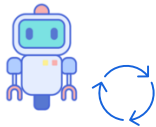


High-level policy

$$\hat{\mathbf{u}} = \pi_u(\cdot | s)$$

Offline RL



Inference Training

**We need:**

High-level Dataset

$$(s_0, \hat{u}_0, \hat{R}_0, s_1, \dots)$$

**We have:**



Offline Trajectory Data

$$(s_0, a_0, R_0, s_1, \dots)$$

Low-level policy

$$\mathbf{a} = \pi_l(\cdot | s, \hat{\mathbf{u}})$$