

## ۱. BBO برای Policy Optimization در RL

در الگوریتم‌های RL مبتنی بر پالیسی مانند PPO یا REINFORCE، هدف بهینه‌سازی یک تابع پالیسی است.

- می‌توان از BBO برای یافتن بهترین پارامترهای پالیسی استفاده کرد.
- پالیسی‌های مختلف به‌عنوان جزایر (habitats) در BBO در نظر گرفته می‌شوند.
- عملکرد پالیسی‌ها مقدار reward معادل شاخص برازندگی (fitness) است.
- مهاجرت بین پالیسی‌ها باعث تبادل استراتژی‌ها می‌شود.

## ۲. استفاده از RL برای هدایت جستجوی BBO

در این رویکرد:

- RL به‌عنوان meta-controller عمل می‌کند که جهت جستجوی الگوریتم BBO را هدایت می‌کند (منظور از meta-controller یا فراکنترل‌کننده، سیستمی است که نقش کنترل‌کننده بر رفتار یک الگوریتم دیگر (در این کاربرد، الگوریتم BBO) را دارد).
- RL تعیین می‌کند که چه پارامترهایی از الگوریتم BBO مانند نرخ مهاجرت، نرخ جهش و ... تغییر کنند.
- با این روش BBO به صورت تطبیقی و پویا تنظیم می‌شود.

## ۳. BBO به عنوان جایگزین روش‌های gradient-based در RL

در مسائل غیرقابل تفکیک (non-differentiable)، استفاده از مشتق‌گیری امکان‌پذیر نیست.

- ✓ در این موارد، می‌توان از BBO به‌عنوان جایگزین روش‌های گرادیانی استفاده کرد تا پالیسی بدون نیاز به گرادیان یاد بگیرد.

## ۴. آموزش اولیه RL (warm-start) با استفاده از BBO

- ابتدا الگوریتم BBO برای تولید پالیسی‌های اولیه قوی استفاده می‌شود.
- سپس RL از این پالیسی‌ها به‌عنوان نقطه شروع استفاده کرده و آن‌ها را بیشتر بهبود می‌دهد.

## ۵. BBO برای بهینه‌سازی معماری یا پارامترهای شبکه RL

در الگوریتم‌هایی مثل DQN یا Actor-Critic، می‌توان از BBO برای:

- انتخاب بهترین ساختار شبکه (تعداد لایه‌ها، نورون‌ها)
- تنظیم هاپرپارامترهای RL مثل learning rate، discount factor و ...