhất và chọn ngẫu nhiên trong số các cá thể còn lại cho đến khi đủ  $M - k$  cá thể.
- Phép đột biến được thực hiện trên **mutant\_pop** và các cá thể sinh ra được bổ sung vào **next\_pop** để làm quần thể cho thế hệ kế tiếp.

## NS-ESc



## NSR-ESc

**Algorithm 4** NSR-ESc

**Input:** noise standard deviation  $\sigma$ , population size  $M$ , max generations  $G$ , behavior characterization  $b(\pi_\theta)$ , number of individuals to remain  $k$

**Initialize:** a population of  $M$  random policy parameters  $pop = [\theta_0^1; \theta_0^2; \dots; \theta_0^M]$ , archive  $\mathcal{A}$

**for**  $j = 1$  **to**  $M$  **do**

Compute  $b(\pi_{\theta_j^j})$ ; Add  $b(\pi_{\theta_j^j})$  to  $\mathcal{A}$

**end for**

**for**  $t = 0$  **to**  $G - 1$  **do**

$next\_pop = []$ ,  $mutant\_pop = []$

**for**  $i = 1$  **to**  $M$  **do**

Compute  $b(\pi_{\theta_i^t})$

Compute  $N_t = N(\theta_i^t, \mathcal{A})$ ,  $F_t = F(\theta_i^t, \mathcal{A})$ ,  $score_t = \frac{N_t + F_t}{2}$

**end for**

**Ranking** current population by score:  $ranked\_pop = \text{ranking}(pop)$

**Select** top  $k$   $\theta_i^t$  from  $ranked\_pop$ , add to  $next\_pop$  and  $mutant\_pop$

**While**  $mutant\_population$  size  $< M - k$  **do**

Randomly select  $\theta_i^t$  from the rest of  $ranked\_pop$

Add  $\theta_i^t$  to  $mutant\_pop$

**end while**

**for**  $i = 1$  **to**  $M - k$  **do**

Select  $\theta_i^t$  from  $mutant\_pop$

Sample  $\varepsilon_i \sim N(0, \sigma^2 I)$

Compute  $\theta_i^t = \theta_i^t + \varepsilon_i$ ; Add  $\pi_{\theta_i^t}$  to  $next\_pop$

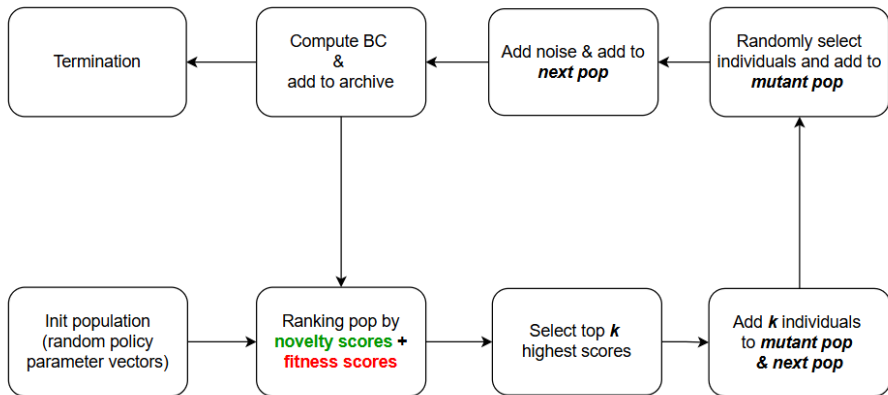
Compute  $b(\pi_{\theta_i^t})$ ; Add  $b(\pi_{\theta_i^t})$  to  $\mathcal{A}$

**end for**

Update  $pop = next\_pop$

**end for**

# NSR-ESc



# Mục lục

- 1 Giới thiệu đề tài
- 2 Cơ sở lý thuyết
  - DQN
  - ES
  - NS
  - QD
- 3 Phương pháp tiếp cận
  - NS-ES
  - QD-ES algo: NSR-ES
  - NS-ES customized
  - NSR-ES customized
- 4 Experiment results



# Environment

- CartPole-v0
- MountainCar-v0
- BreakoutNoFrameskip-v4

# Experiment results

Bảng: Best result for each method

Environment	Method	Best score	Game play
CartPole-v0	DQN	500	<a href="#"><u>video</u></a>
	NS-ES	500	<a href="#"><u>video</u></a>
	NSR-ES	500	<a href="#"><u>video</u></a>
	NS-ESc	500	<a href="#"><u>video</u></a>
	NSR-ESc	500	<a href="#"><u>video</u></a>

# Experiment results

Bảng: Best result for each method

Environment	Method	Best score	Complete	Game play
MountainCar-v0	DQN	-91	YES	<a href="#"><u>video</u></a>
	NS-ES	-129.33	YES	<a href="#"><u>video</u></a>
	NSR-ES	-183	YES	<a href="#"><u>video</u></a>
	NS-ES <sub>c</sub>	-115.66	YES	<a href="#"><u>video</u></a>
	NSR-ES <sub>c</sub>	-109.66	YES	<a href="#"><u>video</u></a>

# Experiment results

Bảng: Best result for each method

Environment	Method	Best score	Game play
BreakoutNoFrameskip-v4	DQN	38	<a href="#"><u>video</u></a>
	NS-ES	16	<a href="#"><u>video</u></a>
	NSR-ES	21	<a href="#"><u>video</u></a>
	NS-ESc	13	<a href="#"><u>video</u></a>
	NSR-ESc	22	<a href="#"><u>video</u></a>

# Comments

- Trên cả 3 môi trường, DQN cho kết quả tốt nhất
- Trong số 4 thuật toán ES, NSR-ESc cho kết quả tốt nhất
- Do cách cài đặt của các thuật toán khác nhau nên bảng kết quả không đưa ra thời gian thực thi, tuy nhiên khi thực nghiệm, nhóm có một số nhận xét như sau:
  - Đối với môi trường CartPole, mặc dù kết quả sau cùng của 5 thuật toán đều bằng nhau nhưng trong quá trình thực nghiệm, thuật toán NSR-ESc có tốc độ hội tụ nhanh nhất
  - Đối với môi trường MountainCar, thuật toán DQN hội tụ nhanh hơn so với các thuật toán khác. Ngoài ra, trong số 4 thuật toán ES thì thuật toán NS-ES có thời gian thực thi lâu nhất
- Mặc dù kết quả đạt được không tốt bằng DQN, tuy nhiên các thuật toán ES đã cho thấy tiềm năng ứng dụng của nó khi giải quyết các bài toán về DRL