How to Leverage Non-Curated **Offline Data for Efficient RL?**

-- World Models to the Rescue

Generalist World Model Pre-Training for Efficient RL

Y. Zhao, A. Scannell, Y. Hou, T. Cui, L. Chen,

D. Büchler, A. Solin, J. Kannala, J. Pajarinen

Introduction

- Offline data boosts RL's sample efficiency.
- However, existing methods commonly need reward-labeled offline data or expert data.
- In practice, there is ample **non-expert** and **reward-free** offline data, but how to use it?

Table 1. Comparison with different policy learning methods that leverage offline data.					
Offline RL	Off2On RL	RLPD	MT IL	MT Offline RL	WPT (ours)
×	×	X	1	×	1
1	\checkmark	✓	×	\checkmark	✓
×	×	×	✓	✓	1
×	✓	\checkmark	×	×	1
×	×	×	✓	×	1
	Comparison with Offline RL ✓ ✓ × ×	Comparison with different policy le Offline RL Off2On RL X X X X X X X X X X X X X X X X X X X X X X X X X X X X	Comparison with different policy learning methodOffline RLOff2On RLRLPDXXX✓✓XXXXXXXXXXXXXXXXXXXXXXXX	Comparison with different policy learning methods that leve Offline RL Off2On RL RLPD MT IL X X X ✓ ✓ ✓ ✓ ✓ X X ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓ X ✓ ✓ ✓	Comparison with different policy learning methods that leverage offline data.Offline RLOff2On RLRLPDMT ILMT Offline RLXXX✓X✓✓X✓✓X✓✓✓✓X✓✓✓✓X✓✓✓✓X✓✓✓X✓✓✓X✓✓✓X✓✓✓X✓✓✓X✓✓✓X✓✓✓X✓✓✓X✓✓✓

Method



Results

1. A single pre-trained world model boosts RL training across multiple embodiments.



2. The pre-trained world model enables fast task adaptation.



THE UNIVERSITY

EDINBURGH

RL Fine-tuning

- An *embodiment-agnostic* world model is trained on the non-curated offline data. $\mathcal{L}(\theta) = \mathbb{E}_{p_{\theta}, q_{\theta}} \left[\sum_{t=1}^{-} \underbrace{-\ln p_{\theta}(o_t | z_t, h_t)}_{\text{pixel reconstruction loss}} + \beta \cdot \underbrace{\text{KL}(q_{\theta}(z_t | h_t, o_t) || p_{\theta}(z_t | h_t))}_{\text{latent state consistency loss}} \right]$
- During fine-tuning, we introduce experience retrieval to collect 1) task-relevant data from the offline data, 2) *execution guidance* to help exploration.
- Policy is updated by:

$$\mathcal{L}(v_{\phi}) = \mathbb{E}_{p_{\theta}, \pi_{\phi}} \left[-\sum_{t=1}^{H-1} \ln v_{\phi}(V_t^{\lambda} \mid s_t) \right] \quad \mathcal{L}(\pi_{\phi}) = \mathbb{E}_{p_{\theta}, \pi_{\phi}} \left[\sum_{t=1}^{H-1} \left(-v_t^{\lambda} - \eta \cdot \mathbf{H}[a_t \mid s_t] \right) \right]$$

3. Experience rehearsal and execution guidance boost performance.



Limitations

UNIVERSITY OF

- The world model uses the recurrent state space model, which limits the scalability.
- The method doesn't utilize action-free data.

NIVERSIT

IMPERIAL



