

B ALGORITHMS

Algorithm 1 BC-SDRPRL: Pre-Training and Initialization

```

1: STEP I: Offline Pre-Training (Behavioral Cloning)
   Input:
      $B^E$ : Expert trajectory ( $S^E, A^E$ )
     T: Pre-training steps
2: Randomly initialize: Actor network  $a^A(s | \theta_\pi)$  and its target  $a^{A'}(s | \theta_{\pi'})$  weights.
3: for pre-training steps = 1, ..., T do
   Compute  $L_{BC}$  using eq. (1)
4: end for
   Output:
      $a^A(s | \theta_\pi)$ : Actor-network
      $a^{A'}(s | \theta_{\pi'})$ : Target actor network

5: STEP II: Offline IOHMM Training (Specialization)
   Input:
      $B^E$ : Expert trajectory ( $S^E$ )
     n: Number of hidden states
6: while expectation-maximization tolerance >  $1 * 10^1$  do
   Compute expectation-maximization  $\zeta$  using eq. (3)
7: end while
   Output:
      $\lambda$ : Trained IOHMM model parameters

8: STEP III: Offline Value Initialization & State Classification
   Input:
      $B_e$ : Expert trajectory ( $S^E$ )
     W: Warmup steps
9: Randomly initialize: Value network  $V(s | \theta_V)$  weights.
10: for warmup steps = 1, ..., W do
   Compute  $L_V$  using eq. (10)
11: end for
   Output:
      $V(s | \theta_V)$ : Value network
      $\hat{x}$ : Hidden states classification based on the value function
      $x^*$ : Specialized hidden state (abnormal state)

12: STEP IV: Online Critic Initialization
   Input:
     env: Environment model
     W: Warmup steps
13: Randomly initialize: Critic network  $Q(s, \pi(s | \theta_\pi) | \theta_Q)$  and its target  $Q'(s, \pi(s | \theta_\pi) | \theta_{Q'})$  weights.
14: Compute  $x_t$  using eq. (4).
15: for warmup steps = 1, ..., W do
16:   if  $x_t$  is  $x_t^*$  then
     Compute  $L_Q$  using eq. (7)
17:   end if
18: end for
   Output:
      $\hat{Q}(x^*, s, \pi | \theta_Q)$ : Critic network

```

Algorithm 2 BC-SDRPRL: Training and Inference

```

1: STEP I: Desired Hidden State Identification
   Input:
      $s_t$ : State of the system
      $\lambda$ : IOHMM trained parameters
      $x^*$ : Specified hidden state
2: Compute  $x_t$  using eq. (4).
3: if  $x_t$  is  $x^*$  then
    $\phi = \text{True}$ 
4: end if
   Output:
      $\phi$ : Boolean activation for the residual policy learning

5: STEP II: Deep Residual Policy Reinforcement Learning
   Input:
      $B^A$ : Initialize empty RL agent buffer
     env: Environment model
      $Q(s, \pi(s | \theta_\pi) | \theta_Q)$ ,  $a^A(s | \theta_\pi)$ : Pre-trained networks
      $\phi$ : Boolean activation for the residual policy learning
      $\tau$ : Target network soft update hyperparameter
     L: Training steps
6: Compute  $x_t$  using eq. (4).
7: for training steps = 1, ..., L do
8:   if  $\phi$  then
9:     Take superposed action combining expert's action and residual action based on the current state  $s_t$  as in eq. (2)
10:    Observe reward  $r$  and next state  $s_{t+1}$ 
11:    Add  $s_t, a_t^A, r_t, s_{t+1}$  in  $B^A$ 
12:    Sample batch of  $s_t, a_t^A, r_t, s_{t+1}$  from  $B^A$ 
13:    Compute  $L_Q$  and  $L_A$  using eq. (7) and eq. (9), respectively.
14:    Update actor and critic networks using Adam optimizer.
15:    Update target networks using eq. (15):
        
$$\theta_{\pi'} \leftarrow \tau \theta_\pi + (1 - \tau) \theta_{\pi'}$$

        
$$\theta_{Q'} \leftarrow \tau \theta_Q + (1 - \tau) \theta_{Q'} \tag{15}$$

16:   else
17:     Take expert's action
18:   end if
19: end for
   Output:
      $\pi_\theta^*(s_t)$ : Optimal control policy from eq. (2)

```
