# Paper Critique

Shuvrajeet Das, DA24D402

**Course:** DA7400, Fall 2024, IITM
**Paper:** [Deep Laplacian-based Options for Temporally-Extended Exploration]
**Date:** [04-09-2024]

Make sure your critique Address the following points:
1. The problem the paper is trying to address
2. Key contributions of the paper
3. Proposed algorithm/framework
4. How the proposed algorithm addressed the described problem
Note: Be concise with your explanations. Unnecessary verbosity will be penalized. Please don't exceed 2 pages.

---

## 1 The problem the paper is trying to address

The paper addresses the problem of exploration in reinforcement learning (RL), particularly in complex environments. In RL, agents must explore their environments to gather information that allows them to learn effective policies for maximizing rewards. However, exploration is challenging because the agent needs to balance between exploiting known rewarding actions and exploring new actions that might lead to better long-term outcomes.

## 2 Key contributions of the paper

1. **Deep Covering Eigenoptions (DCEO):** Introduced DCEO, a novel deep RL algorithm extending Laplacian-based options to high-dimensional environments using neural networks.

2. **Scalable Approximation of Laplacian Eigenfunctions:** DCEO utilizes the generalized Laplacian objective:

$$\min_{f_1,\ldots,f_d} \sum_{i=1}^{d} c_i f_i^\top L f_i \quad \text{s.t.} \quad f_i^\top f_j = \delta_{ij}$$

   where $L$ is the graph Laplacian, $f_i$ are the eigenfunction approximations, and $c_i$ are coefficients. This avoids costly eigendecomposition.

3. **Fully Online Option Discovery and Learning:** DCEO integrates option discovery and reward maximization into a continuous process, enabling adaptability in non-stationary environments.

4. **Generalization to Diverse Environments:** Demonstrated the effectiveness of DCEO across various environments, including pixel-based and 3D navigation tasks.

5. **Temporally-Extended Exploration:** DCEO encourages structured exploration by maximizing intrinsic rewards derived from eigenfunctions:

$$r_{f_i}(s, s') = f_i(s') - f_i(s)$$

   leading to improved state coverage.

# 3 Proposed algorithm/framework

---

**Algorithm 1** Fully Online DCEO Algorithm

---

1: **for** $i = 1$ **to** $T$ **do**
2:     **if** $i = 1$ **or** $i \perp \ell$ **then**
3:        $\tau \leftarrow$ **True**             $\triangleright$ Option termination
4:        $o \leftarrow -1$             $\triangleright$ No active option
5:        Reset environment and observe state $s$
6:     **end if**
7:     **if** $\tau$ **then**
8:        $\tau \sim U(0,1) < 1/D$    ✓ $\tau$
9:     **end if**
10:    **if** $\tau$ **then**
11:       **if** $U(0,1) < \theta$ then **then**
12:          **if** $U(0,1) < \mu$ then **then**
13:             $o \leftarrow U(0);\ \tau \leftarrow$ **False**; $\alpha \sim \pi_o(\cdot|s)$
14:          **else**
15:             $o \leftarrow -1;\ \tau \leftarrow$ **True**; $\alpha \sim U(\mathcal{A})$
16:          **end if**
17:       **else**
18:          $\alpha \leftarrow \max_{\alpha \in \mathcal{A}} Q(s, \alpha)$
19:       **end if**
20:     **else**
21:       $\alpha \sim \pi_o(\cdot|s)$
22:     **end if**
23:    Execute $\alpha$, observe $r$, $s'$, $l \perp \ell$        $\triangleright$ $\ell$ is episode termination
24:    Store transition $(s, \alpha, r, s')$ in buffer $B$; $s \leftarrow s'$
25:    Sample a minibatch of transitions $(s_j, \alpha_j, r_j, s_{j+1})$
26:    Train each option with intrinsic reward
27:    Minimize the generalized Laplacian
28:    Train main learner on extrinsic reward
29: **end for**

---

## 3.1 Scalability with Deep Function Approximation:

**Approximate Eigenfunctions:** The algorithm uses neural networks to approximate the eigenfunctions of the graph Laplacian instead of relying on costly eigendecomposition. This allows the method to scale to large, complex environments, including those with pixel-based inputs.

## 3.2 Online and Continuous Option Discovery:

**Single Continuous Learning Cycle:** Unlike previous methods that required a separate option discovery phase, DCEO integrates option discovery and reward maximization into a single, continuous learning process. This means that the agent is always learning and refining its options while simultaneously learning to maximize rewards, making it adaptable to dynamic and non-stationary environments.

## 3.3 Generalization Across Tasks and Environments:

**Effective in Diverse and Complex Environments:** The paper demonstrates that DCEO is effective in a variety of environments, from simple grid worlds to complex 3D navigation tasks and Atari games