

Learn Goal-Conditioned Policy with Intrinsic Motivation for Deep Reinforcement Learning

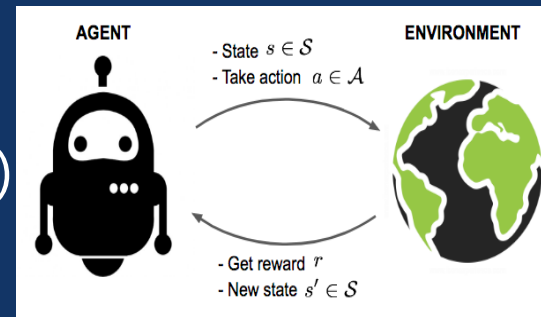
Presenter: Jinxin Liu (liujinxin@westlake.edu.cn)

Co-Authors: Donglin Wang; Qiangxing Tian; Zhengyu Chen

CHAPTER 1

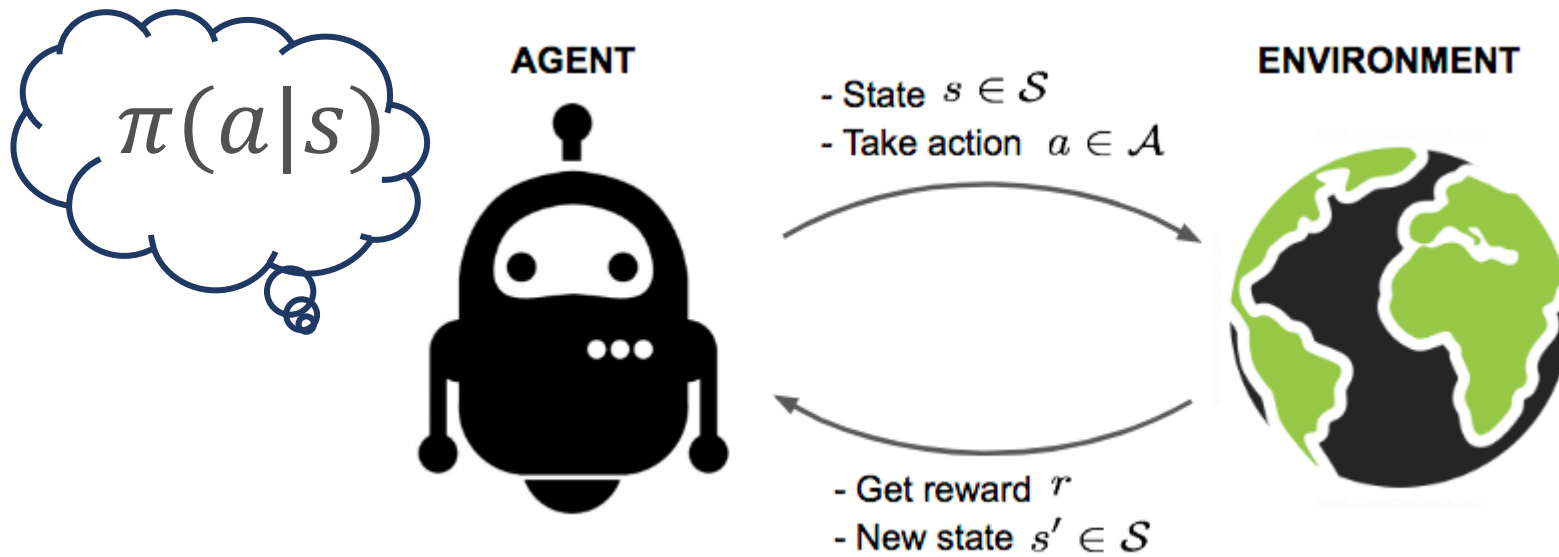
Deep Reinforcement Learning (RL)

A brief explanation of deep reinforcement learning.



Reinforcement Learning (RL)

Mathematical formalism:
Reinforcement Learning



$$\max_{\pi} E_{\pi}(\sum_t [r_t | \pi])$$

where reward $r_t := r(s_t, a_t)$

Representation:
Deep Networks

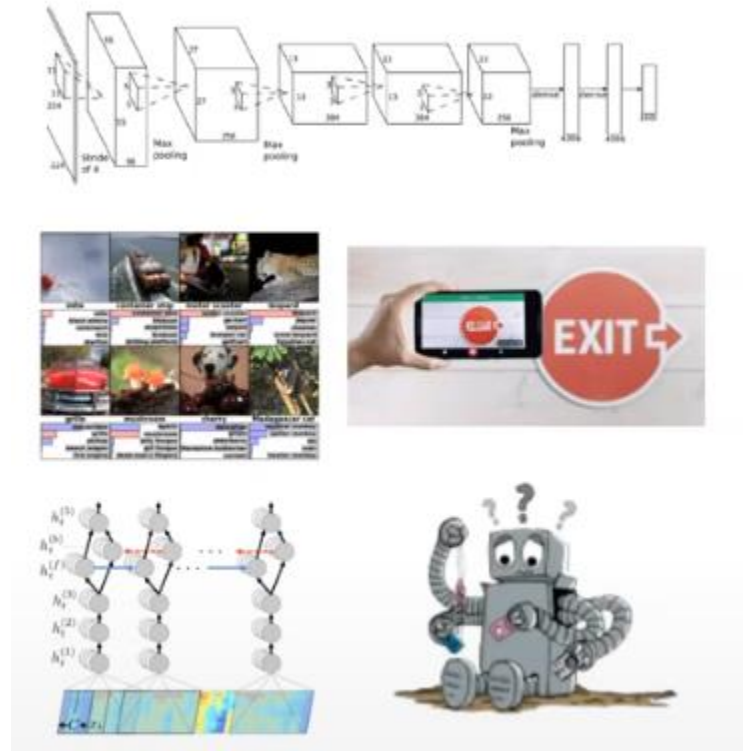


Image credit: RAIL

Reinforcement Learning (RL)

It (more or less) works



Stack a Lego block



Image Credit:
Google

Scalable Deep Reinforcement Learning for Robotic Manipulation



Image credit: *Push a mug onto a coaster*
<https://bair.berkeley.edu/blog/2019/05/20/solar/>

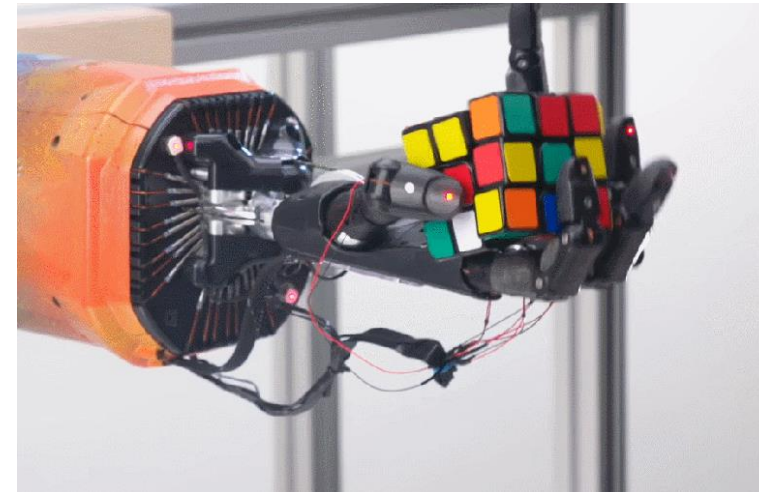


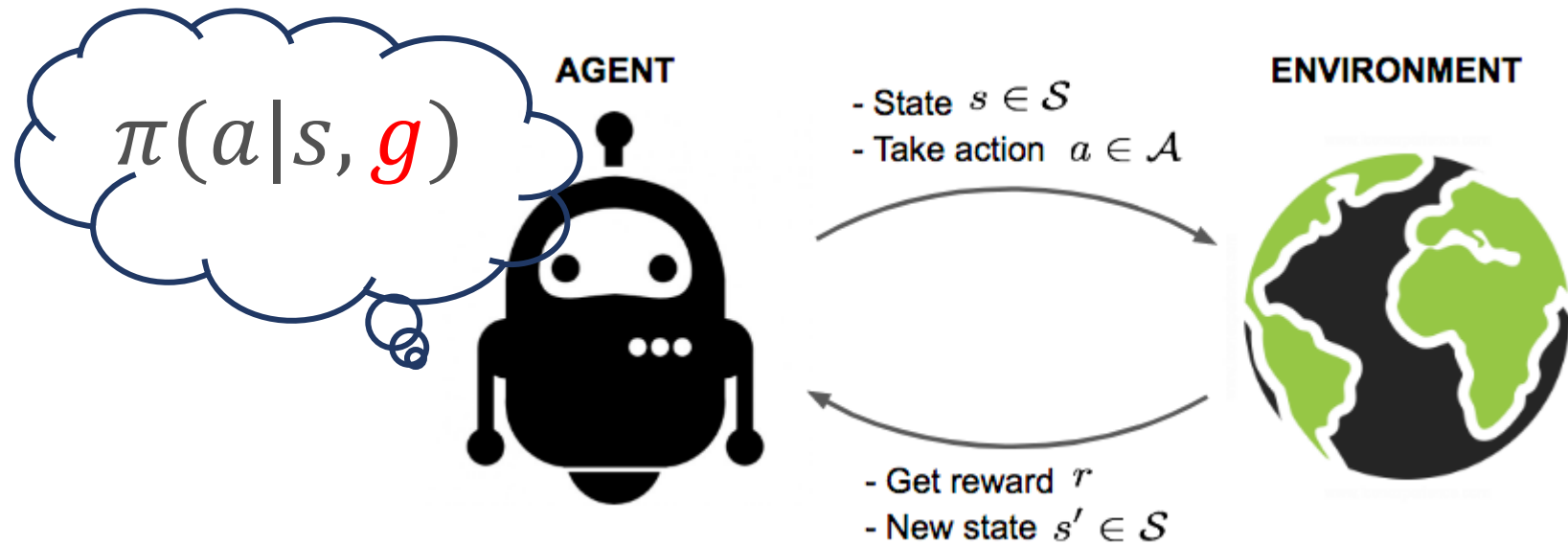
Image Credit:
OpenAI

Solving Rubik's Cubewith a Robot Hand

Multi-Goal Reinforcement Learning

Mathematical formalism:

Multi-Goal Reinforcement Learning



$$\max_{\pi} \mathbf{E}_g \mathbf{E}_{\pi} (\sum_t [r_t | \pi])$$

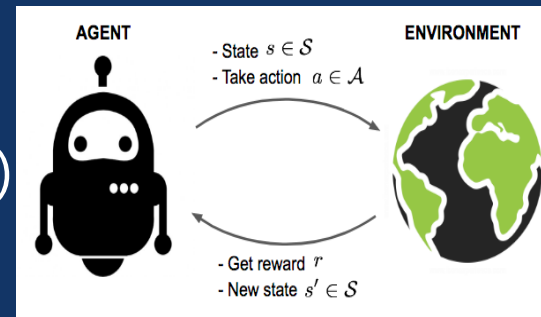
where reward $r_t := r(s_t, a_t, g)$

CONTENTS

CHAPTER 1

Deep Reinforcement Learning (RL)

A brief explanation of deep reinforcement learning.



CHAPTER 2

Unsupervised RL

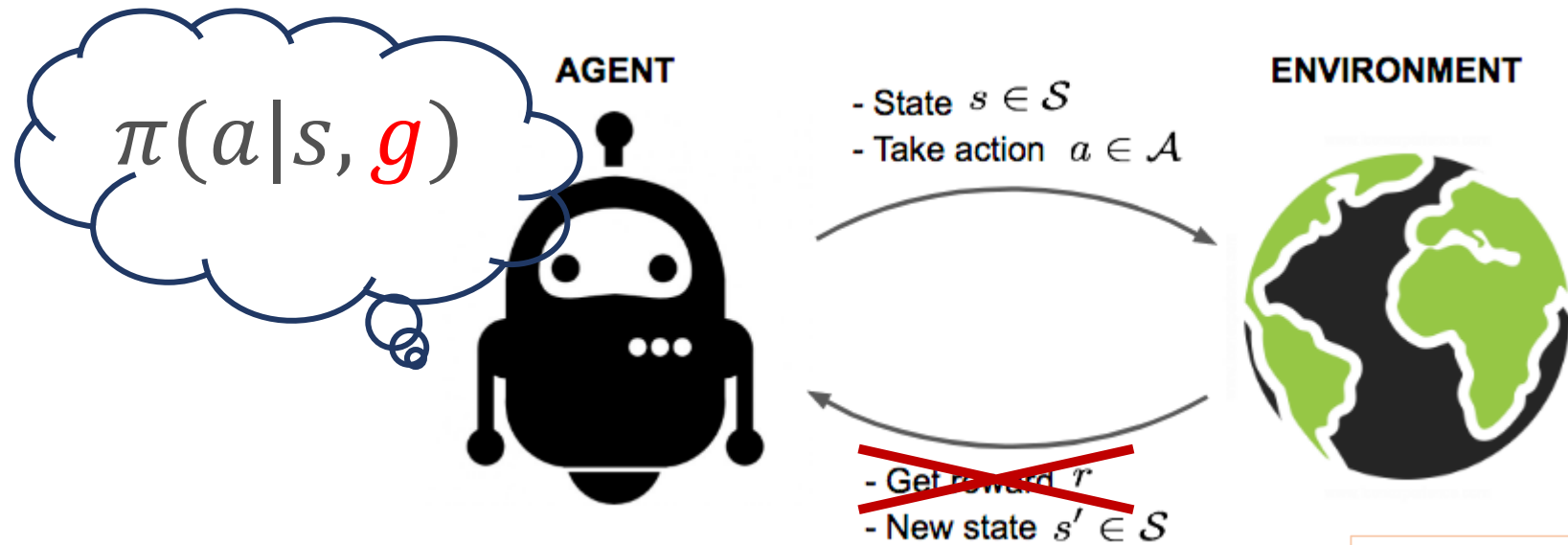
A brief explanation of unsupervised reinforcement learning.



Unsupervised/Self-supervised RL

Mathematical formalism:

Unsupervised Reinforcement Learning



$$\max_{\pi} \mathbf{E}_g \mathbf{E}_{\pi} (\sum_t [r_t | \pi])$$

where reward $r_t := r(s_t, a_t, s_{t+1}, g)$ **X**

1. Learn reward function

$$r_t := r(s_t, a_t, s_{t+1}, g)$$

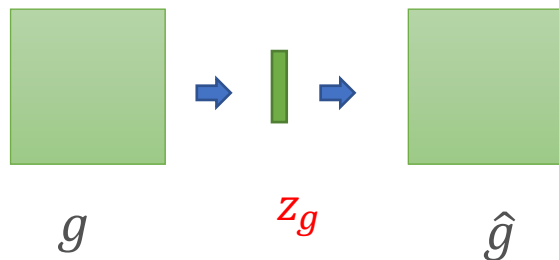
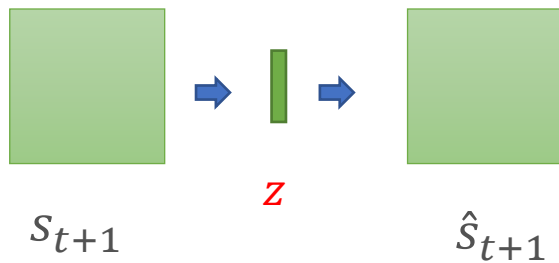
2. Learn goal distribution

$$\mathbf{E}_g$$

Unsupervised/Self-supervised RL

1. Learn reward function

$$r_t := r(s_t, a_t, s_{t+1}, g) \\ = r(s_{t+1}, g)$$



$$r(s_{t+1}, g) := -\cos(z, z_g)$$

2. Learn goal distribution

Assume the spaces of perceptual goals and states are same.

Sample goals from the historical trajectories of the policy to be trained.

perceptually-specific goal based approach

Unsupervised/Self-supervised RL

Latent-variable based approach

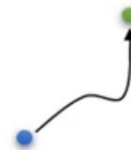
Sharma, Archit, et al.

"Dynamics-aware unsupervised discovery of skills."

Eysenbach, Benjamin, et al.

"Diversity is all you need: Learning skills without a reward function."

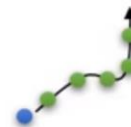
$$\max_{\pi, p(\mathcal{T})} \mathcal{H}[\mathbf{s}_H] - \mathcal{H}[\mathbf{s}_H | \mathcal{T}] = \mathcal{I}[\mathbf{s}_H; \mathcal{T}]$$



$$\max_{\pi, p(\mathcal{T})} \mathcal{H}[\tau] - \mathcal{H}[\tau | \mathcal{T}] = \mathcal{I}[\tau; \mathcal{T}]$$



$$\max_{\pi, p(\mathcal{T})} \mathcal{H}[\mathbf{s}] - \mathcal{H}[\mathbf{s} | \mathcal{T}] = \mathcal{I}[\mathbf{s}; \mathcal{T}]$$



Unsupervised/Self-supervised RL

Eysenbach, Benjamin, et al.

"Diversity is all you need: Learning skills without a reward function."

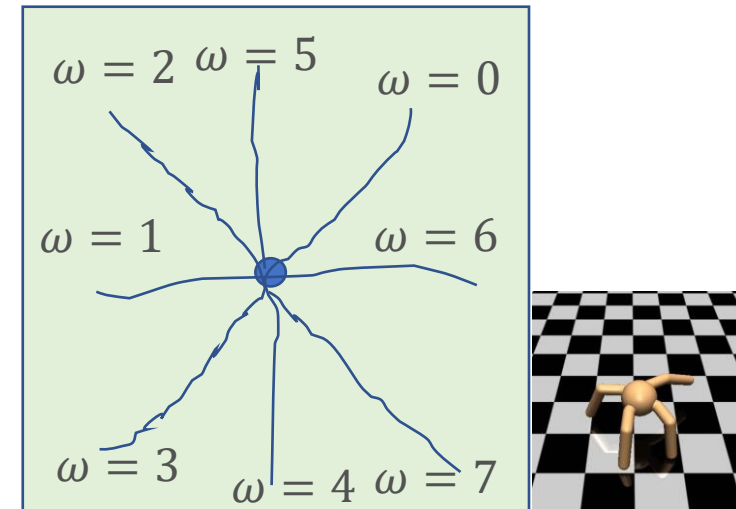
$$\text{DIAYN: } \max I(\omega; s)$$

$$\pi(a|s, \omega), \omega \sim p(\omega)$$

$$\max_{\pi, p(s)} I(\omega; s) = H[\omega] - H[\omega|s]$$

$$I(\omega; \tau) \geq \underbrace{H[\omega]}_{\text{prior (fixed)}} + \underbrace{\mathbb{E}_{\omega, s} [q_{\phi}(\omega|s)]}_{\text{reward function: } r}$$

prior (fixed) *reward function: r*

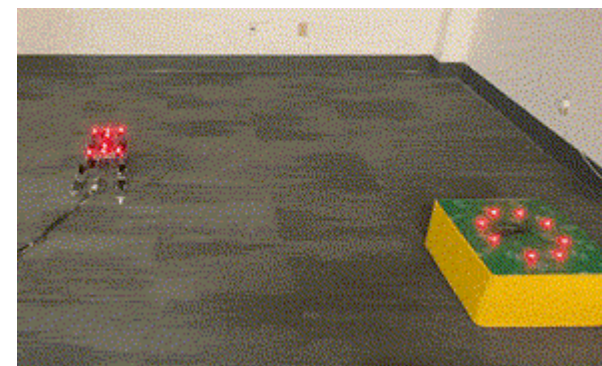
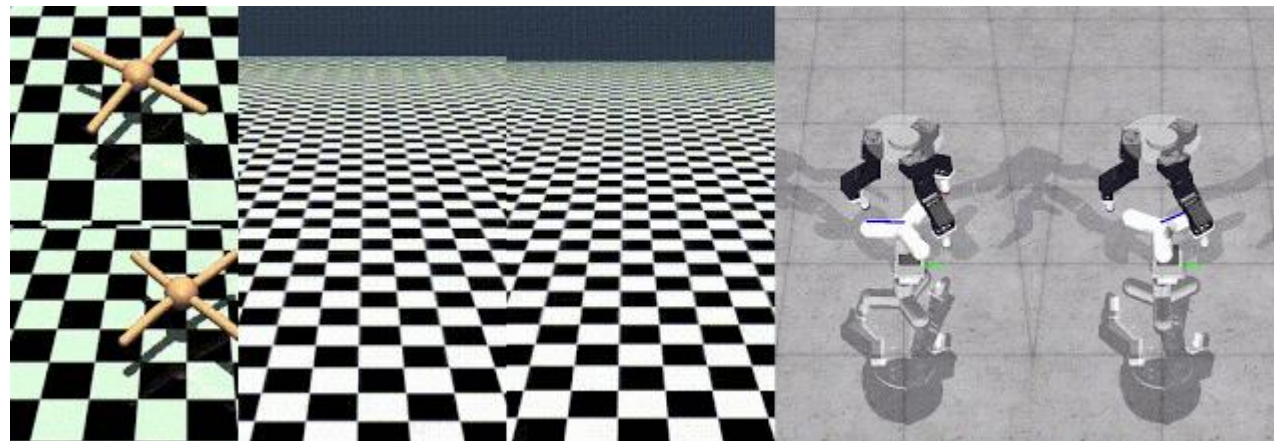
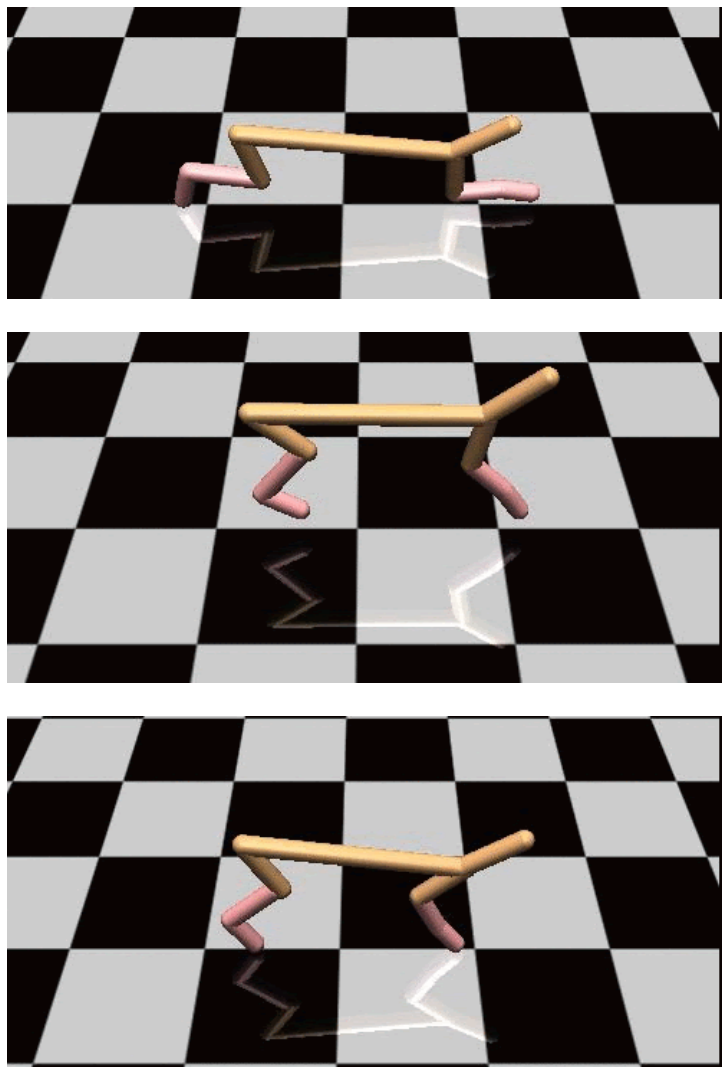


● *Initial State*

Agent

Latent-variable based approach

Unsupervised/Self-supervised RL



Learned diverse skills.

Image Credit: GOOGLE

Unsupervised/Self-supervised RL

perceptually-specific goal based approach

$$r(s_{t+1}, g) := -\cos(z, z_g)$$

Prior non-parametric measure function may
limit the repertoires of behaviors and
impose manual engineering burdens.

Latent-variable based approach

$$\max I(\omega; s) \\ \pi(a|s, \omega), \omega \sim p(\omega)$$

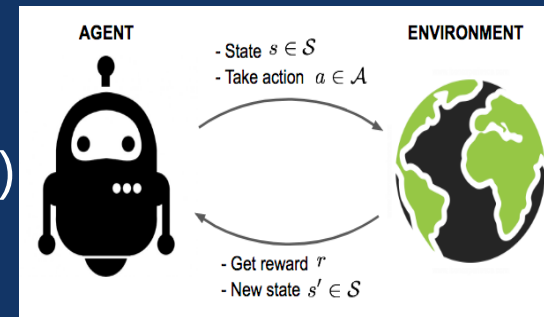
Such policy is conditioned on latent variables
rather than perceptually-specific goals.



CHAPTER 1

Deep Reinforcement Learning (RL)

A brief explanation of deep reinforcement learning.



CHAPTER 2

Unsupervised RL

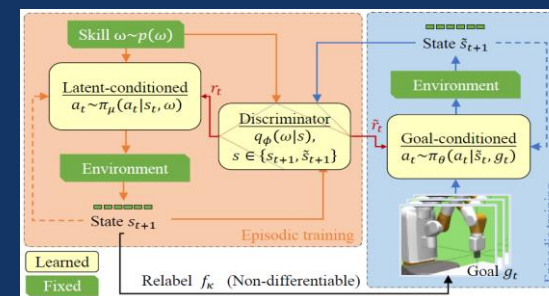
A brief explanation of unsupervised reinforcement learning.



CHAPTER 3

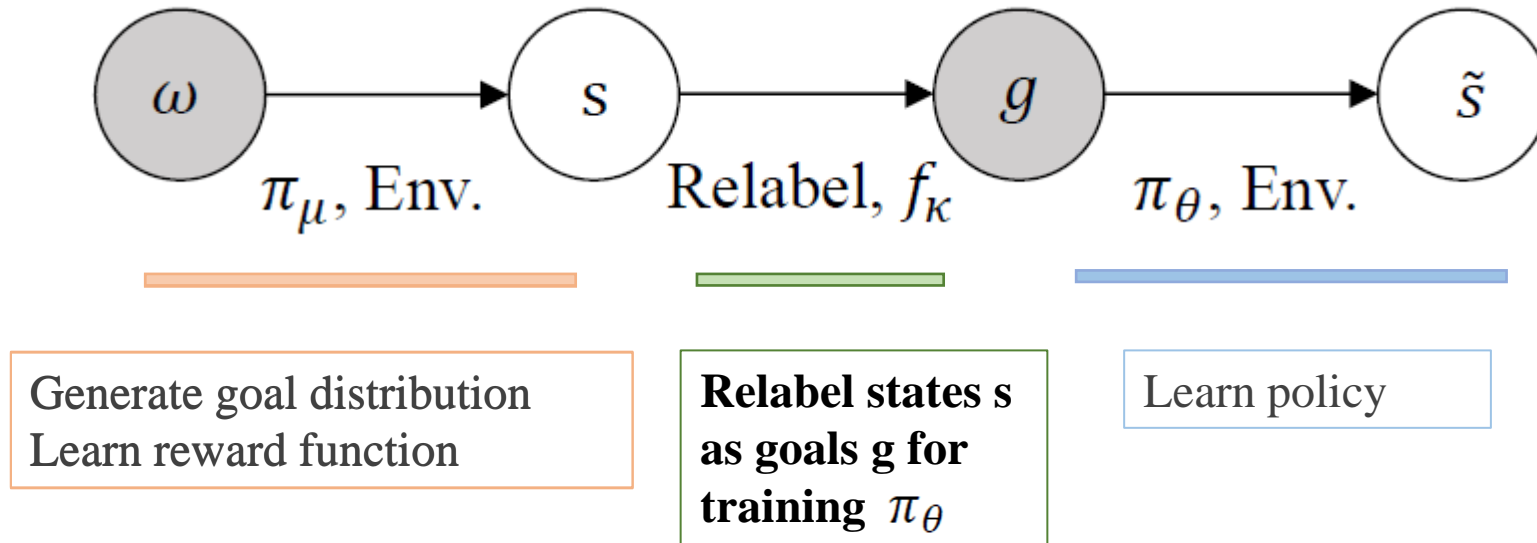
Method: GPIM

Learn Goal-Conditioned Policy with Intrinsic Motivation for Deep Reinforcement Learning



Proposed GPIM method

Learn Goal-conditioned Policy with Intrinsic Motivation



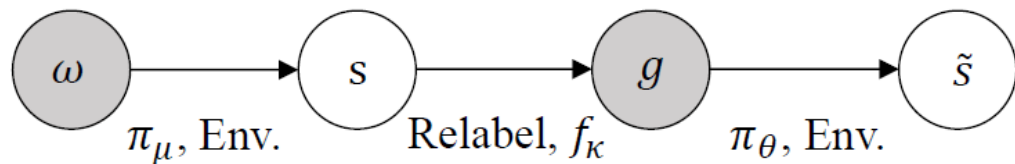
$$\mathcal{F}(\mu, \theta) = \mathcal{I}(s; \omega) + \mathcal{I}(\tilde{s}; g)$$

$$\mathcal{F}(\mu, \theta) \geq \mathcal{I}(s; \omega) + \mathcal{I}(\tilde{s}; \omega)$$

$$= 2\mathcal{H}(\omega) + \mathbb{E}_{p_m(\cdot)} [\log p(\omega|s) + \log p(\omega|\tilde{s})]$$

Proposed GPIM method

Learn Goal-conditioned Policy with Intrinsic Motivation



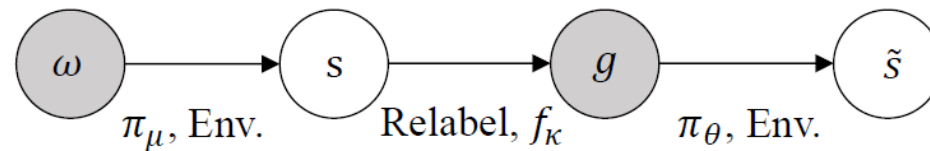
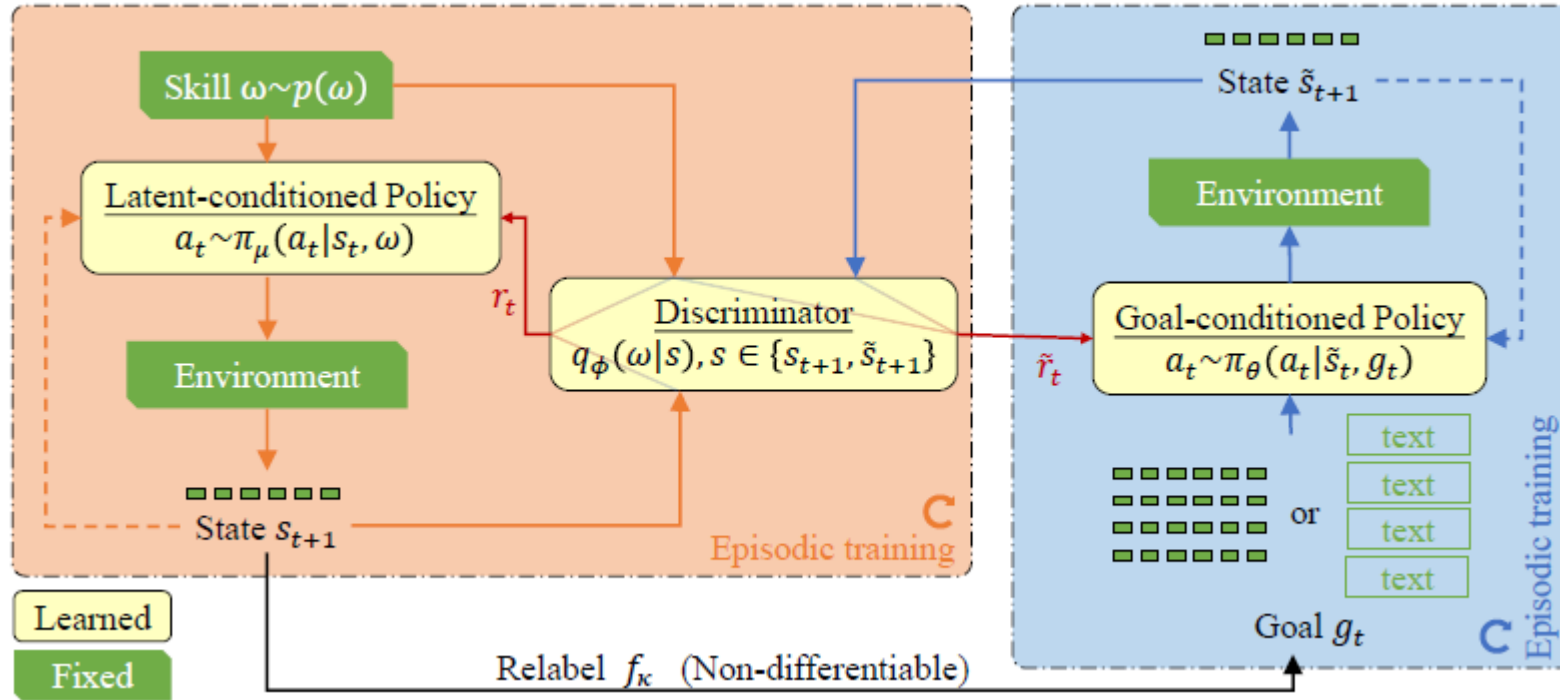
$$\begin{aligned}
 \mathcal{F}(\mu, \theta) &\geq \mathcal{I}(s; \omega) + \mathcal{I}(\tilde{s}; \omega) \\
 &= 2\mathcal{H}(\omega) + \mathbb{E}_{p_m(\cdot)} [\log p(\omega|s) + \log p(\omega|\tilde{s})] \\
 &\geq 2\mathcal{H}(\omega) + \mathbb{E}_{p_m(\cdot)} [\log q_\phi(\omega|s) + \log q_\phi(\omega|\tilde{s})]
 \end{aligned}$$

Reward function (for both policies):

$$r_t = \log q_\phi(\omega|s_{t+1}) - \log p(\omega)$$

Proposed GPIM method

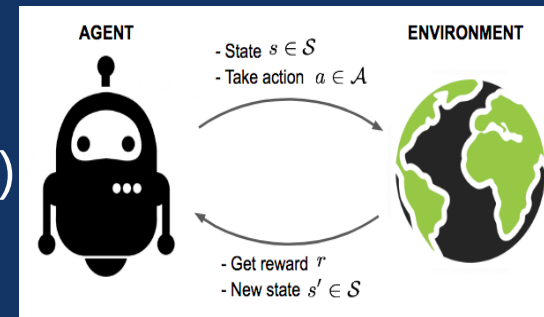
Learn Goal-conditioned Policy with Intrinsic Motivation



CHAPTER 1

Deep Reinforcement Learning (RL)

A brief explanation of deep reinforcement learning.



CHAPTER 2

Unsupervised RL

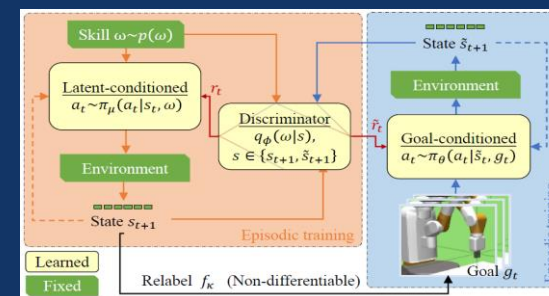
A brief explanation of unsupervised reinforcement learning.



CHAPTER 3

Method: GPIM

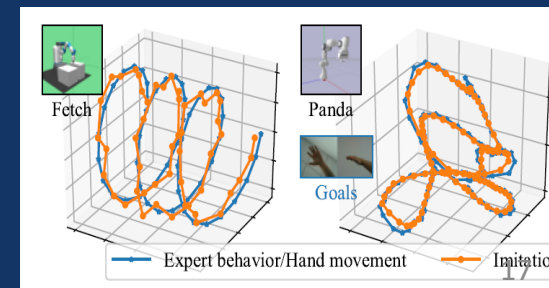
Learn Goal-Conditioned Policy with Intrinsic Motivation for Deep Reinforcement Learning



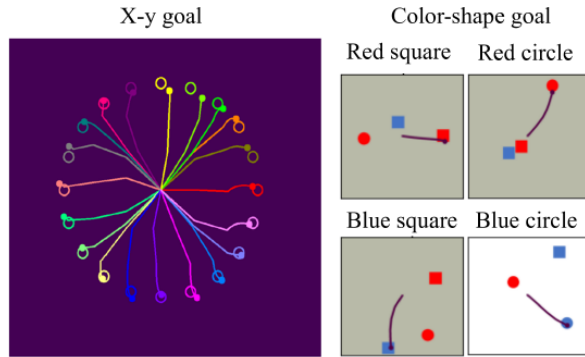
CHAPTER 4

Experiments

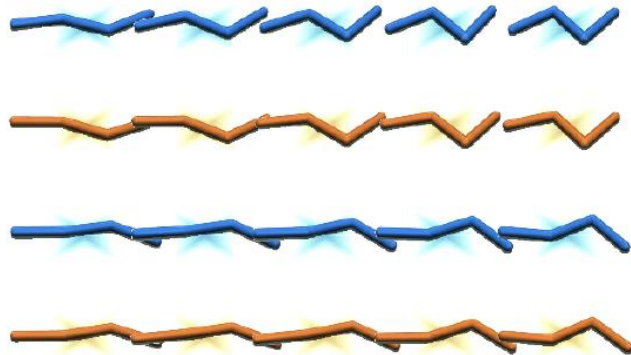
Briefly introduce the experiments in GPIM.



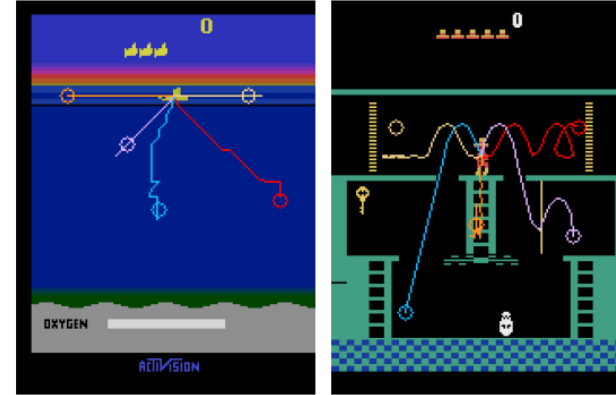
Experiments: Visualization



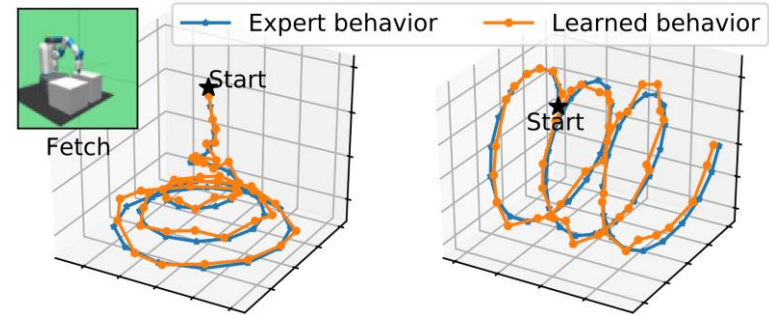
Fixed goals



Dynamic goals

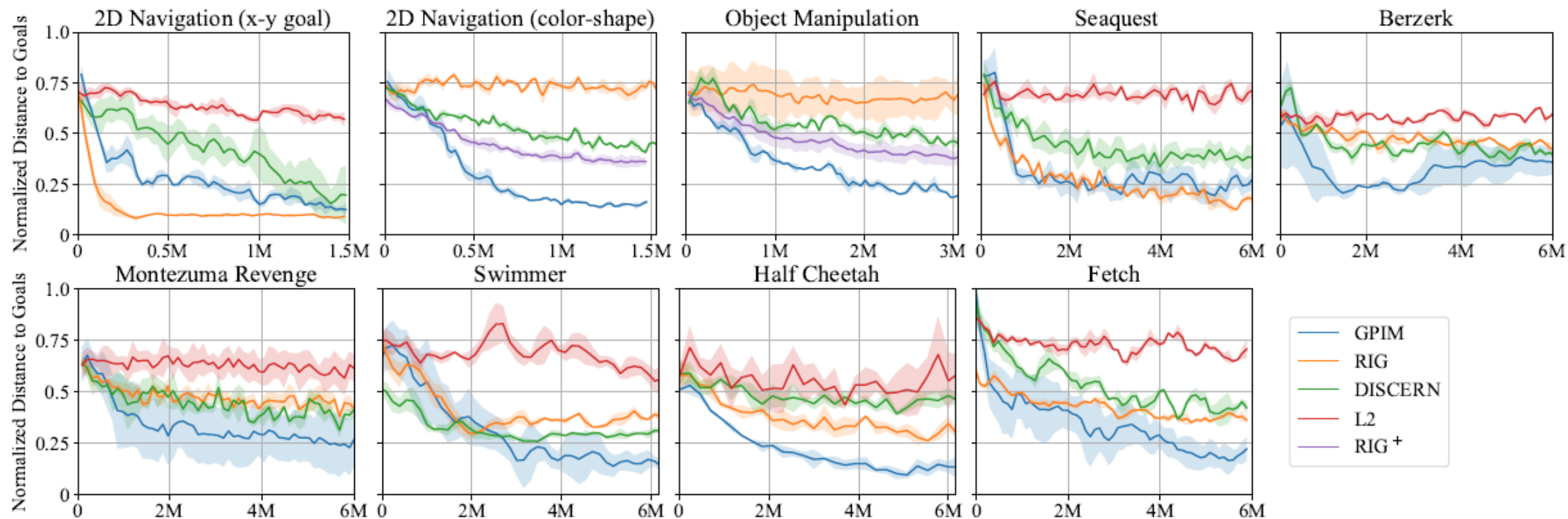


High dimensional goals



Temporally-extended goals

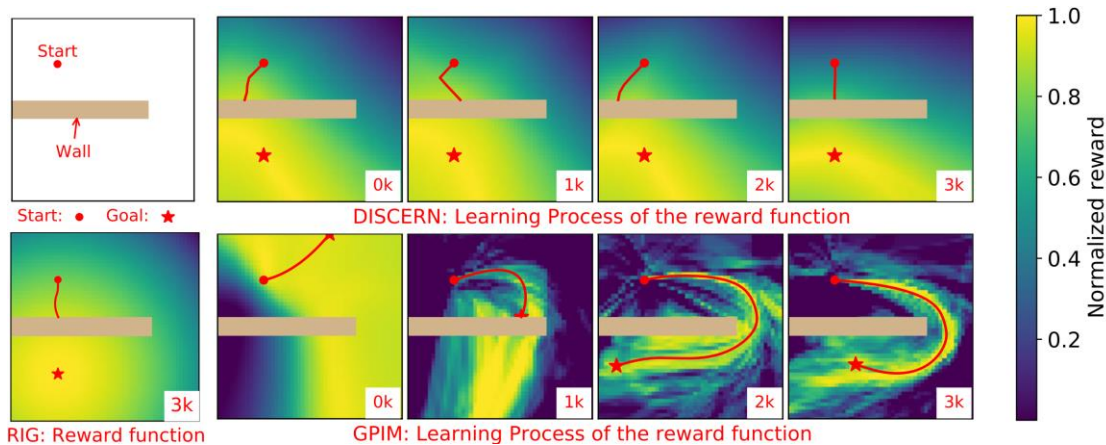
Experiments: Compared with Baselines



Performance (normalized distance to goals vs. actor steps) of our GPIM and baselines.

Experiments: Dynamical Distance Learning

The performance of unsupervised RL methods depends on the diversity of autonomously generated goals and **the expressiveness of the learned reward function**, which is conditioned on the generated goals.



Our method (GPIM) builds up the reward function after exploring the environment, *the dynamic of which itself further shapes the reward function*. We can see that our model provides the reward function better expressiveness of the task by compensating for the dynamic.

Conclusion

1. We introduce a latent-conditioned policy with a procedural relabeling function to generate tasks for training the goal-conditioned policy.
2. We theoretically describe the performance guarantee of our (unsupervised) objective compared with the standard multi-goal RL.
3. We also conduct extensive experiments on a variety of robotic tasks to demonstrate the effectiveness and efficiency of our method, which outperforms prior unsupervised methods.

Thank you