

# Learning to Explore via Meta-Policy Gradient

Tianbing Xu\*, Qiang Liu (UT, Austin), Liang Zhao\*, Jian Peng  
(UIUC)

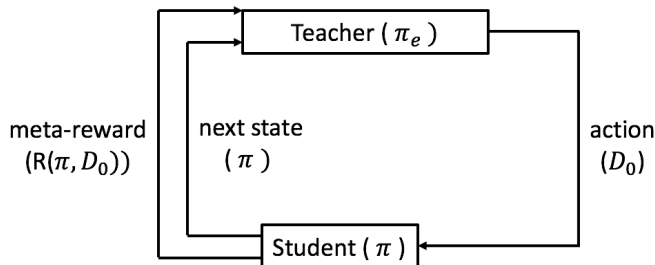
Baidu Research\*, CA, USA

July 7, 2018

# Demo

Video Result: [▶ Continuous Control Tasks](#)

# Teacher-Student Exploration-Exploitation Interactions



- ▶ Teacher (exploration policy  $\pi_e$ )  
Learning to generate high quality data based on student's performance improvement
- ▶ Student (exploitation policy  $\pi$ )  
Learning from teacher's demonstrations to improve the performance

# Meta-Reward

The student performance improvement:

$$\mathcal{R}(\pi, D_0) = R(\pi') - R(\pi) \quad (1)$$

$\pi'$  look-ahead policy of student

$$\pi' = DDPG(\pi, D_0)$$

$R(\pi)$  the cumulative reward of roll-out generated by policy  $\pi$ .

# Learning to Explore via Meta-Policy Gradient

Teacher's Objective:

$$J(\pi_e) = E_{D_0 \sim \pi_e} [\mathcal{R}(\pi, D_0)]$$

Teacher's policy gradient (Meta-Policy Gradient):

$$\nabla_{\theta^{\pi_e}} J = E_{D_0 \sim \pi_e} [\mathcal{R}(\pi, D_0) \nabla_{\theta^{\pi_e}} \log P(D_0 | \pi_e)] \quad (2)$$

where,

$$\nabla_{\theta^{\pi_e}} \log P(D_0 | \pi_e) = \sum_t \nabla_{\theta^{\pi_e}} \log \pi_e(a_t | s_t)$$

# Learning to Explore Algorithm

---

**Algorithm 1** Learning to Explore

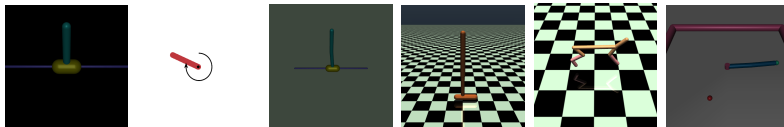
---

- 1: **for** iteration  $t$  **do**
- 2:   Generate  $D_0$  by executing teacher's policy  $\pi_e$ .
- 3:   Update actor policy  $\pi$  to  $\pi'$  using DDPG based on  $D_0$
- 4:   Generate  $D_1$  from  $\pi'$  and estimate the reward of  $\pi'$ . Calculate the meta reward:  $\hat{\mathcal{R}}(\pi, D_0) = \hat{R}_{\pi'} - \hat{R}_{\pi}$ .
- 5:   Update Teacher's Policy  $\pi_e$  with meta policy gradient

$$\theta^{\pi_e} \leftarrow \theta^{\pi_e} + \eta \nabla_{\theta^{\pi_e}} \log \mathcal{P}(D_0 | \pi_e) \hat{\mathcal{R}}(\pi, D_0)$$

- 6:   Add both  $D_0$  and  $D_1$  into the Replay Buffer  $B \leftarrow B \cup D_0 \cup D_1$ .
  - 7:   Update  $\pi$  using DDPG based on Replay Buffer, that is,  $\pi \leftarrow \text{DDPG}(\pi, B)$ . Compute the new  $\hat{R}_{\pi}$ .
  - 8: **end for**
-

# Experiments on Continuous Control Tasks



**Figure:** Illustrative screen-shots of environments we experiment with Meta and DDPG

# Meta-Exploration Policy Explores Efficiently

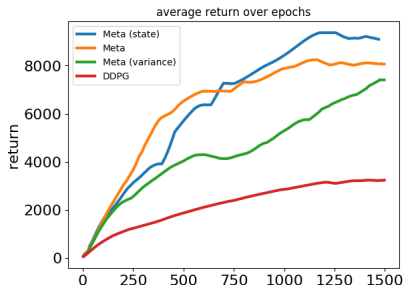


Figure: Comparison between meta exploration policies and DDPG

Meta:

$$\pi_e \sim N(f(s, \theta^{\pi_e}), \Sigma)$$

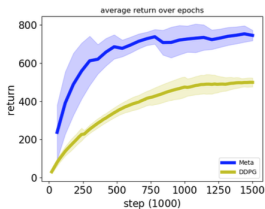
Meta (State): adding more Q-function related features

Meta (Variance):

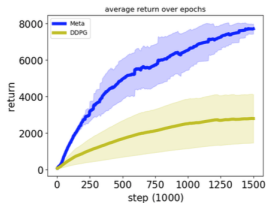
$$\pi_e \sim N(\mu(s, \theta^\pi), \Sigma)$$



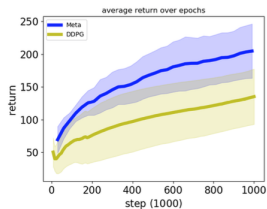
# Sample Efficiency with Higher Return



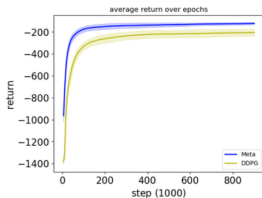
(a) InvertedPendulum



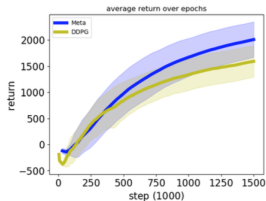
(b) InvertedDoublePendulum



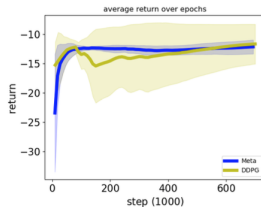
(c) Hopper



(d) Pendulum



(e) HalfCheetah



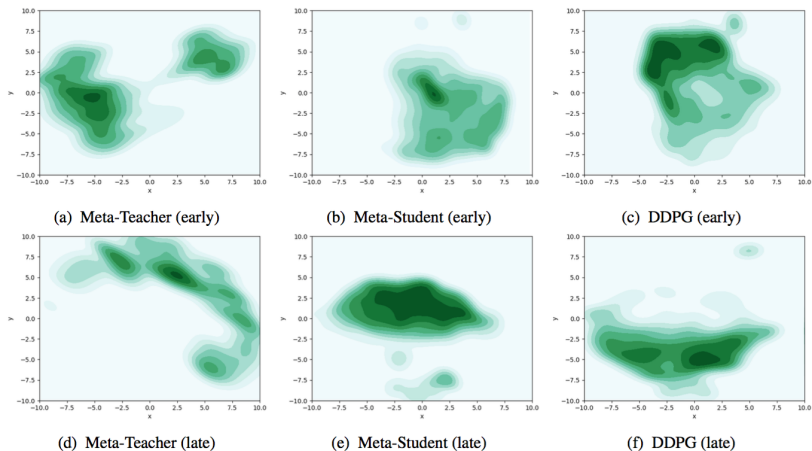
(f) Reacher

# Sample Efficiency with Higher Return

Table: Reward achieved in different environments

env-id	Meta	DDPG
InvertedDoublePendulum-v1	<b>7718 <math>\pm</math> 277</b>	2795 $\pm$ 1325
InvertedPendulum-v1	<b>745 <math>\pm</math> 27</b>	499 $\pm$ 23
Hopper-v1	<b>205 <math>\pm</math> 41</b>	135 $\pm$ 42
Pendulum-v0	<b>-123 <math>\pm</math> 10</b>	-206 $\pm$ 31
HalfCheetah-v1	<b>2011 <math>\pm</math> 339</b>	1594 $\pm$ 298
Reacher-v1	-12.16 $\pm$ 1.19	<b>-11.67 <math>\pm</math> 3.39</b>

# Guided Exploration with Diverse and Adaptive Meta Policies



**Figure:** State Visitation Density Contours of Meta and DDPG in Early (the first row) and Late (the second row) Learning Stages.

# Conclusion

- ▶ Developed a Meta-RL method for More Efficient Exploration
- ▶ *Guided* Exploration by Teacher with Diverse and Adaptive Meta Policies
- ▶ Teacher Explores Spaces *Globally*, far away from the Student's States

Thank you!