



# **Parrot: Data-Driven Behavioral Priors for Reinforcement Learning**

Avi Singh\*, Huihan Liu\*, Gaoyue Zhou, Albert Yu, Nicholas Rhinehart, Sergey Levine University of California, Berkeley

ICLR 2021

## Background





 $MDP: \mathcal{M} = (S, A, \rho, P, r, \gamma)$ 

sample inefficiency: agent needs millions of interactions with the environment to be well trained

## **Motivation**

In other machine learning fields, such as NLP or CV, pre-training on large, previously collected datasets to bootstrap learning for new tasks has emerged as a powerful paradigm to reduce data requirements when learning a new task.

#### try similarly useful pre-training for RL agents

#### Benefit:

- reduce interactions and make a faster trainning process
- applicable to similar tasks



ParN<sub>o</sub>C

Without Behavioral Prior

With Behavioral Prior

模式识别与神经





The main idea: utilize manipulation data from a diverse range of prior tasks to train our behavioral prior, and then use it to bootstrap exploration for new tasks.

Convert  $\mathcal{M} = (S, A, \rho, P, r, \gamma)$  to  $\mathcal{M} = (S, Z, \rho, P, r, \gamma)$ , then action *a* is obtained from *s* and *z* 



#### **Real NVP**

$$\begin{cases} y_{1:d} &= x_{1:d} \\ y_{d+1:D} &= x_{d+1:D} \odot \exp(s(x_{1:d})) + t(x_{1:d}) \end{cases}$$

$$\iff \begin{cases} x_{1:d} = y_{1:d} \\ x_{d+1:D} = (y_{d+1:D} - t(y_{1:d})) \odot \exp\left(-s(y_{1:d})\right) \\ \frac{\partial y}{\partial x^T} = \begin{bmatrix} \mathbb{I}_d & 0 \\ \frac{\partial y_{d+1:D}}{\partial x_{1:d}^T} & \text{diag}\left(\exp\left[s\left(x_{1:d}\right)\right]\right) \end{bmatrix}$$





(a) Forward propagation



enable an invertible mapping, which makes the RL agent still retains full control over the action space of the original MDP

$$a = f_{\phi}(z; s)$$
$$p_{\text{prior}}(a|s) = p_z(f_{\phi}^{-1}(a; s)) \left| \det \left( \frac{\partial f_{\phi}^{-1}(a; s)}{\partial a} \right) \right|$$

maximize the likelihood term to learn behavioral prior





$$z'_{d+1:D} = z_{d+1:D} \odot \exp(v(z_{1:d};\phi(s))) + t(z_{1:d};\phi(s)))$$

The v, t and  $\psi$  are functions implemented using neural networks





For a new task, initialize the RL policy to the base distribution used for training the prior, so that at the beginning of training,  $\pi_{\theta}(z|s) = p_z(z)$ 

Update  $\pi_{\theta}(z|s)$  with (s, z, s', r)



#### **Algorithm 1 RL with Behavioral Priors**

- 1: Input: Dataset  $\mathcal{D}$  of state-action pairs (s, a) from previous tasks, new task  $M^*$
- 2: Learn  $f_{\phi}$  by maximizing the likelihood term in Equation 2
- 3: for step k in  $\{1, ..., N\}$  do
- 4:  $s \leftarrow \text{current observation}$
- 5: Sample  $z \sim \pi_{\theta}(z|s)$
- 6:  $a \leftarrow f_{\phi}(z;s)$
- 7:  $s', r \leftarrow \text{Execute } a \text{ in } M^*$
- 8: Update  $\pi_{\theta}(z|s)$  with (s, z, s', r)
- 9: end for
- 10: **Return**: Policy  $\pi_{\theta}(z|s)$  for task  $M^{\star}$ .

 $p_{\text{prior}}(a|s) = p_z \left( f_{\phi}^{-1}(a;s) \right) \left| \det \left( \frac{\partial f_{\phi}^{-1}(a;s)}{\partial a} \right) \right|$ 





Method	difference
Combining RL with demonstrations	Requires collecting demonstrations for the specific task
Generative modeling and RL	Non-invertible, lack of the ability of retaining full control over the action space
Hierarchical learning	focuses on modeling the temporal structure
Meta-learning/meta-IL	Need to interact with the prior tasks, with access to rewards and additional samples





Problem setting:

- Initial positions of all objects are randomized, and must be inferred from visual observations
- Not all objects in the scene are relevant to the current task
- A reward of +1 is provided when the objective for the task is achieved, and the reward is zero otherwise.





Prior-explore: While collecting data, an action is executed from the prior with probability  $\epsilon$ , else an action is executed from the learned policy





- The size of the dataset positively correlates with performance, but about 10K trajectories are sufficient for obtaining good performance, and collecting additional data yields diminishing returns
- Note that initializing with even a smaller dataset size (like 5K trajectories) yields much better performance than learning from scratch





The authors suspect this is due to the fact that pick and place tasks involve a completely new action (that of opening the gripper), which is never observed by the prior if it is trained only on grasping

# Thanks