

Reinforcement Learning for PDE Control Problems

Nick Winovich[†], Bart van Bloemen Waanders[†]

Deepanshu Verma[‡], Lars Ruthotto[‡]

[†] Sandia National Laboratories¹
Center for Computing Research
Scientific Machine Learning Group

[‡] Emory University
Department of Mathematics
Department of Computer Science
Scientific Computing Group

SIAM-UQ 2022

¹ Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.

1 PDE Control Problems

- Motivating Examples
- Control Theory
- Prototype PDE System

2 Reinforcement Learning

- Actor-Critic Agent Models
- Proximal Policy Optimization
- Numerical Results

1 PDE Control Problems

- Motivating Examples
- Control Theory
- Prototype PDE System

2 Reinforcement Learning

- Actor-Critic Agent Models
- Proximal Policy Optimization
- Numerical Results

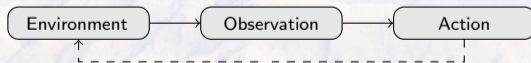
Control Problems in Scientific Computing

Control problems encapsulate a wide range of applications in scientific computing relevant to large-scale issues such as:

- Containment of wildfires, response to natural disasters
- Reducing impact of oil spills and contamination events

High-level points:

- Need for policy guidance/assistance with real-time decisions
- Control problems often arise naturally in the form of a sequence of “action-reaction” iterations:



Find an *optimal control policy* $a^*(t) = P(u(t), t)$ that induces a trajectory u^* of the system $\frac{\partial u}{\partial t}(t) = G(u(t), a(t), t)$ which minimizes the cost function:

$$J(t) = I(u(T), T) + \int_t^T L(u(t), a(t), t) dt \quad (1)$$

e.g. with a spatial loss defined by $L = \int_{\Omega_T} |u^+|^2 dm$ for a given target region $\Omega_T \subset \Omega$.

Given parameters θ , find a policy $a(x, y, t)$ such that the solution $u(x, y, t)$ to:

$$\begin{cases} \mathcal{L}_\theta[u] = F_\theta[a] & \text{in } \Omega \times [0, T] \\ u = g_d & \text{on } \Omega_d \\ \frac{\partial u}{\partial n} = g_n & \text{on } \Omega_n \end{cases} \quad (2)$$

minimizes the cost function $J = \int_0^T \int_{\Omega_T} |u^+|^2 dm dt$ for a given target region $\Omega_T \subset \Omega$.

- **Computational Time:** fast evaluation of policy required for real-time applications.
- **Scale:** states are governed by PDEs \rightarrow *infinite-dimensional*
- **Uncertainty:** precise source location and velocity field are unknown.
- **Local vs. Global solution method:**
 - **Local Method:** optimal policy for *fixed* problem configuration, e.g. NLP methods, PMP.
 - **Global Method:** optimal policy for *every* problem configuration, e.g. Dynamic Programming (HJB PDE), *Reinforcement Learning (RL)* ; suitable for real-time applications
- **Curse-of-Dimensionality:** HJB solvers suffer from CoD \rightarrow cost increases with dimension

PDE System

$$\frac{\partial u}{\partial t} + \mathbf{v} \cdot \nabla u - D \cdot \Delta u = f^{(t)} \quad \text{in} \quad \Omega \times [0, T] \quad \text{with} \quad D = 0.5$$

$$u = 0 \quad \text{on} \quad \{x = 0\} \cup \{y = 0\} \cup \{y = 1\} \quad \text{and} \quad \frac{\partial u}{\partial n} = 0 \quad \text{on} \quad \{x = 1\}$$

Initial Source

$$f^{(0)} = 5.0/\sigma \cdot \exp(-(|x - x_0| + |y - y_0|)/\sigma) \quad \text{with} \quad \sigma = 0.01$$

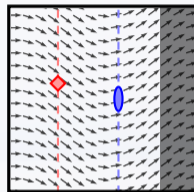
$$x_0 \sim \text{Uniform}(0.1, 0.25) \quad \text{and} \quad y_0 \sim \text{Uniform}(0.1, 0.9)$$

Velocity Field

$$\mathbf{v}_\delta(x, y) = \left(\sqrt{\eta^2 - \delta^2 \cdot \sin^2(2\pi \cdot [x - \phi_0])}, -\eta \cdot \sin(2\pi \cdot [x - \phi_0]) \right)$$

$$\text{where } \eta = 12.5, \quad \delta = 0.75, \quad \text{and} \quad \phi_0 \sim \text{Uniform}(0.0, -0.3)$$

Environment

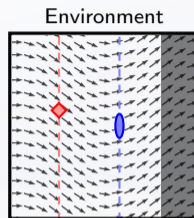


Sink Update

$$f^{(t+1)} = f^{(t)}(x, y) - \alpha \cdot \exp(-(|x - 0.5|/\sigma_x + |y - A_t|/\sigma_y))$$

with $\sigma_x = 0.025$, $\sigma_y = 0.05$, and $\alpha = 2.5$

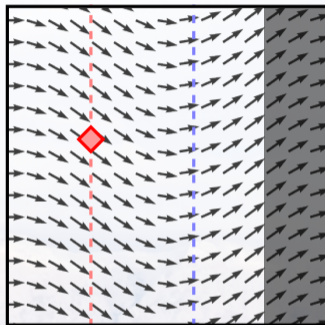
- Update RHS every $\Delta t = 0.02$ time-step (25 total updates)



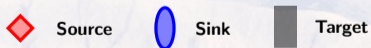
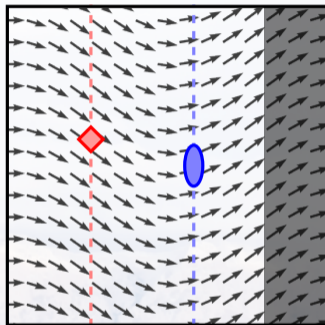
Control Decision

“Select the y -coordinate $A_t \in [0, 1]$ of the next sink location based on the current system state $U_t = \{u(x, y, t)\}_{(x, y) \in \Omega}$ ”

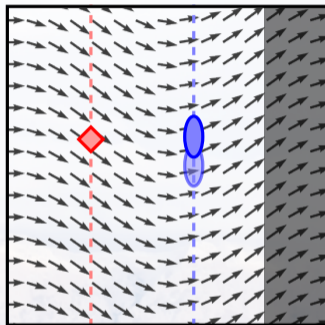
Prototype PDE Environment: Time-step 0



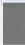


Prototype PDE Environment: Time-step 1

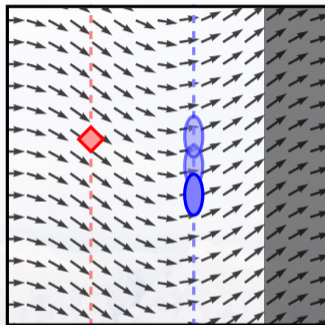


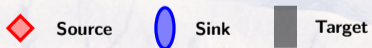
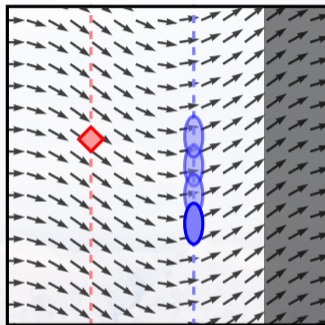
Prototype PDE Environment: Time-step 2



 **Source**  **Sink**  **Target**

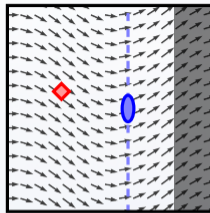
Prototype PDE Environment: Time-step 3





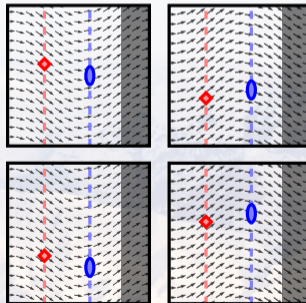
Local Solutions

- Fixed source location and velocity field
- Solutions are computed w.r.t. a fixed problem instance



Global Solutions

- Varying source location and velocity field
- Offline calibration yields approximate solutions for an entire family of problem instances
- Minimal overhead for performing run-time inference



1 PDE Control Problems

- Motivating Examples
- Control Theory
- Prototype PDE System

2 Reinforcement Learning

- Actor-Critic Agent Models
- Proximal Policy Optimization
- Numerical Results

Environment

The dynamics of a PDE control problem can be used to define the transition map and reward function for a MDP.

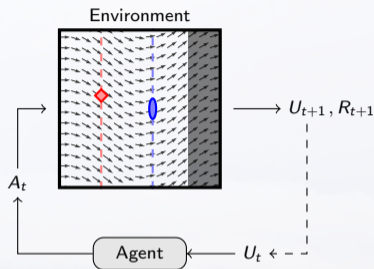
Observations

Decisions are often made based off of incomplete information regarding the underlying physical system.

Agent

The agent processes observation data in order to select an appropriate action for achieving the best possible outcome.

- The agent's policy is parameterized by values θ corresponding to a neural network architecture.
- This leads to a differentiable policy assignment rule.



- Observation data U_t may differ from the internal system state S_t
- Details regarding the calculation of the reward R_{t+1} are unknown to the agent

Objective Function

Given a parameterized policy π_θ and initial state distribution d_0 , we define the *episodic-return objective*:

$$J(\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \Delta t \cdot R_{t+1} \mid S_0 \sim d_0 \right]$$

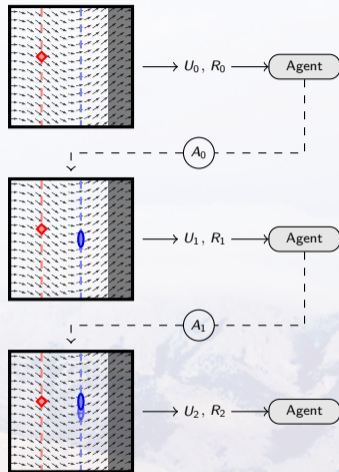
Policy Gradient Theorem

The gradient of the objective function can be expressed as follows:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \Delta t \cdot q_{\pi_\theta}(S_t, A_t) \cdot \nabla_\theta \log \pi_\theta(A_t | U_t) \mid S_0 \sim d_0 \right]$$

where $q_\pi(s, a) = \mathbb{E}_{\pi_\theta}[G_t | S_t = s, A_t = a]$.

- Optimal parameters are then approximated using gradient ascent.



Reference: <https://www.deepmind.com/learning-resources/reinforcement-learning-lecture-series-2021> [Lecture 9]

Variance Reduction

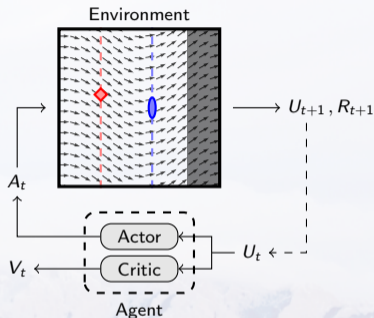
$$\mathbb{E}_{\pi_{\theta}} [B(U_t) \cdot \nabla_{\theta} \log \pi_{\theta}(A_t|U_t)] = 0$$

for any baseline $B(U_t)$ independent of the choice of action A_t .

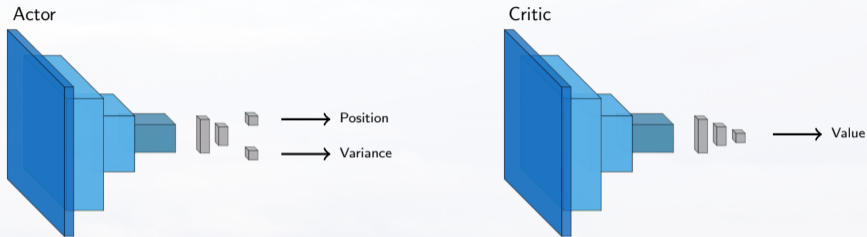
For a value estimate $V_{\pi_{\theta}}$ dependent only on the current state, we have:

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T \Delta t \cdot [q_{\pi_{\theta}}(S_t, A_t) - V_{\pi_{\theta}}(U_t)] \cdot \nabla_{\theta} \log \pi_{\theta}(A_t|U_t) \right] \\ &\approx \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T \Delta t \cdot [G_t - V_{\pi_{\theta}}(U_t)] \cdot \nabla_{\theta} \log \pi_{\theta}(A_t|U_t) \right] \end{aligned}$$

- Variance in the gradient approximation can be significantly reduced by providing an accurate value estimate.
- The action-value can be approximated by the observed return G_t .



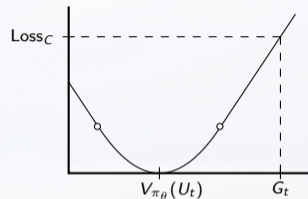
Reference: <https://www.deepmind.com/learning-resources/reinforcement-learning-lecture-series-2021> [Lecture 9]



- Both the actor and critic network architectures begin with convolutional layers to efficiently parse the spatially-structured observation data retrieved from the PDE environment.
- The actor provides predictions for the parameters of a probability distribution characterizing the desired action, while the critic network estimates the value of the current state.

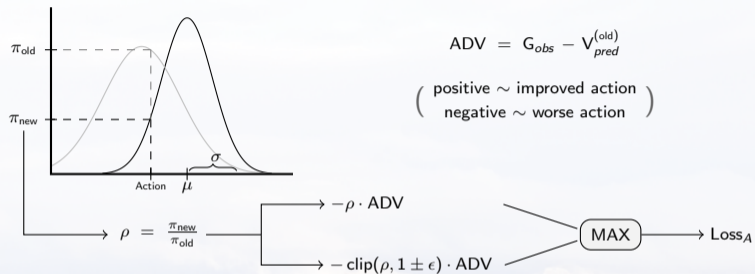
$$\nabla_{\theta} J(\theta) \approx \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T \Delta t \cdot \left[\underset{\substack{\uparrow \\ \text{observed return}}}{G_t} - \underset{\substack{\uparrow \\ \text{critic prediction}}}{V_{\pi_{\theta}}(U_t)} \right] \cdot \nabla_{\theta} \log \pi_{\theta}(A_t | U_t) \right]$$

$$\text{Loss}_C = \underbrace{|G_t - V_{\pi_{\theta}}(U_t)|^2}_{\text{"advantage"}}$$

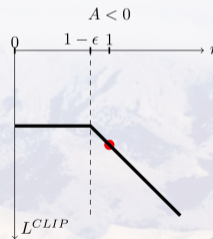
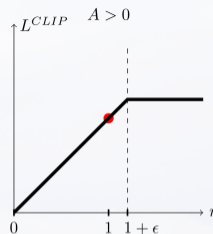


- The critic is trained to provide an accurate baseline for variance reduction.
- Huber Loss can be used to improve robustness of estimate with respect to outliers.

Actor Loss: Proximal Policy Optimization

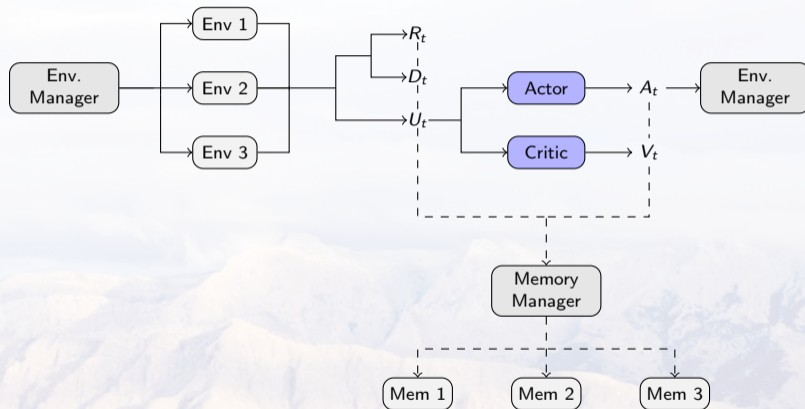


- PPO is designed to avoid over-tuning parameters so that the updated actor distributions remain relatively close to the previous distributions.
- Simplified version of previous work on *trust region policy optimization*.

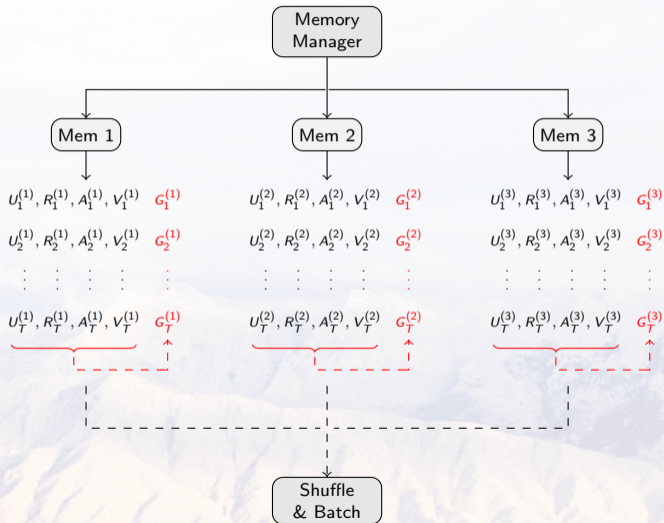


Reference: Proximal Policy Optimization Algorithms [<https://arxiv.org/pdf/1707.06347.pdf>]

Parallel Environment Workflow



Parallel Environment Workflow



Algorithm 1 Actor-Critic Agent: Proximal Policy Optimization

```
for  $n$  in  $\{1, \dots, \text{num\_episodes}\}$  do
    environment_manager.reset()
    states, actions, values, rewards  $\leftarrow$  run_parallel_envs()

    returns  $\leftarrow$  memory_manager.compute_returns(rewards)
    memory_manager.shuffle_and_batch()

    for  $k$  in  $\{1, \dots, \text{num\_batches}\}$  do
        ADV  $\leftarrow$  (returns[k] - values[k])
        ratio  $\leftarrow \pi_{\theta^A}(\text{actions}[k] \mid \text{states}[k]) / \text{probs}[k]$ 
        clipped  $\leftarrow -\text{ADV} \cdot \text{clip}(\text{ratio}, 1 - \epsilon, 1 + \epsilon)$ 
        unclipped  $\leftarrow -\text{ADV} \cdot \text{ratio}$ 

         $L_A \leftarrow \max[\text{clipped}, \text{unclipped}]$ 
         $L_C \leftarrow \frac{1}{2} \text{ADV}^2$  or  $\text{huber\_loss}(\text{values}[k], \text{returns}[k])$ 
        minimize( $L_A, \theta^A$ )
        minimize( $L_C, \theta^C$ )

    end for
end for
```

▷ Run simulations

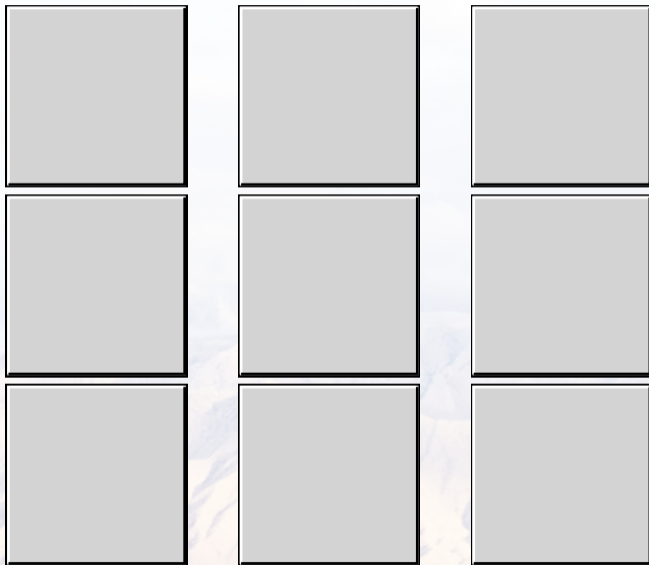
▷ Process results

▷ Evaluate actions

▷ Compute losses

▷ Update parameters

Initial Predictions for Control Policies



Rewards

-70.70	-36.21	-20.66
-87.85	-47.75	-28.22
-66.77	-54.20	-53.73

Network Predictions for Control Policies



Rewards

-66.09	-32.51	-14.10
-61.22	-36.60	-25.22
-31.11	-38.41	-47.21

Summary

- RL provides a potential framework for solving PDE control problems semi-globally.
- Run-time inference requires minimal overhead after offline calibration.
- Parallel implementations can be leveraged to reduce training time.

Future Work

- How well does this approach extend to more complex systems?
- How does this compare with local methods and semi-global dynamic programming?
- Can knowledge of the physical system be incorporated to improve model calibration?
 - HJB equations, FEM formulation

Thank you for your time.

Questions?