

## Abstract

This paper proposes an actor-critic algorithm for controlling the temperature of a battery pack using a cooling fluid. This is modeled by a coupled 1D partial differential equation (PDE) with a controlled advection term that determines the speed of the cooling fluid. The Hamilton-Jacobi-Bellman (HJB) equation is a PDE that evaluates the optimality of the value function and determines an optimal controller.

- We propose an algorithm that treats the value network as a Physics-Informed Neural Network (PINN) to solve the continuous-time HJB equation
- We derive a control function from the HJB equation.

Our experiments show that a hybrid-policy method that updates the value network using the HJB equation and updates the policy network identically to PPO achieves the best results in the control of this PDE system.

## The 1D pack cooling problem

Modelled by the following coupled PDE

$$\begin{aligned} u_t(x, t) &= -D(x, t)u_{xx}(x, t) + h(x, t, u(x, t)) + \frac{1}{R(x, t)}(u - w) \\ w_t(x, t) &= -\sigma(t)w_x(x, t) + \frac{1}{R(x, t)}(w - u) \end{aligned}$$

with the following boundary conditions

$$\begin{aligned} u_x(0, t) &= u_x(1, t) = 0 \\ w(0, t) &= U(t) \end{aligned}$$

where

$u$ : Temperature distribution across battery pack  
 $w$ : Temperature distribution across cooling fluid  
 $D$ : Thermal diffusion constant  
 $R_T$ : Thermal resistance  
 $h$ : Internal heat generation in the battery pack  
 $U$ : temperature of the cooling fluid at the boundary  
 $\sigma$ : Transport speed of the cooling fluid (controller)

Objective of the controller:

$$\text{Maximize } \int_0^{\infty} \int_0^1 \gamma^t (-u(x, t)^2) dx dt$$

Environment parameters:

$$\begin{aligned} h(x, t, u(x, t)) &= e^{0.1u(x, t)} \\ u(x, 0) &= \sum_{n=0}^9 C_n \cos(\pi n x) \\ \Delta x &= 0.01, \Delta t = 0.01 \\ U(t) &= -5.0, D(x, t) = 0.01, R(x, t) = 2.0 \end{aligned}$$

## The HJB equation

Consider a controlled dynamical system modeled by the following equation:

$$\dot{x} = f(x, u), \quad x(t_0) = x_0$$

In control theory, the optimal value function  $V^*(x)$  is useful towards finding a solution to control problems:

$$V^*(t) = \sup_{\sigma} \frac{1}{\Delta t} \int_{t_0}^t \gamma^{\Delta t} L(x(\tau; t_0, x_0, \sigma(\cdot)), \sigma(\tau)) d\tau$$

where  $L(x, \sigma)$  is the reward function,  $\Delta t$  is the time step size for numerical simulation, and  $\gamma$  is the discount factor.

**Theorem 2.1.** A function  $V(x)$  is the optimal value function if and only if:

- $V \in C^1(\mathbb{R}^n)$  and  $V$  satisfies the Hamilton-Jacobi-Bellman (HJB) Equation

$$(\gamma - 1)V(x) + \sup_{\sigma \in U} \{L(x, \sigma) + \gamma \Delta t \nabla_x V^T(x) f(x, \sigma)\} = 0$$

- For all  $x \in \mathbb{R}^n$ , there exists a controller  $\sigma^*(\cdot)$  such that:

$$\begin{aligned} L(x, \sigma^*(x)) + \gamma \Delta t \nabla_x V^T(x) f(x, \sigma^*(x)) \\ = \sup_{\sigma} \{L(x, \sigma) + \gamma \Delta t \nabla_x V^T(x) f(x, \sigma)\} \end{aligned}$$

## Discretize the PDE in space to form an ODE

$$\begin{aligned} \dot{U} &= -DAU + h(U) + \frac{1}{R}(W - U) \\ \dot{W} &= -\sigma(t)BW + \frac{1}{R}(U - W) \end{aligned}$$

where  $AU$  approximates  $u_{xx}$  and  $BW$  approximates  $w_x$  using finite differences. Use this discretization to derive the HJB equation.

## HJB control of the pack cooling problem

**Theorem 4.1.** Let  $u(\cdot, t), w(\cdot, t) \in L_2[0, 1]$ . With  $\sigma(t) \in [0, 1]$  and the reward function  $L(u(\cdot, t), w(\cdot, t), \sigma(t)) = -\|u(\cdot, t + \Delta t)\|_2^2$ , the HJB equation for the 1D pack cooling problem is:

$$(\gamma - 1)V - \|u(\cdot, t + \Delta t)\|_2^2 + \langle V_u(u(\cdot, t), w(\cdot, t)), u(\cdot, t) \rangle + \frac{1}{R} \langle V_w(u(\cdot, t), w(\cdot, t)), u(\cdot, t) - w(\cdot, t) \rangle + \max(0, -\langle V_w(u(\cdot, t), w(\cdot, t)), w_x(\cdot, t) \rangle) = 0$$

where  $\|\cdot\|$  is the  $L_2[0, 1]$  norm and  $\langle \cdot, \cdot \rangle$  is the  $L_2[0, 1]$  inner product.

**Corollary 4.2.** Let  $u(\cdot, t), w(\cdot, t) \in L_2[0, 1]$ . With  $\sigma(t) \in [0, 1]$  and the reward function  $L(u(\cdot, t), w(\cdot, t), \sigma(t)) = -\|u(\cdot, t + \Delta t)\|_2^2$ , provided the optimal value function  $V^*(u, w)$  with  $V_w^*(\cdot, t) \in L_2[0, 1]$ , the optimal controller for the 1D pack cooling problem is:

$$\sigma^*(t) = \begin{cases} 1, & \langle V_w^*(u(\cdot, t), w(\cdot, t)), w_x(\cdot, t) \rangle < 0, \\ 0, & \text{otherwise} \end{cases}$$

where  $\langle \cdot, \cdot \rangle$  is the  $L_2[0, 1]$  inner product.

## New loss functions for the value network

Derived from the HJB equation:

$$\begin{aligned} MSE_f &= \frac{1}{T} \sum_{t=0}^{T-1} \left( (\gamma - 1)V - \|u(\cdot, t + \Delta t)\|_2^2 + \langle V_u(u(\cdot, t), w(\cdot, t)), u(\cdot, t) \rangle \right. \\ &\quad \left. + \frac{1}{R} \langle V_w(u(\cdot, t), w(\cdot, t)), u(\cdot, t) - w(\cdot, t) \rangle + \max(0, -\langle V_w(u(\cdot, t), w(\cdot, t)), w_x(\cdot, t) \rangle) \right)^2 \end{aligned}$$

At  $u(x, T) = 0, w(x, T) = -R(x, t)$ , we have:  $u(x, T) = 0$  and  $u_t(x, T) = 0$ . Thus,  $V(0, -R(x, t)) = 0$ .

$$MSE_u = (V(0, -R(x, t)))^2 = (V(0, -2))^2$$

At  $u(x, T) = 0, w(x, T) = -R(x, t)$ ,  $V$  achieves its global maximum.

$$MSE_n = \|\nabla_u V(0, -R(x, t))\|_2^2 + \|\nabla_w V(0, -R(x, t))\|_2^2$$

Controller:

$$\bar{\sigma}(t) = \begin{cases} 1, & \langle V_w(u(\cdot, t), w(\cdot, t)), w_x(\cdot, t) \rangle < 0, \\ 0, & \text{otherwise} \end{cases}$$

## Proposed algorithms

### HJB value iteration

- Initiate value network parameter  $\phi$
- Run the controller  $\bar{\sigma}(t)$  in the environment for  $T$  timesteps and observe samples  $\{(s_t, a_t, R_t, s_{t+1})\}_{t=1}^T$
- Compute the value network loss as:  $J(\phi) = MSE_f + MSE_u + MSE_n$
- Update  $\phi \leftarrow \phi - \alpha \nabla_{\phi} J(\phi)$
- Run steps 2-4 for multiple iterations

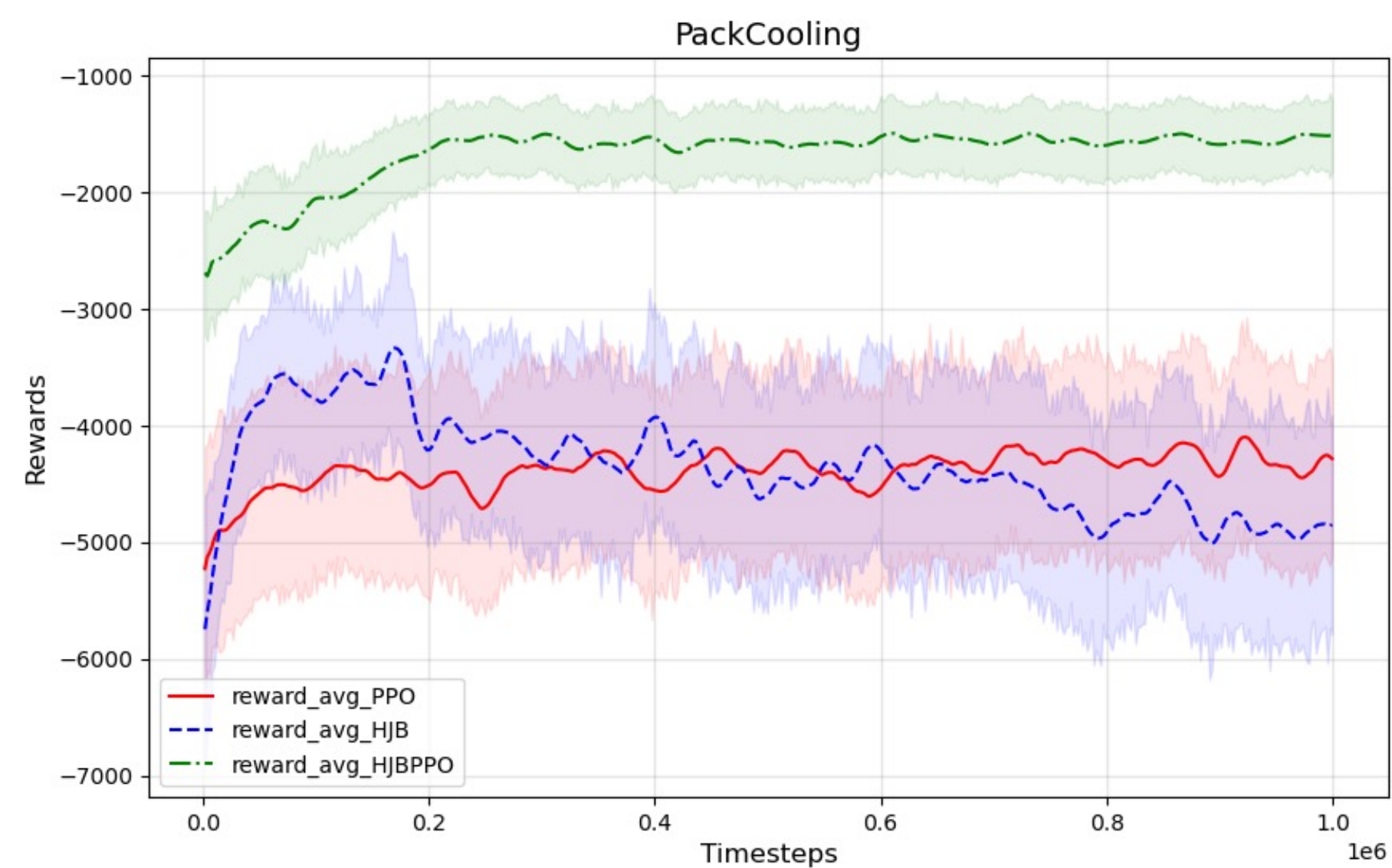
### HJBPPPO – Action Selection

- Retrieve state  $s_t$ , policy network parameter  $\theta$  and value network parameter  $\phi$
- Sample  $i \in \{0, 1\}$
- if  $i = 0$  then select the controller  $\bar{\sigma}(t)$
- else run policy  $\pi_{\theta}(\cdot | s_t)$
- end

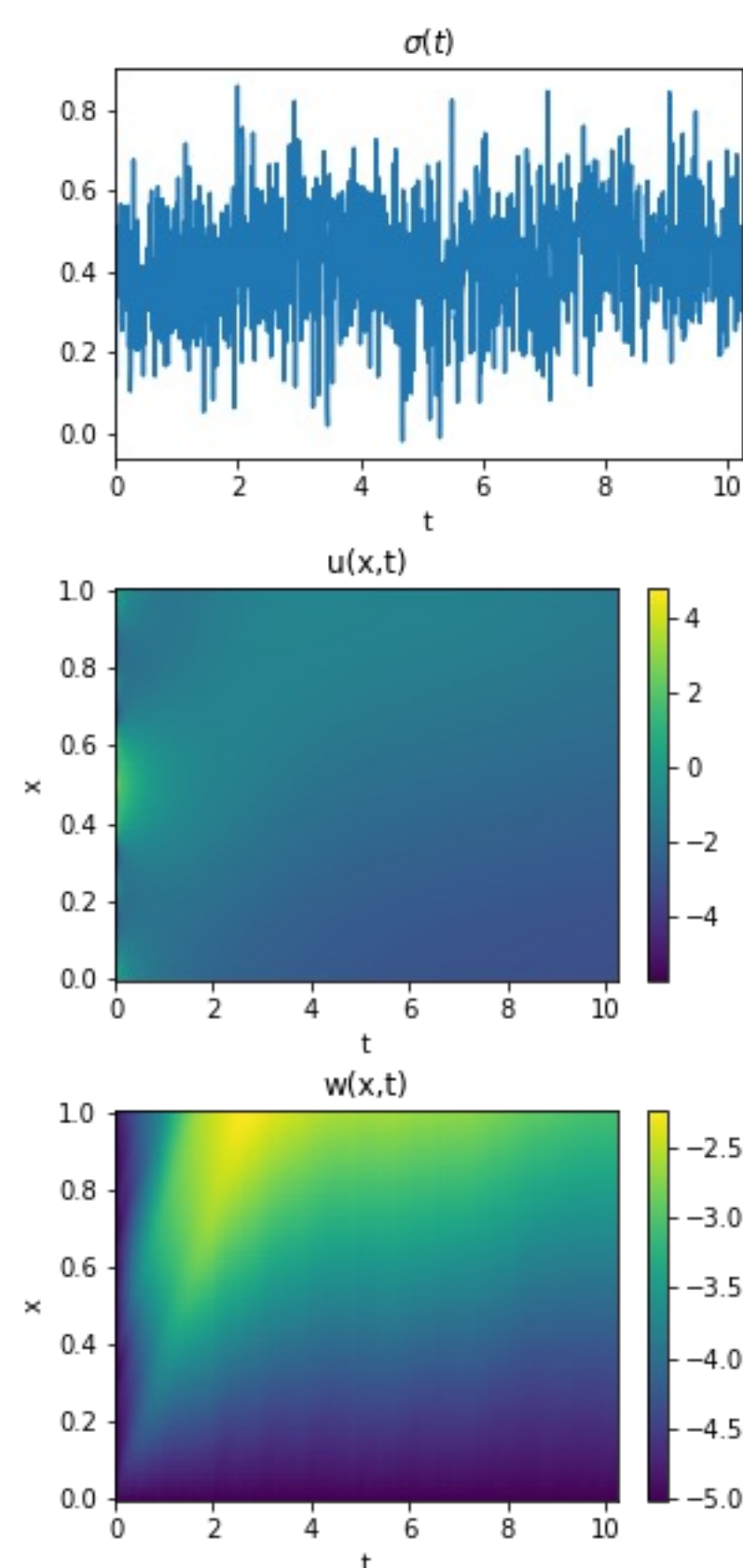
### Hamilton Jacobi Bellman Proximal Policy Optimization (HJBPPPO)

- Initiate policy network parameter  $\theta$  and value network parameter  $\phi$
- Run action selection as given earlier for  $T$  timesteps and observe samples  $\{(s_t, a_t, R_t, s_{t+1})\}_{t=1}^T$
- Compute the advantage  $A_t = \delta_t + (\gamma\lambda)\delta_{t+1} + \dots + (\gamma\lambda)^{T-t-1}\delta_{T-1}$  where  $\delta_t = R_t + \gamma V_{\phi}(s_{t+1}) - V_{\phi}(s_t)$ ,  $\gamma$ : Discount factor ( $\approx 0.99$ ) and  $\lambda$ : Smoothing factor ( $\approx 0.95$ )
- Compute  $r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)}$
- Compute the objective function of the policy network:  $L(\theta) = \frac{1}{T} \sum_{t=0}^{T-1} \min[r_t(\theta) A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) A_t]$  where  $\epsilon$ : clipping parameter ( $\approx 0.2$ )
- Update  $\theta \leftarrow \theta - \alpha_1 \nabla_{\theta} L(\theta)$
- Compute the value network loss as:  $J(\phi) = MSE_f + MSE_u + MSE_n$
- Update  $\phi \leftarrow \phi - \alpha_2 \nabla_{\phi} J(\phi)$
- Run steps 2-5 for multiple iterations

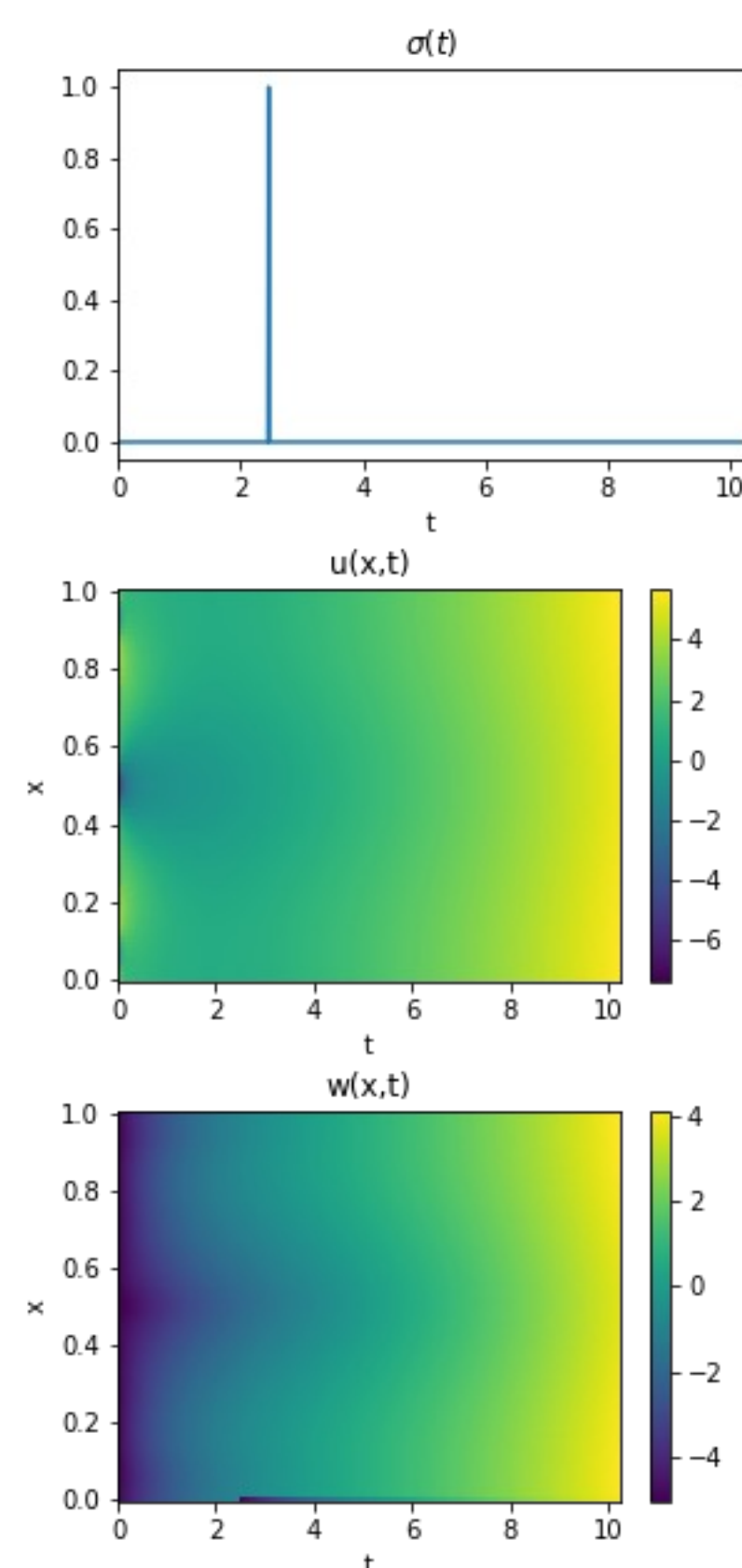
## Results



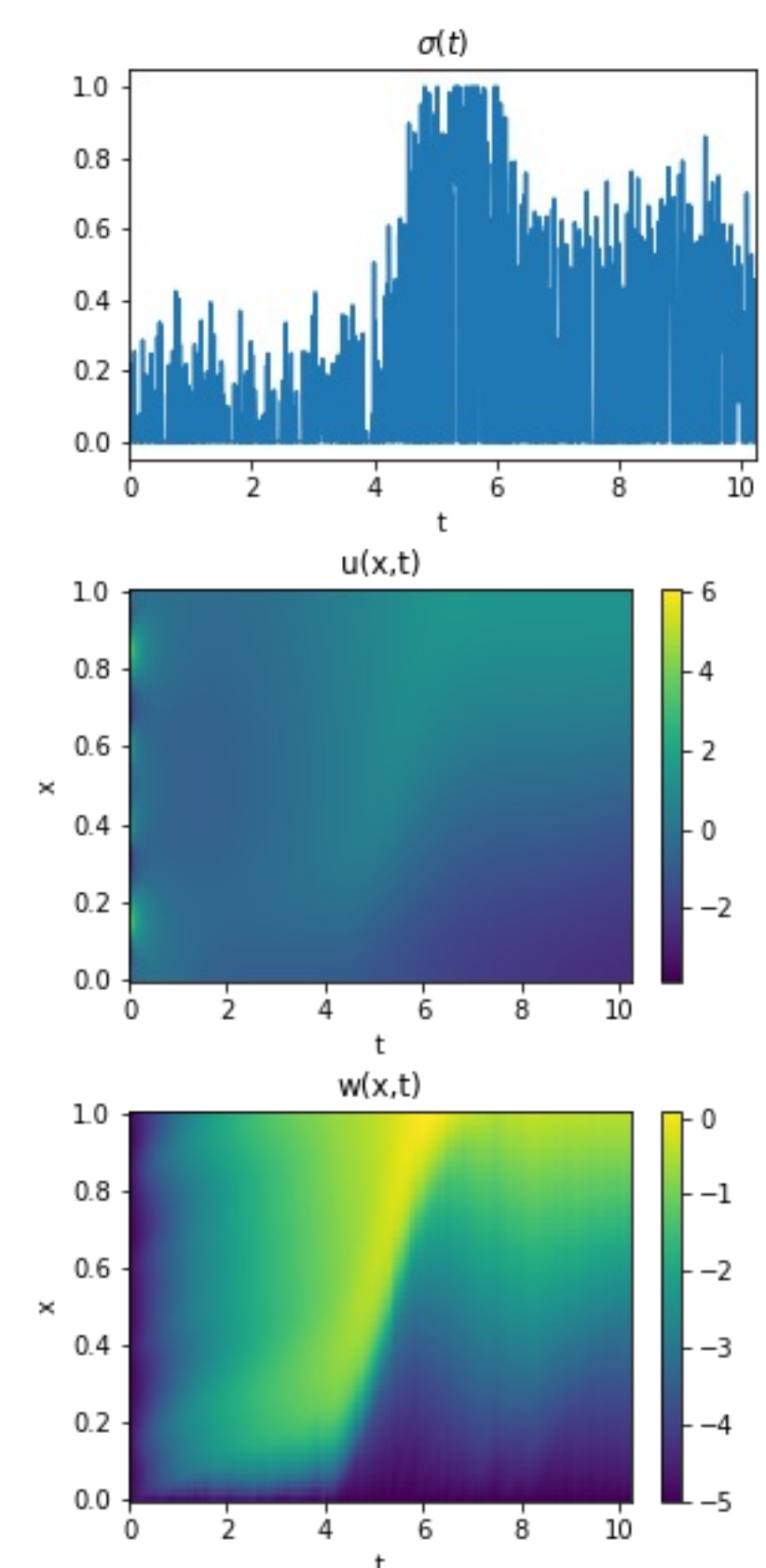
Reward curves of PPO (red), HJB value iteration (blue), and HJBPPPO (green) averaged over 5 seeds. Shaded area indicates 0.2 standard deviations.



Trajectory of PPO. Cumulative reward: -3970.02



Trajectory of HJB value iteration. Cumulative reward: -7294.51



Trajectory of HJBPPPO. Cumulative reward: -881.55