

# Learning to Integrate Exploration Strategies for Reinforcement Learning via an Option Framework

(LESSON)

ICML 2023

Woojun Kim\*, Jeonghye Kim\*, Youngchul Sung  
(\* denotes equal contribution)



# Previous RL exploration method

Exploration refers to the method of efficiently exploring the environment within limited resources and time in order to gather crucial information for the agent's decision-making process.

RND

Temporally extended  
exploration

epsilon-greedy

## Task specific

There is yet no single exploration method that is shown to be universally effective across all tasks.

## Invariable during training period

The required exploration strategy can vary over time during the training period within a given task, but existing methods do not reflect these characteristics.

# LESSON - A Unified Framework for Multiple Exploration Strategies

LESSON is a unified framework that utilizes the **option-critic architecture**<sup>[1]</sup> and an off-policy structure to **effectively integrate multiple exploration strategies** for adaptive exploration strategy selection during the learning phase.

## Option-critic architecture

**Options (Temporally-extended actions)** : a form of action generalization that captures high-level behavior by combining multiple sub-actions.

**Call-and-Return Option Execution Model** : a hierarchical control framework that improves the efficiency and effectiveness of the agent's navigation by breaking down complex tasks into a series of options.

- **Call** : The agent **selects an option** that represents a sub-goal or desired behavior based on its current state.
- **Return** : The agent **executes the selected option until its termination condition is met**.

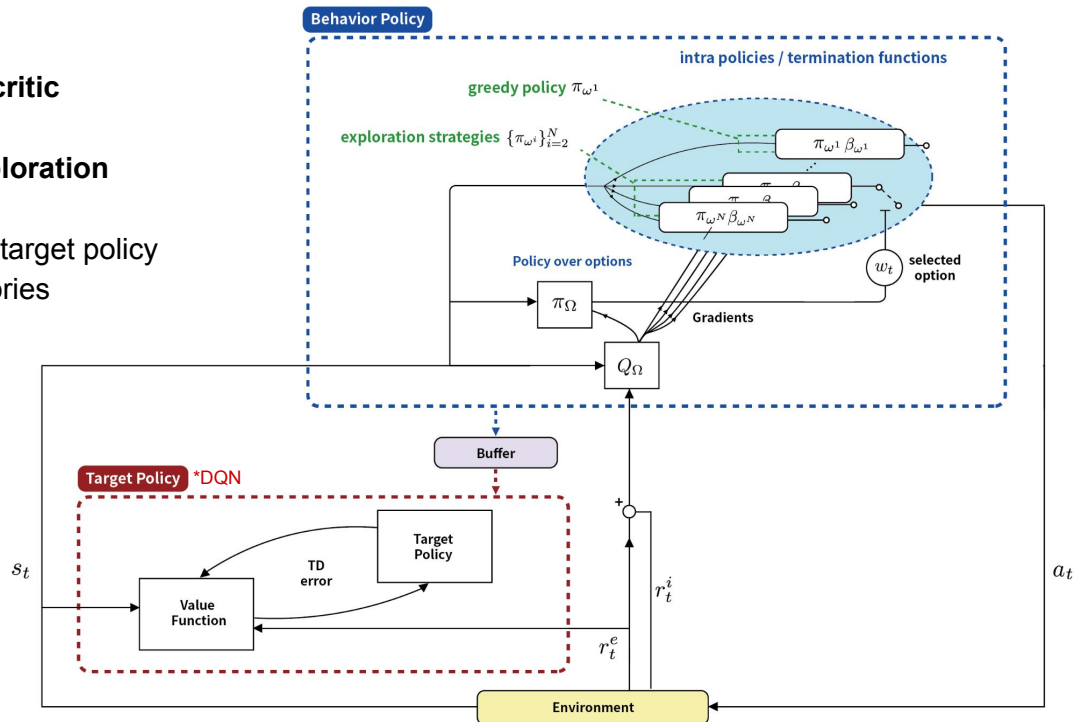
# Method - Behavior Policy Construction via Option-Critic

## Key ideas of LESSON are

- (1) to replace the behavior policy with an option-critic architecture whose (N : number of option)
- (2) intra-policies are defined by N-1 component exploration strategies and the greedy policy, and
- (3) to train both the option-critic architecture and the target policy with two different objectives based on the trajectories generated by the option-critic architecture

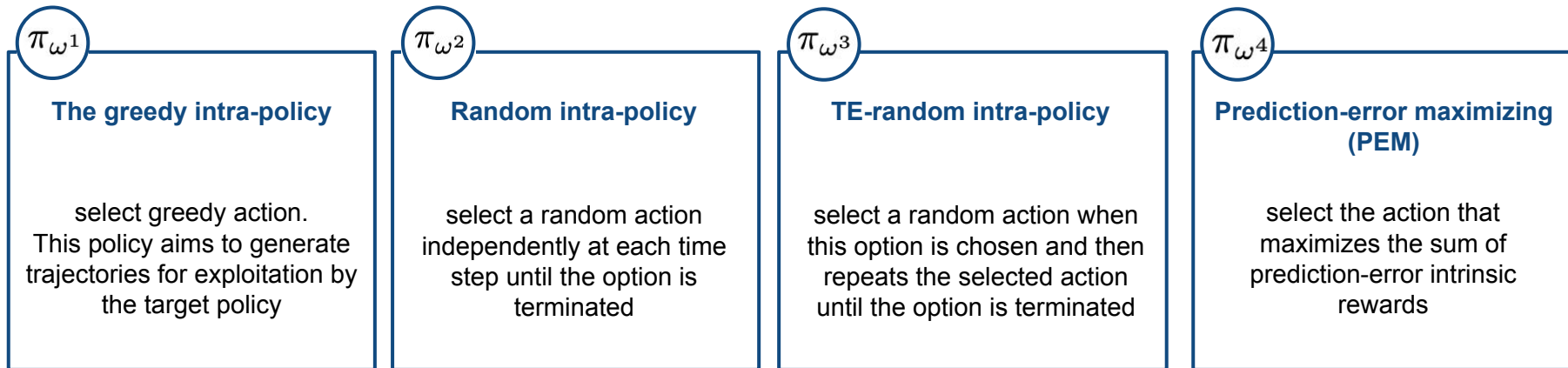
## Important factor for exploration-exploitation trade-off in sampling

- Inclusion of the greedy policy
- Design of the objective function



# Method - Behavior Policy Construction via Option-Critic

## Design of intra-policies



$$\pi_{\omega^1} + \pi_{\omega^2} = \epsilon\text{-greedy}$$

$$\pi_{\omega^1} + \pi_{\omega^3} = \epsilon z\text{-greedy}$$

$$\pi_{\omega^1} + \pi_{\omega^4} = \text{RND}$$

# Method - Learning Option-Critic

With the predefined intra-policies, we need to learn the **option selection policy**  $\pi_\Omega$  and the **termination functions**  $\{\beta_\omega\}$

**Objective function for behavior policy :** 
$$J(\pi_\Omega, \{\beta_\omega\}) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \left( \underbrace{r_{t+1}^e}_{\text{exploitation}} + \alpha \underbrace{r_{t+1}^i}_{\text{exploration}} \right) \right]$$

behavior policy should not only sample for exploration but also for exploitation for a trade-off between these two

**Option-value function :** 
$$Q_\Omega(s_t, \omega_t) = \mathbb{E} \left[ \sum_{l=t}^{\infty} \gamma^{l-t} (r_{l+1}^e + \alpha r_{l+1}^i) \mid s_t, \omega_t \right]$$

**Learning Option-Value Function**

$$\mathcal{L}(\theta_\Omega) = \mathbb{E}_{(s_t, \omega_t, r_t^e + \alpha r_t^i, s_{t+1}) \sim \mathcal{D}} \left[ (y_t - Q_\Omega(s_t, \omega_t; \theta_\Omega))^2 \right], \text{ where}$$
$$y_t = r_t^e + \alpha r_t^i + \gamma \left( (1 - \beta_{\omega_t}(s_{t+1})) Q_\Omega(s_{t+1}, \omega_t; \theta_\Omega^-) + \beta_{\omega_t}(s_{t+1}) \max_{\omega'} Q_\Omega(s_{t+1}, \omega'; \theta_\Omega^-) \right)$$

**Learning Termination Functions**

$$\frac{\partial Q_\Omega}{\partial \theta_{\beta_\omega}} = -\mathbb{E} \left[ \nabla_{\theta_{\beta_\omega}} \beta_\omega(s_{t+1}; \theta_{\beta_\omega}) (Q_\Omega(s_{t+1}, \omega_t) - \max_{\omega} Q_\Omega(s_{t+1}, \omega)) \right]$$

# Experiments - Performance Comparison

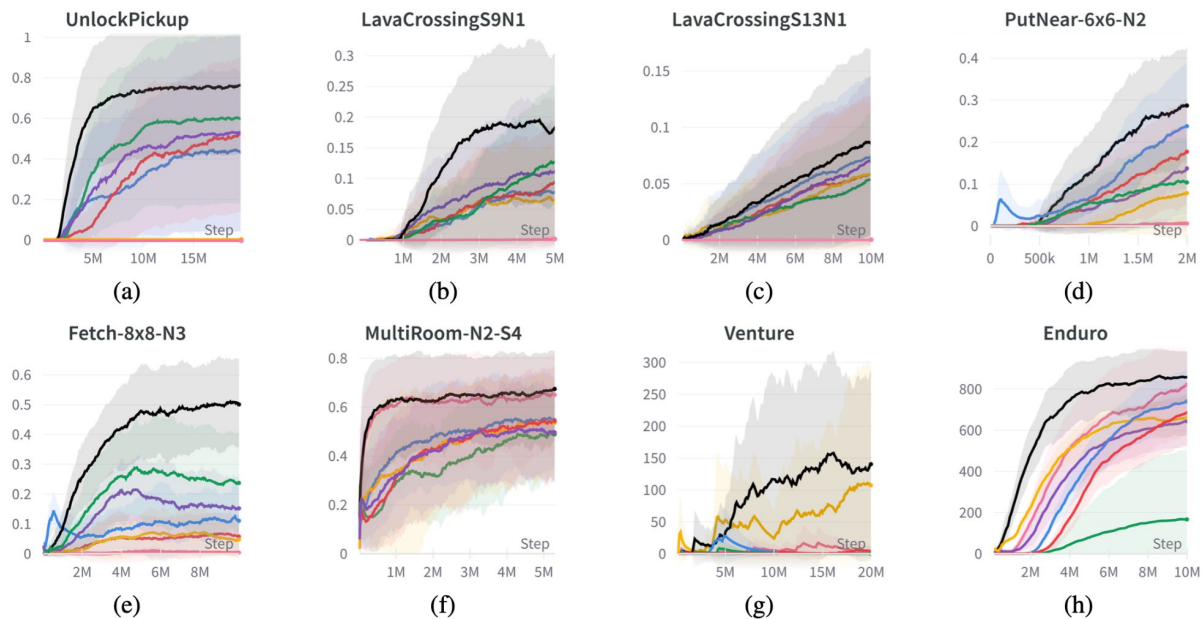


## Baselines

- $\epsilon$ -greedy (vanilla DQN)
- two simple DQN variants
  - $\epsilon z$ -greedy,  $\epsilon r$ -greedy
- RND-based DQN
- Equal weight combining (EWC)
- $\epsilon$ -BMC, which learns  $\epsilon$

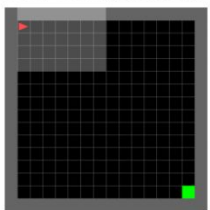
## Environments

- 14 MiniGrid environments
- 4 Atari environments

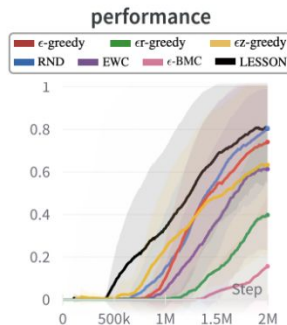


# Experiments - Exploration Behavior Analysis

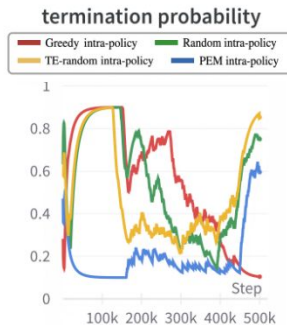
Empty-16x16 (bottom-goal)



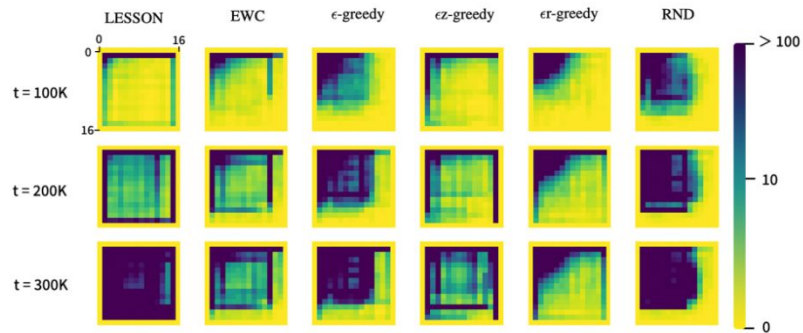
(a)



(b)



(c)



(d)

- Adaptive exploration-exploitation trade-off achieved through adaptive selection of intra-policies over time. (c)
- **LESSON combines visitation patterns of RND and ez-greedy to cover entire state space. (d)**



# The LESSON We Learned

- It is crucial to **adaptively select exploration strategies** according to environmental characteristics and the learning phase to achieve an **efficient exploration-exploitation balance**.
- In order to create a unified framework, we propose to incorporate the option-critic architecture with intra-policies consisting of a **greedy policy and a set of exploration strategies**, and meticulously **designing the objective function**.

**Thank you for listening**

