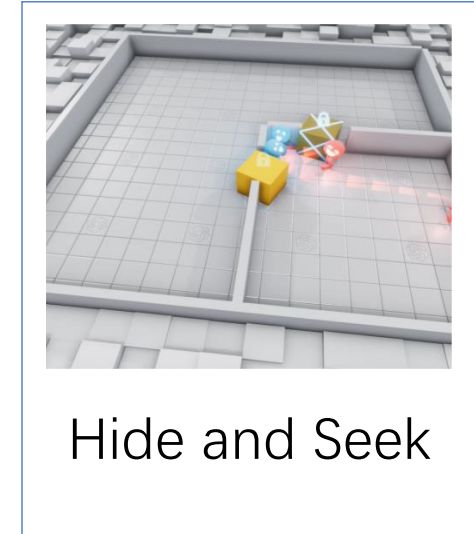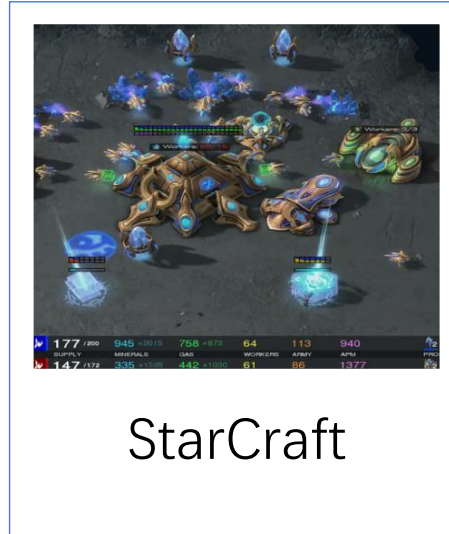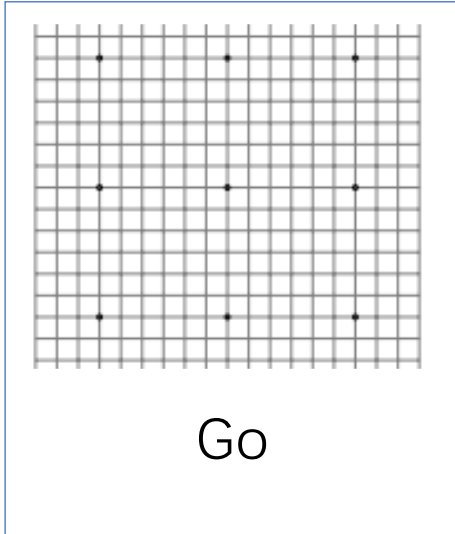# A Research on the Method of Automatic Curriculum Learning for Reinforcement Learning

Weizhe Chen

# Motivation

- Reinforcement Learning has succussed in multiple domains
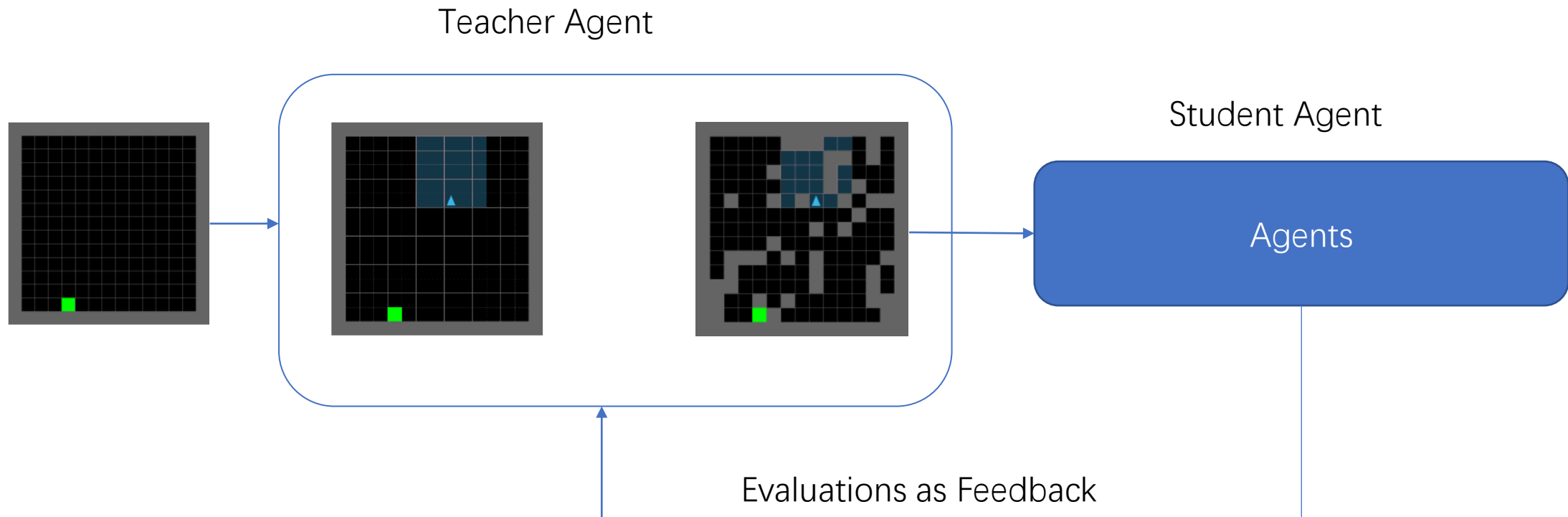


Go



StarCraft



Hide and Seek

- Curriculum learning is one of the most important auxiliary tool behind
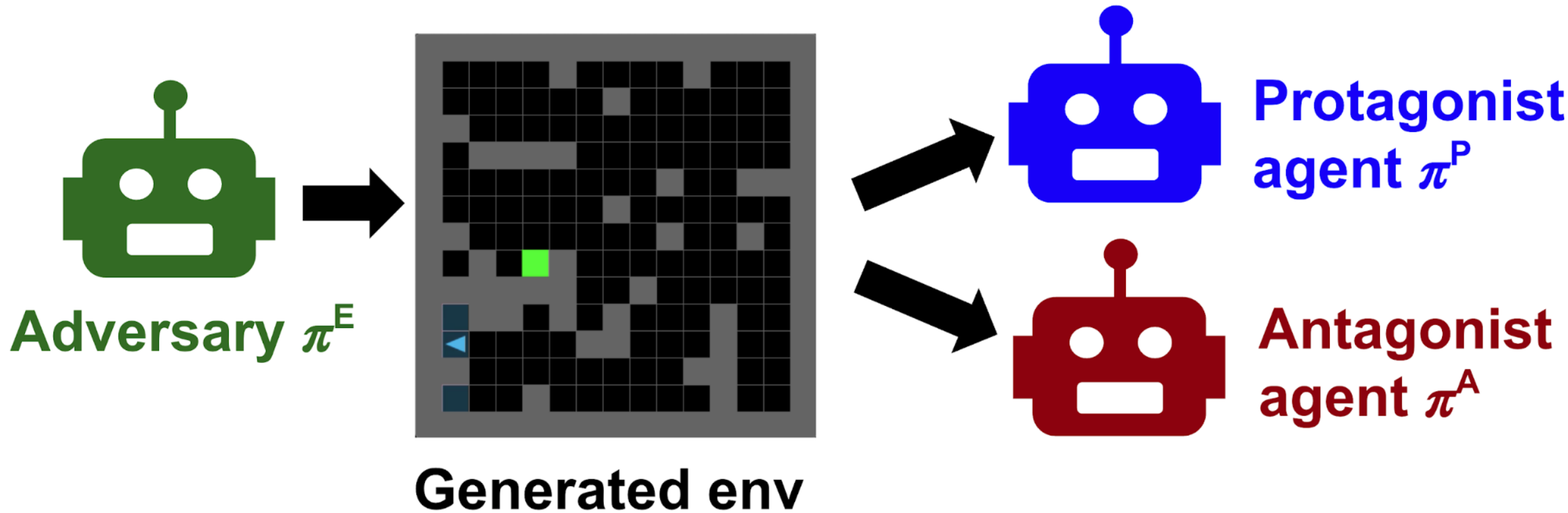  - Turn a hard enough problem into a series of sub-task and learn them sequentially

# The computing power is huge

- Dota2: 50k CPU, 500 GPU
- 2 vs 2 football: 4096+128 CPU, 16 TPU, 50 days

# Unsupervised Environment Design

Teacher Agent

Student Agent

Agents

Evaluations as Feedback

# PAIRED



**Adversary** $\pi^E$

**Generated env**

**Protagonist agent** $\pi^P$

**Antagonist agent** $\pi^A$

[1]. Dennis M, Jaques N, Vinitsky E, et al. Emergent Complexity and Zero-shot Transfer via Unsupervised Environment Design[J]. Neurips, 2020.
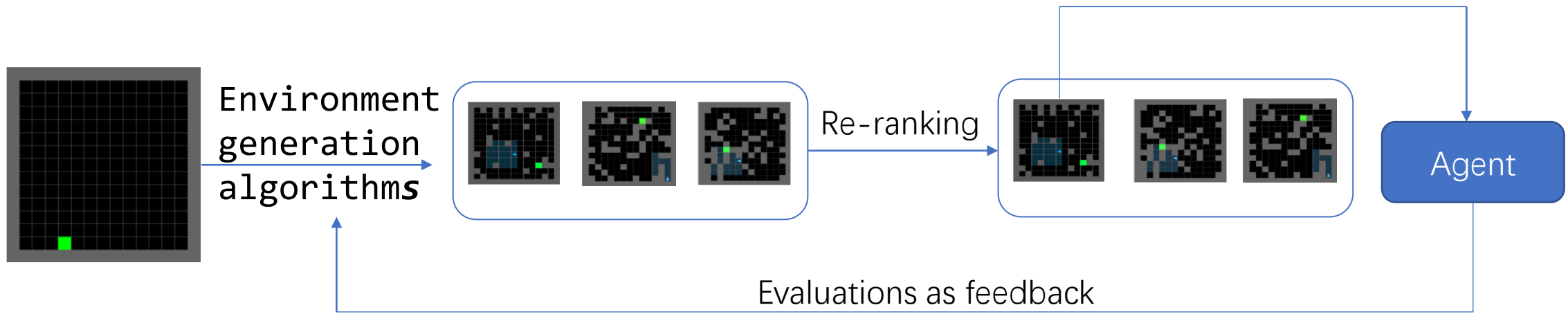
# PAIRED

- Using the additional agent to calculate the regret value.
    - $U = reward(Antagonist) - reward(Protagonist)$
    - $Reward'(Antagonist) = U$
    - $Reward'(Protagonist) = -U$
    - $Reward(Environment) = U$
    - Making the environment to be easier for the agent to learn since the difference of a good policy and the reward of a bad policy is now larger

# Problem Definition



Environment generation algorithm**s** → Re-ranking → Agent → Evaluations as feedback

# Problem Definition

- A special Re-ranking Problem:
    - Re-ranking: one step in the recommendation system to rank the items so that the recommendation can consider more than just the predicted score but also diversity
    - Environment: Item
    - Agent: User

# Markov Metric

- Defined as the difference between states
- Normal Requirement: $|V(s) - V(t)| \leq d(s, t)$
- well-setup metrics are hard to calculate!
- In this paper, we define $d(s, t) = |V(s) - V(t)|$
  - Can use the trained value function

# SOPED

- From intuition, it is easier for agent to learn if the current environment is not changing too much from the previous one
- We re-rank the potential environments by the last environment we used in training
- Similarity Ordered Population-based Environment Design (SOPED)

1   Randomly initialize $\pi^P, \pi^A, \pi^\Lambda$;

2   **while** *Not converged* **do**

3     Generate the environment parameters using the teacher agents population: $\theta_i \sim \Theta_i \sim \pi_i^\Lambda$;

4     Calculate the initial observation $o_i = I^{\theta_i}(s_i)$;

5     Calculate the distance with the latest trained environment $M_{\theta^{r-1}}$, and choosing the corresponding $\theta_i$ as the environment parameters used in this round of trainig $\theta^r$;

6     Create the training environment $M_\theta$ according to $\theta_r$;

7     Sample training trajectories $\tau^A$ and $\tau^P$ in $M_\theta$, and calculate $U(\pi^A)$ 和 $U(\pi^P)$;

8     Calculate ***Regret*** $= U(\pi^A) - U(\pi^P)$;

9     Train policy $\pi^P$ with reward $-$***Regret***;

10    Train policies $\pi^A$ and $\pi^\Lambda$ with reward ***Regret***;

11   **end**

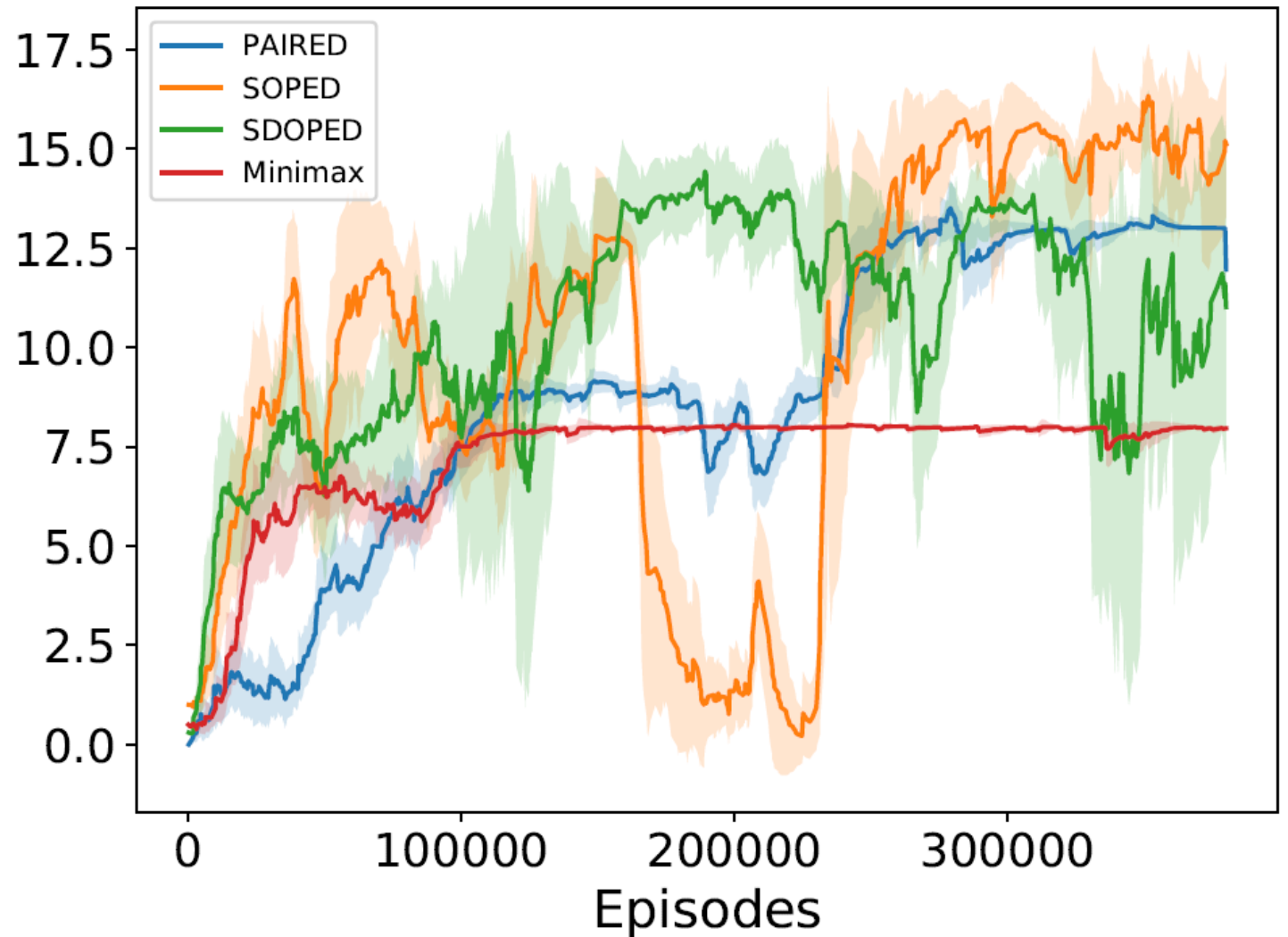12   **return** $\pi^P, \pi^A, \pi^\Lambda$

# SDOPED

- On the other hand, SOPED will more likely make the high0level generator to stuck in a local optimal rather than a global optimal

- We need to consider diversity during our training

- Determinantal Point Process (DPP) is used
  - Given the similarity of each pair of items and the score of the items, DPP can re-rank the sequence so that it both consider the score and the diversity
  - Similarity $s = e^{-0.5d}$, where d is the distance

- Similarity and Diversity Ordered Population-based Environment Design, SDOPED
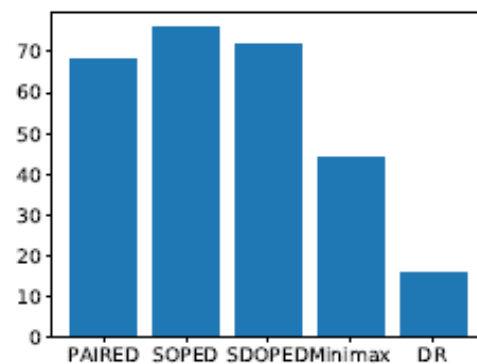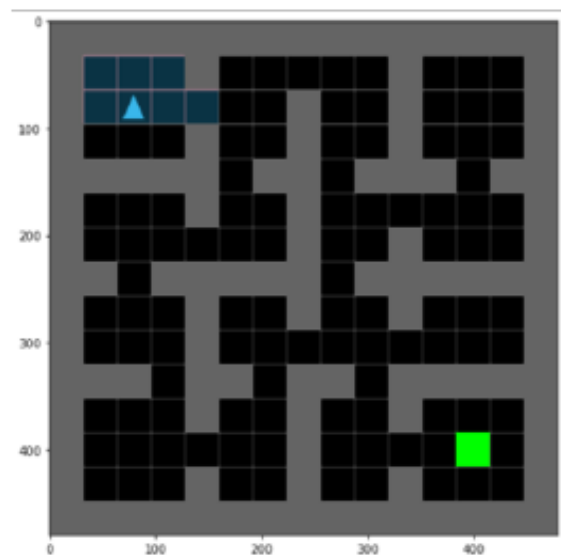
# SDOPED

1   Randomly initialize $\pi^P, \pi^A, \pi^\Lambda$;

2  **while** *Not converged* **do**

3       Generate the environment parameters using the teacher agents population:$\theta_i \sim \Theta_i \sim \pi_i^\Lambda$;

4       Calculate the initial observation $o_i = I^{\theta_i}(s_i)$;

5       Calculate the distance with the latest trained environment $M_{\theta r-1}$ as the score of the item;

6       Calculating the similarity of each pair of environments as S;

7       Calculate the Kernel Matrix $L = Diag(r)\dot{S}Diag(r)\,\Theta_{small} = DPP(L)$;

8       Uniformly sample $\theta \in \Theta_{small}$;

9       Create actual environment $M_\theta$ according to $\theta_r$;

10      Sample training trajectories $\tau^A$ and $\tau^P$ in $M_\theta$, and calculate $U(\pi^A)$ 和 $U(\pi^P)$;

11      Calculate $Regret = U(\pi^A) - U(\pi^P)$;

12      Train policy $\pi^P$ with reward $-Regret$;

13      Train policies $\pi^A$ and $\pi^\Lambda$ with reward $Regret$;

14  **end**

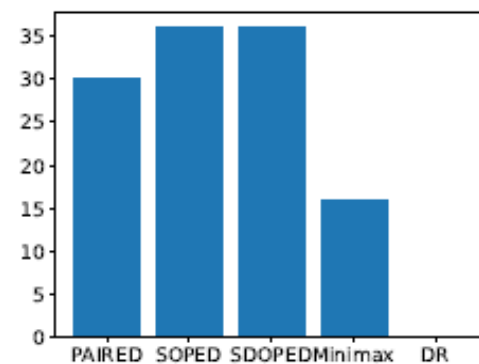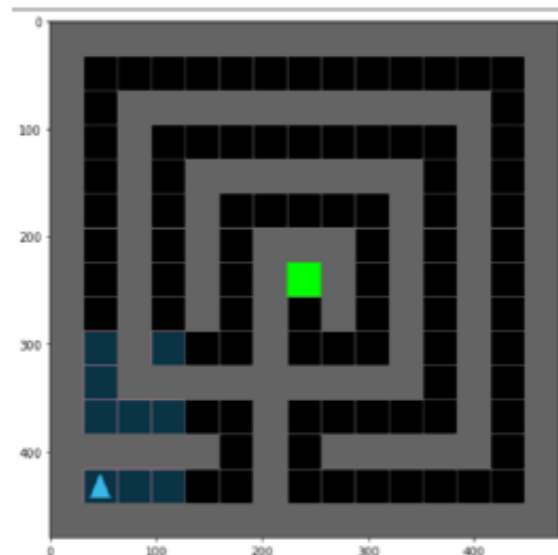15  **return** $\pi^P, \pi^A, \pi^\Lambda$

# Experiment

- Maze Design
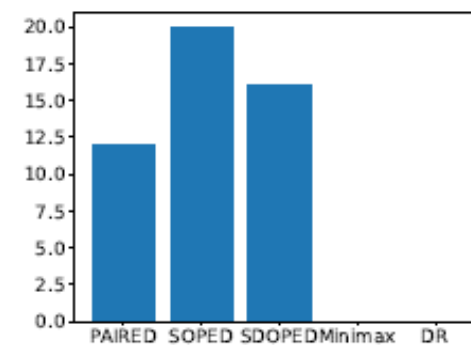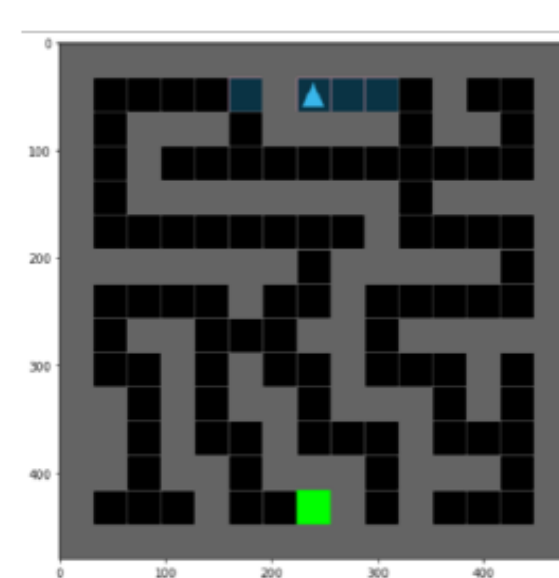  - Solved Path Length:
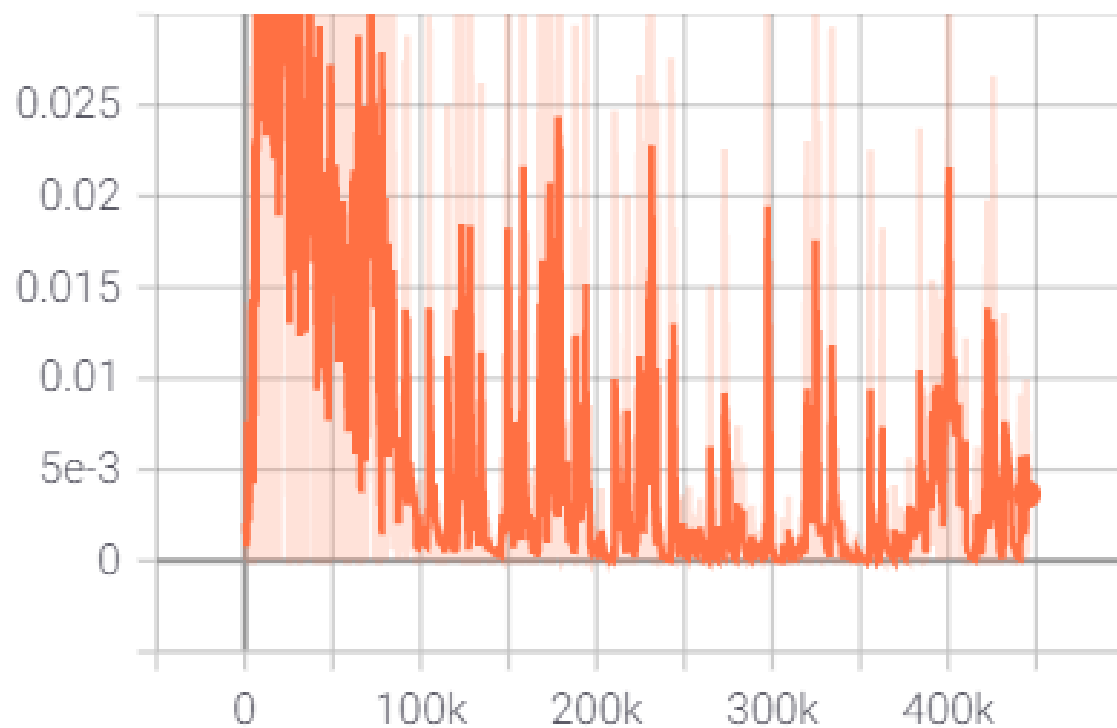
# Experiment



(a). 16 Room

(b). Labyrinth

(c). Maze

# Experiment

- Value loss:

# Takeaway

- Using environment design, we can improve the training efficiency of the agents, and make the final result more robust

- The value predictor of the high-level environment generator can be easily trained, but the environment generator is very hard to train

- Even if the generator seems to have converge, it still could suddenly completely change its output