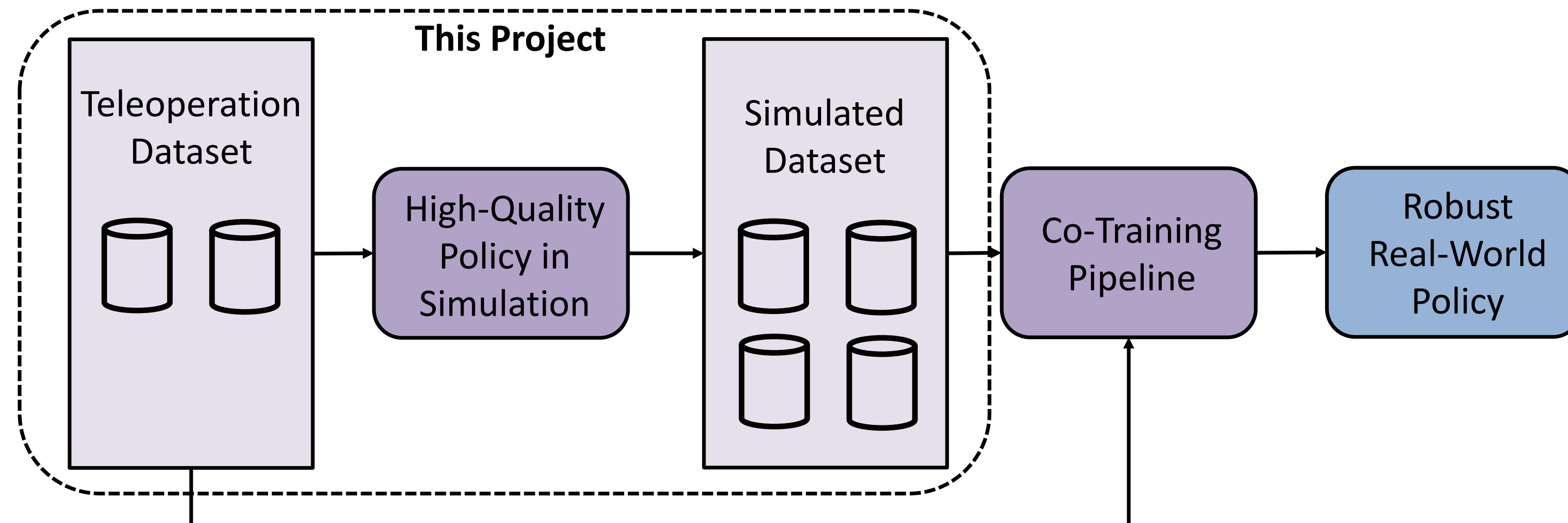


## Background

### Goals

Train an **autonomous** agent to perform **dexterous manipulation** tasks with the DexNex hardware setup.

Use **Sim-and-Real Co-Training** [1] to train a **more robust** policy that performs well on a **varied set** of real-world manipulation tasks.



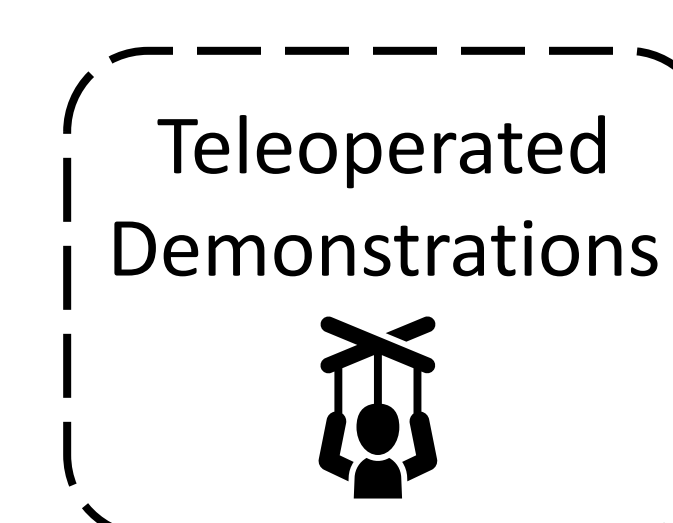
## Model Framework

### A. Specifications

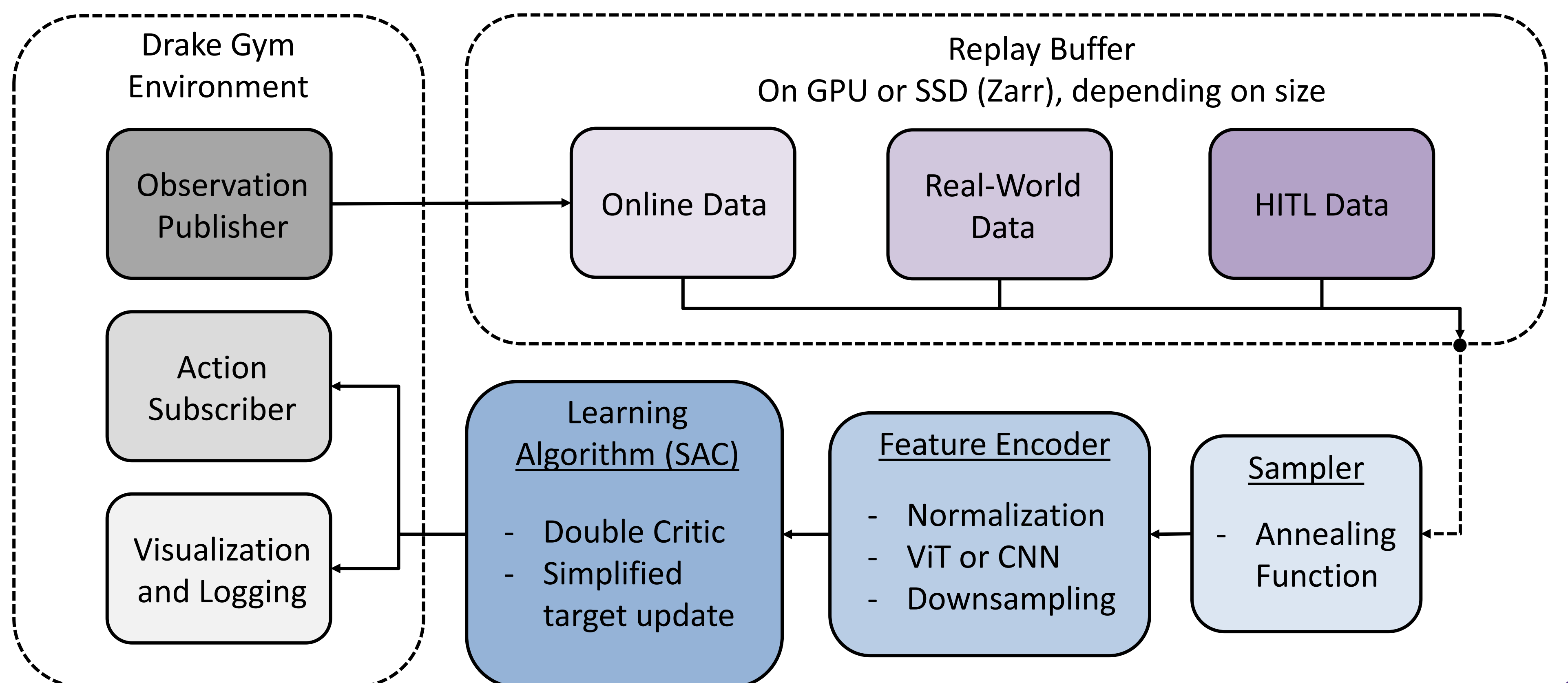
Algorithm Type: SAC [2]

Simulator: Drake [3]

### B. Pre-training



### C. Training



## Pipeline Testing

### Sparse Reward Function

$$r = r_f(p_5) \quad r_f(p_5) = \begin{cases} R, & \text{if } p_5 \text{ in right box} \\ 0, & \text{otherwise} \end{cases}$$

### Algorithm 1 Soft Actor-Critic (SAC)

```

1: Initialize parameter vectors  $\psi, \theta_1, \theta_2, \phi$ 
2: for each iteration do
3:   for each environment step do
4:      $a_t \sim \pi_\phi(a_t | s_t)$ 
5:      $s_{t+1} \sim p(s_{t+1} | s_t, a_t)$ 
6:      $D \leftarrow D \cup \{(s_t, a_t, r(s_t, a_t), s_{t+1})\}$ 
7:   end for
8:   for each gradient step do
9:      $\theta_i \leftarrow \theta_i - \lambda_Q \nabla_{\theta_i} J_Q(\theta_i)$ , for  $i \in \{1, 2\}$ 
10:     $\phi \leftarrow \phi - \lambda_\pi \nabla_\phi J_\pi(\phi)$ 
11:     $Q_{\bar{\psi}_i} \leftarrow \tau Q_{\psi_i} + (1 - \tau) Q_{\bar{\psi}_i}$ , for  $i \in \{1, 2\}$ 
12:   end for
13: end for
  
```

Figure 1: SAC Algorithm

### Shaped Reward Function

$$r = -\lambda_1 \|p_2 - p_1\| - \lambda_2 \|p_4 - p_3\| - \lambda_3 \|p_6 - p_5\| + r_f(p_5)$$

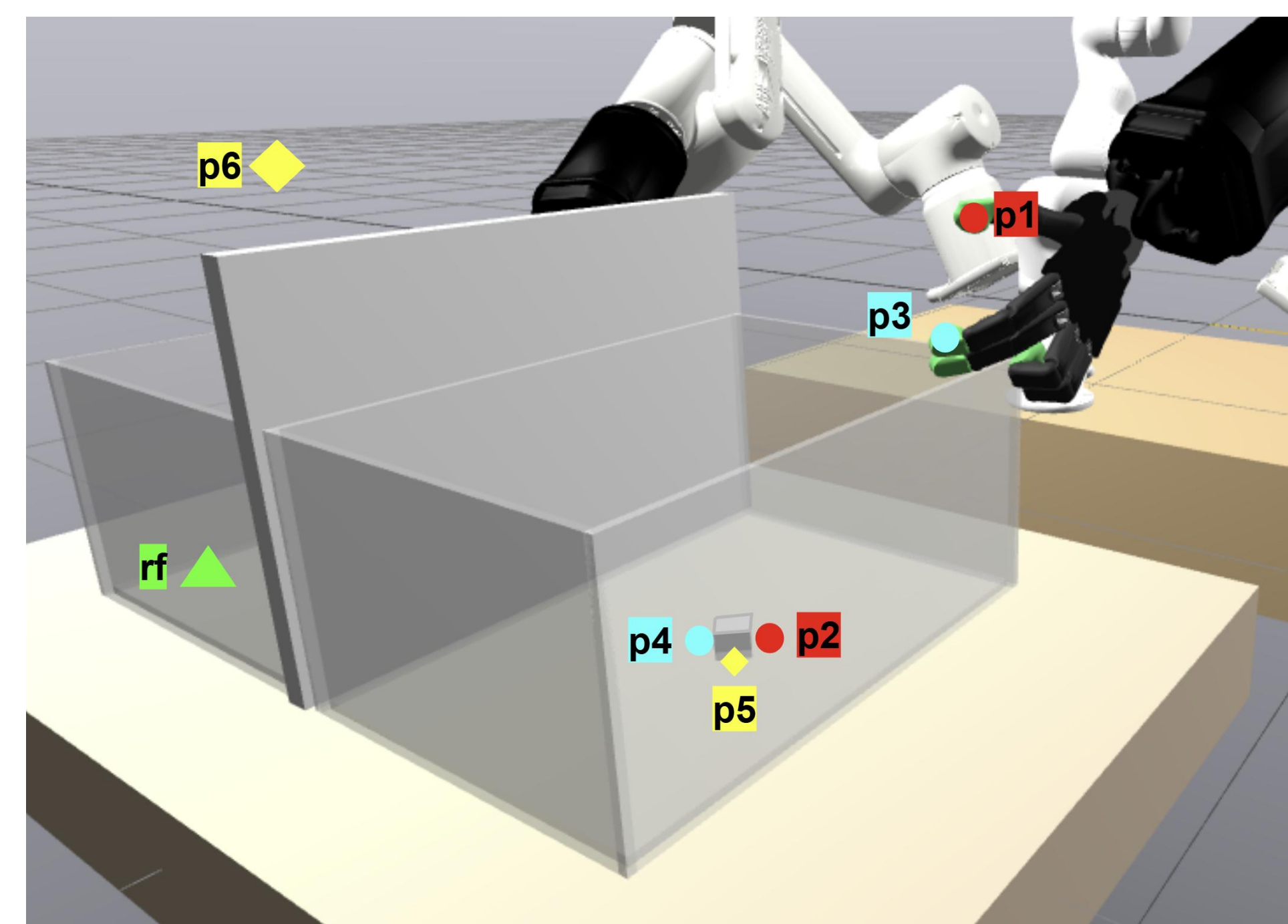


Figure 2: Reward Function Diagram

## Conclusions and Learning

- Difficulty learning with sparse reward
- Critic instability
- Computational bottlenecks
- Convergence to local minima



## Future Work

- Creation of HITL demonstrations
- Domain randomization
- Fix computational bottlenecks
- Simulated dataset generation
- Integration with Co-Training pipeline

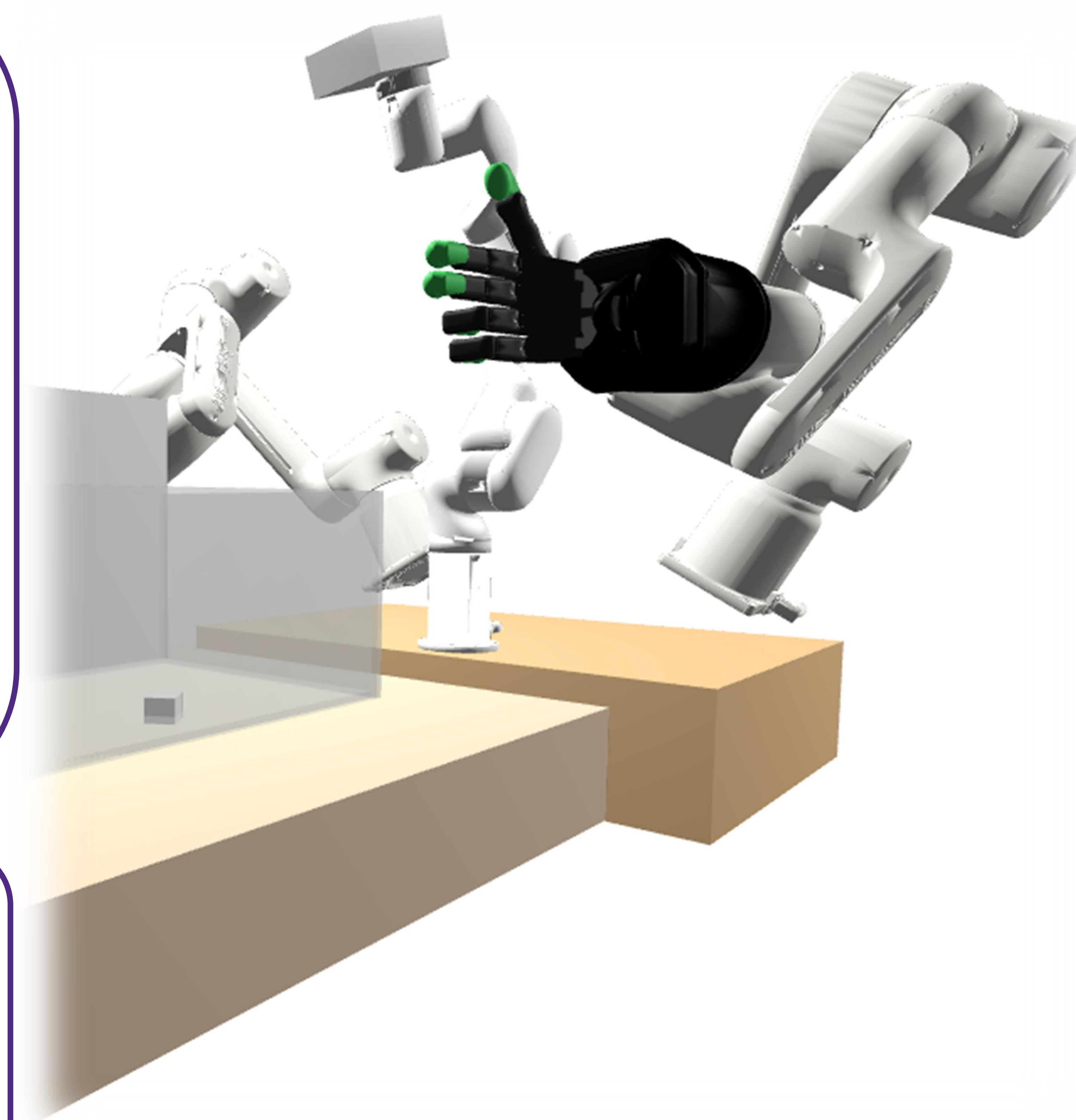
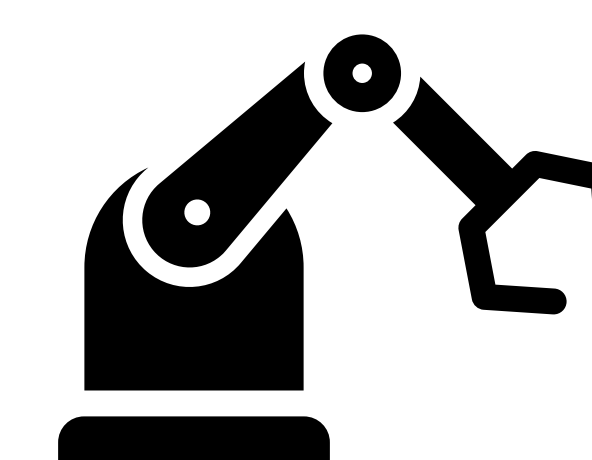


Figure 3: DexNex avatar visualized in Drake