

题目:

Precise and Dexterous Robotic **Manipulation** via Human-in-the-Loop Reinforcement Learning

Related Work: RLPD: 加入先验数据的强化学习

SERL: 加入人类演示的强化学习

Intervention 🐮

HG-DAgger: 人类修正, 但训练有监督网络

RL 🐮

[HIL-SERL: Precise and Dexterous Robotic Manipulation via Human-in-the-Loop Reinforcement Learning](#)

数学模型

Markov Decision Process

$$\text{MDP } \mathcal{M} = \{S, \mathcal{A}, \rho, \mathcal{P}, r, \gamma\}$$

S 状态空间

\mathcal{A} 动作集合

ρ 初始状态分布

\mathcal{P} 转移概率

r 奖励函数

γ 折扣因子 (调整长期与近期)

Policy: $\pi \rightarrow E[\sum_{t=0}^H \gamma^t r(s_t, a_t)] \uparrow$

训练过程 $\rightarrow \pi$

Target Policy \nearrow On - Policy
Behavior Policy \searrow Off - Policy

常见的 Policy: $\pi(a|s)$ 是一个高斯分布, 去训练分布参数

Sparse reward function: 何为**稀疏**

RLPD:

$$\mathcal{L}_Q(\phi) = E_{s,a,s'} \left[\left(Q_\phi(s, a) - (r(s, a) + \gamma E_{a' \sim \pi_\theta} [Q_\phi(s', a')]) \right)^2 \right]$$

$$\mathcal{L}_\pi(\theta) = -E_s \left[E_{a \sim \pi_\theta(a)} [Q_\phi(s, a)] + \alpha \mathcal{H}(\pi_\theta(\cdot|s)) \right],$$

System Overview



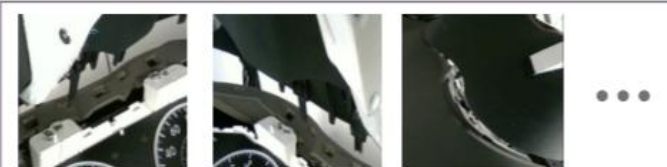
北京航空航天大学
BEIHANG UNIVERSITY

底层控制器
Careful
数据捕捉

Modularized Environment

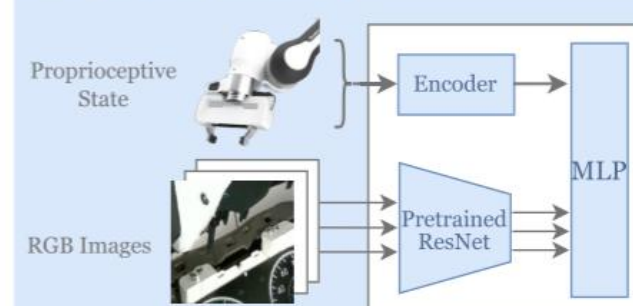


Single and Dual Arm Support

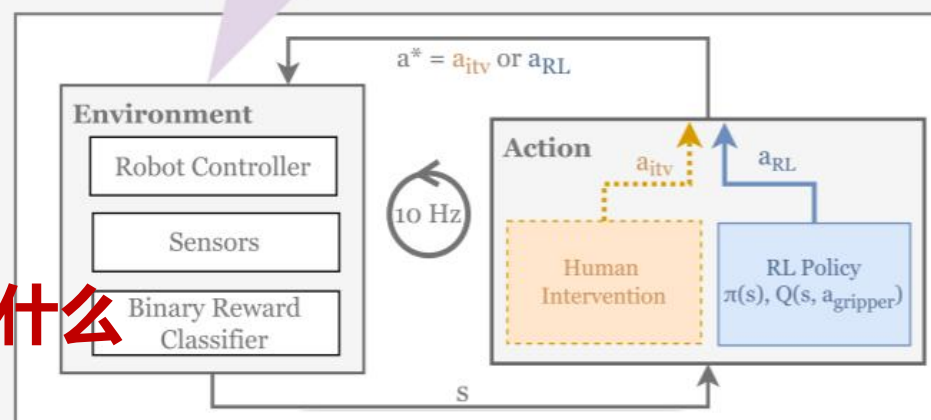


Multiple Camera Support

Network Architecture



Actor Process



Policy Transitions

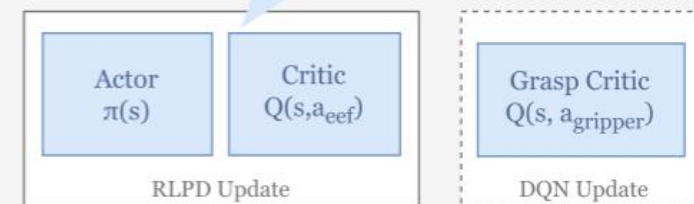
Interventions



Policy Transitions

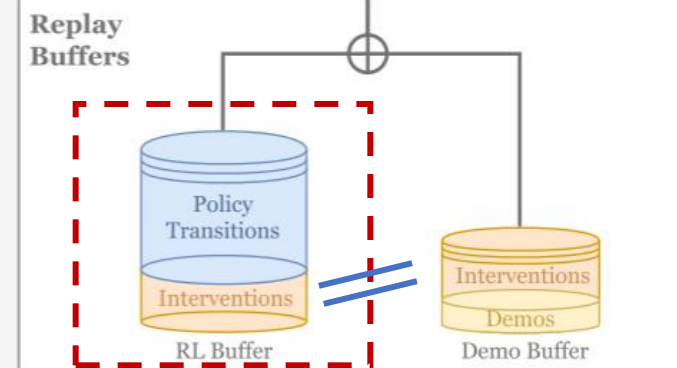
Learner Process

Policy
Parameters



Training Batch

RLPD



SERL的区别

Encoder 是什么

Interventions
>
RL Policy

Sample-efficient Contact-rich tasks

1. Pretrained Vision Backbones

ImageNet \longrightarrow ResNet-10

与proprioceptive数据融合

2. Sparse Reward Function

3. Downstream Robot

相对末端坐标系描述感知信息

阻抗控制器力控

\longrightarrow Wrist Camera

4. Gripper Control

离散动作难训练:

弄一个离散动作MDP使用DQN训练

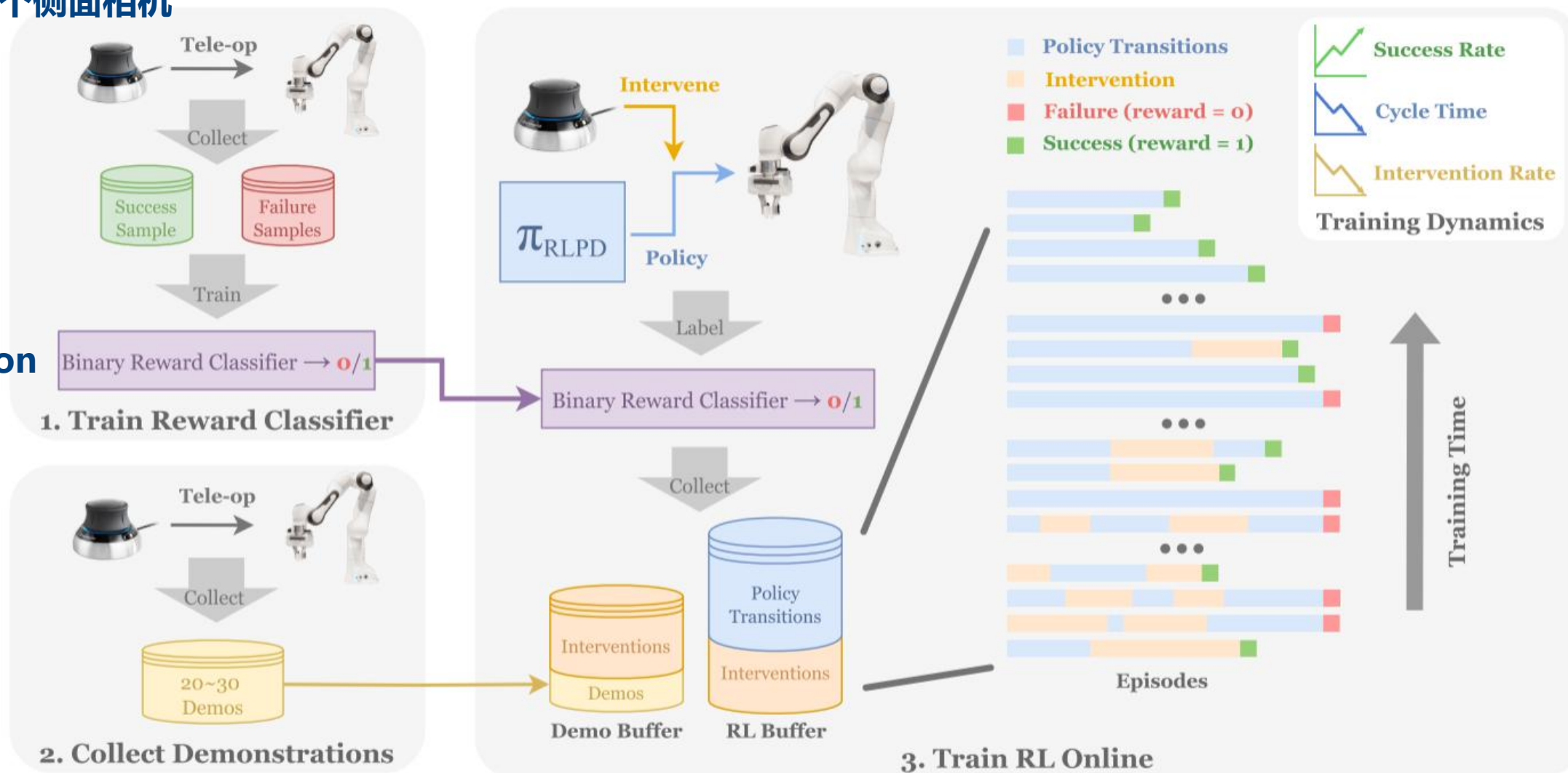
\longrightarrow 设计小负惩罚, 防止不必要的动作

Training Process



腕部相机 + 几个侧面相机

Reward Function



避免早期的长期稀疏干预

并非所有任务都能HIL

本身需要一个reward function

并没有进行长任务和广泛泛化性的实验

什么是BC

什么是BC