

# SELF-TRANSFER LEARNING

Hélène Plisnier,  
Denis Steckelmacher,  
Ann Nowé

## TRANSFER LEARNING IN RL

- ▶ two tasks: a **source** task and a **target** task
- ▶ goal: make learning the target task **faster** than when learning it from scratch

## TRANSFER LEARNING IN RL

- ▶ the source is generally somewhat **different** than the target
  - ▶ examples: goal, dynamics
- ▶ however, no real consensus on **how different**
  - ▶ at least same robot ?

## TRANSFER LEARNING IN RL

- ▶ the agent that has learned the source is the **advisor** of the agent learning the target

## TL ALGORITHMS

- ▶ different transferables
- ▶ different techniques
- ▶ → we transfer a policy  $\pi$  using Policy Shaping
  - ▶ **Probabilistic Policy Reuse (PPR)**
  - ▶ **Policy Intersection (PI)**

## PROBABILISTIC POLICY REUSE

An advisory policy  $\pi_A(s_t)$  gets to decide which action to take with probability  $\psi$ :

$$a_t \sim \begin{cases} \pi_A(s_t) & \text{with probability } \psi \\ \pi_L(s_t) & \text{with probability } 1 - \psi \end{cases}$$

- ▶ looks like  $\epsilon$ -Greedy

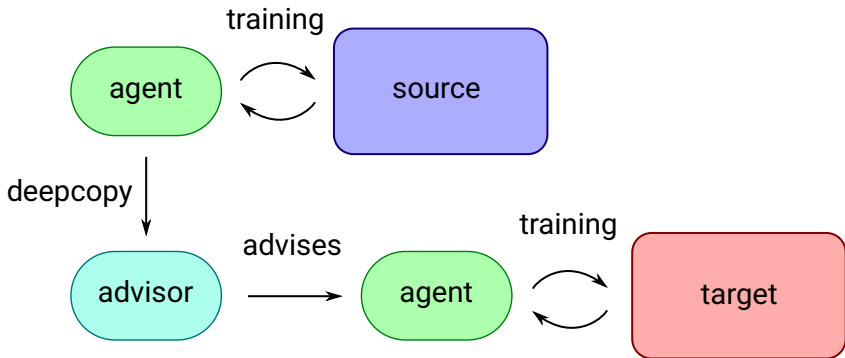
## POLICY INTERSECTION

An advisory policy  $\pi_A(s_t)$  is mixed with the agent's policy  $\pi_L(s_t)$  before an action  $a_t$  is sampled:

$$a_t \sim \frac{\pi_A(s_t)\pi_L(s_t)}{\underbrace{\pi_A(s_t) \cdot \pi_L(s_t)}_{\sum_{a \in A} \pi_A(a|s_t)\pi_L(a|s_t)}}$$

- ▶ acts as an intersection between the advisor and the agent

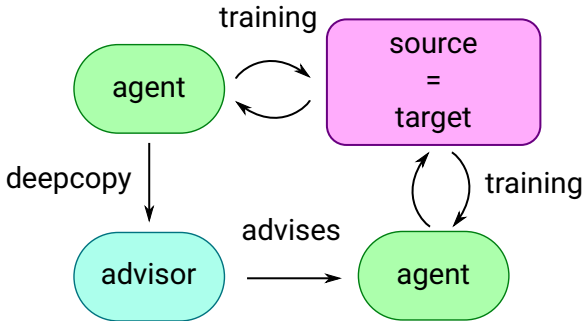
## Transfer





# CONTRIBUTION

## Self-Transfer



- ▶ different levels of training
- ▶ different initialization

**Dual Policy Distillation**, Kwei-Herng Lai, Daochen Zha, Yuening LiandXia Hu, IJCAI 2020

- ▶ launches two agents in two separate instances of the same environment
- ▶ with different initialization
- ▶ they regularly distill their policy in each other's policy

## PYBULLET CONTINUOUS CONTROL ENVS



- ▶ Need to adapt Policy Intersection for continuous actions

## POLICY INTERSECTION FOR CONTINUOUS ACTIONS

---

---

**Require:**  $\pi_L$  is the currently learning actor, and  $\pi_A$  is the frozen actor used as advisor

$A^A$  = set of actions sampled from  $\pi_A(s_t)$

**for** every action  $a_i^A \in A^A$  **do**

    Get probability  $\pi_L(a_i^A | s_t)$

    Compute possibility  $p_i = \frac{\pi_L(a_i^A | s_t)}{\max_{a \in A^A} \pi_L(a | s_t)}$

**if**  $p_i > p \sim U(0, 1)$  **then**

        Execute action  $a_i^A$

**return**

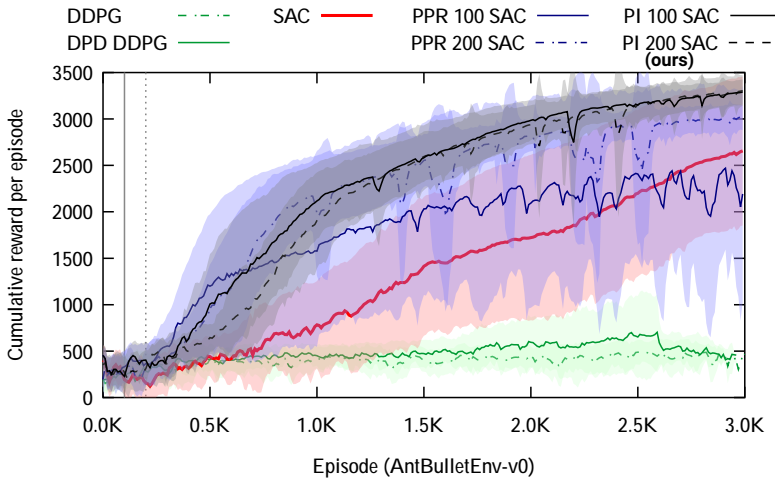
**end if**

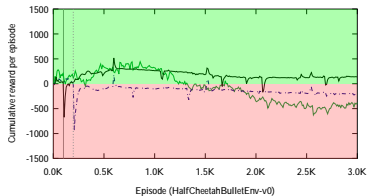
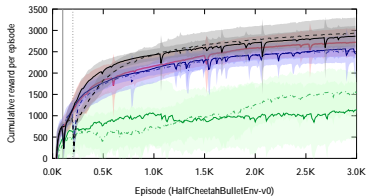
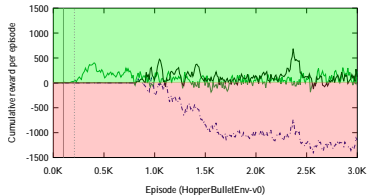
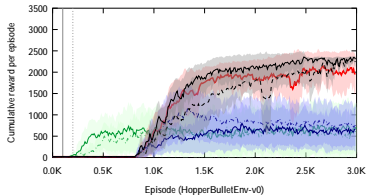
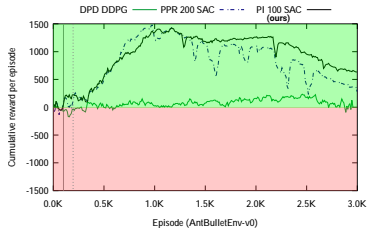
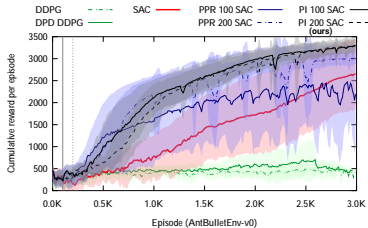
**end for**

Execute action  $a \sim \pi_L(s_t)$

---

# RESULTS





## SPECULATIONS

Why does Self-Transfer helps?

- ▶ the advisor pushes for exploration
- ▶ different initialization
- ▶ compensate for neural networks errors