# Data Driven Reward Initialization for Preference based Reinforcement Learning

Mudit Verma, Subbarao Kambhampati

School of Computing & AI, Arizona State University

{muditverma, rao} @ asu.edu

## Objectives

Preference Based Reinforcement Learning has gained immense popularity for several tacit tasks however even State of the art algorithms have high variance in agent performance.

We investigate reasons for this high variance and tie it to high bias during reward initialization.

We show that this bias occurs with popular initialization techniques like Xavier-Uniform and Orthonormal and present a simple data driven initialization method to mitigate this.

## PbRL Objective

- **Approximating the Human Reward Model**

    > **Objective:** Compute probability for the human preferring trajectory $\tau_1$ over $\tau_2$:

$$P_\psi[\tau_0 \succ \tau_1] = \frac{\exp\left(\sum_t R_h(s_t^0, a_t^0)\right)}{\sum_{i \in \{0,1\}} \exp\left(\sum_t R_h(s_t^i, a_t^i)\right)}$$

    > **Objective:** The human reward model can then by learnt by minimizing the cross-entropy between the predictions made by the supervised learner and the ground truth human labels as follows:

$$\mathcal{L}_{CE} = - \mathop{\mathbb{E}}_{(\tau_0, \tau_1, y) \sim \mathcal{D}}[y(0)\log P_\psi[\tau_0 \succ \tau_1] + y(1)\log P_\psi[\tau_1 \succ \tau_0]]$$

## Motivation & Method

- **Issue of Reward Initialization**

    > Reward functions already suffer from the issue of degenracy. While this can be a source of high variance in the obtained policy, we find that initial bias of reward functions can also have unintended impacts the agent performance.

    > We find that data-independent initialization techqniues are limited to ensuring that only the parameters of the reward model are uniform, while the co-domain of the reward model is still "patchy" and biased.
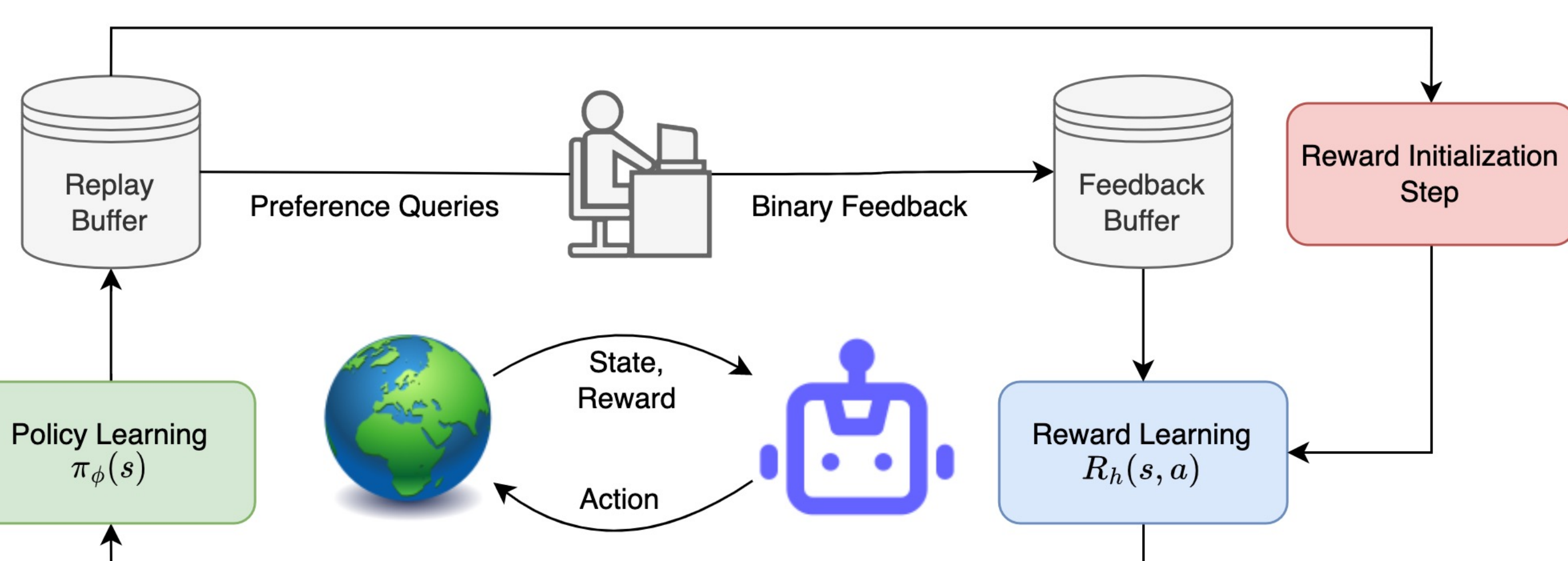
- **Proposed Data Driven Reward Initialization**

    > Pre-training of the Reinforcement Learning agent is common in PbRL setups like PEBBLE and SURF.

    > We make use of the pretraining phase and the randomly collected trajectories to collect "diverse" state-action samples.

    > We can then simply ensure that the reward model predicts some pre-defnined constant over all these states.

$$\mathcal{L}_P = \sum_{\tau_j \sim D_\tau} \sum_{i=1}^{i=|\tau_j|} \left\| R_\psi^h(s_i) - \epsilon \right\|_2$$
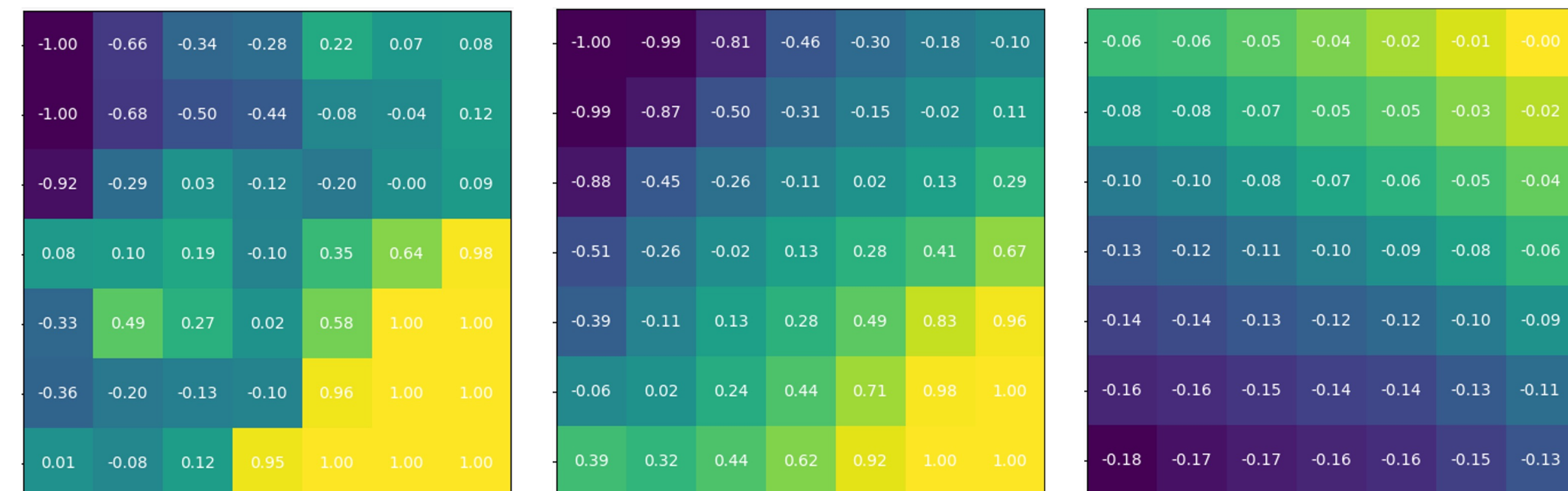
## Experiments & Results

- **Investigate via Empirical Evaluation**

    > **Question 1:** Empirically verify that seed of the experiment impacting reward initialization yields substantially different reward models.
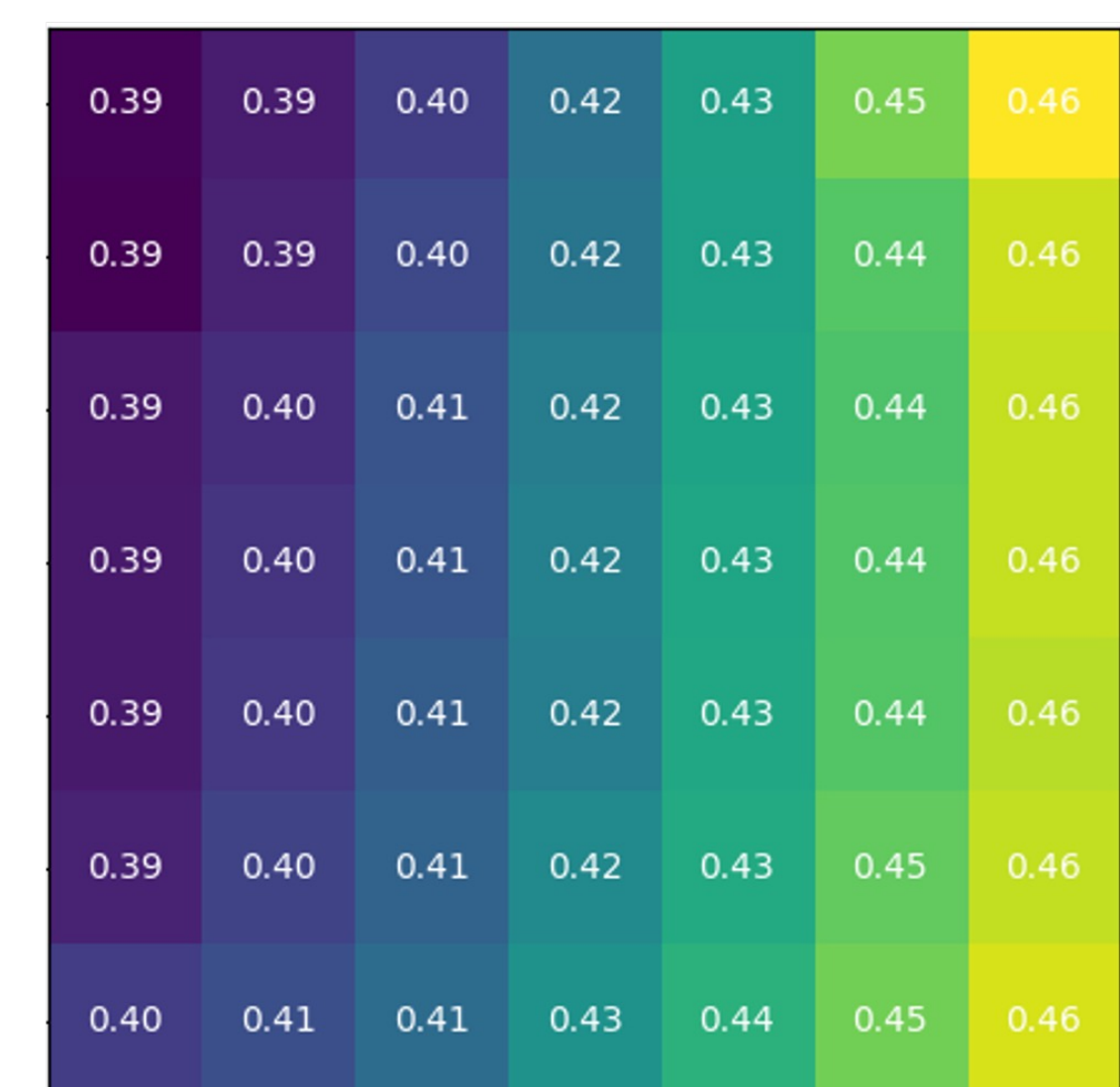
(a) Reward values for 7x7 Gridworld states over different seeds via Uniform Initialization

    > **Question 2:** Whether the proposed data-driven reward initialization scheme improves / stabilises the performance?
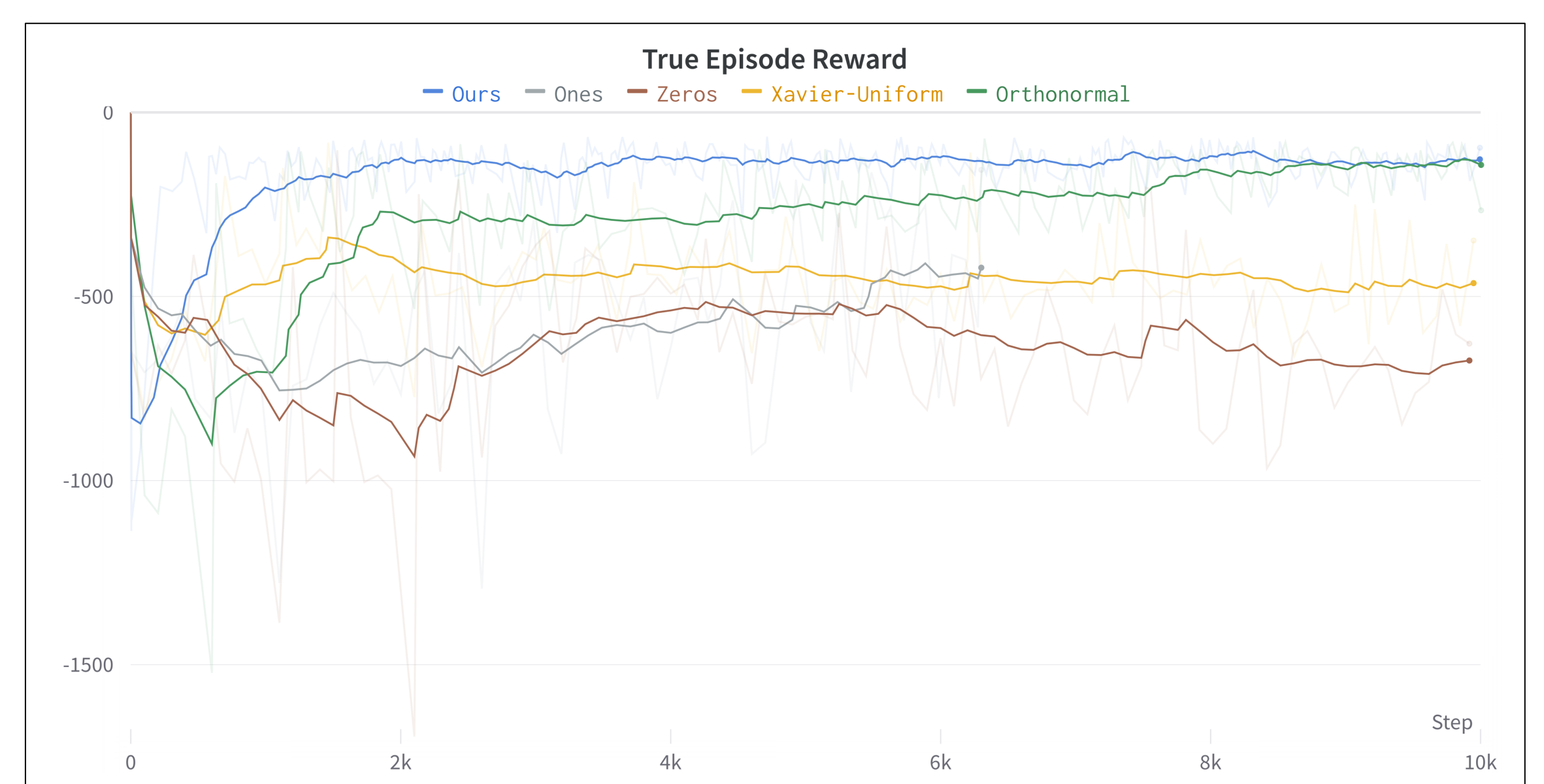
> Empirical results on 7x7 and 15x15 gridworlds with backbone PEBBLE are shown in figures (b-d).

> Using the data-driven initialization has zero cost to the human in the loop, and requires zero additional trajectory sampling (for algorithms that have a pre-training phase).
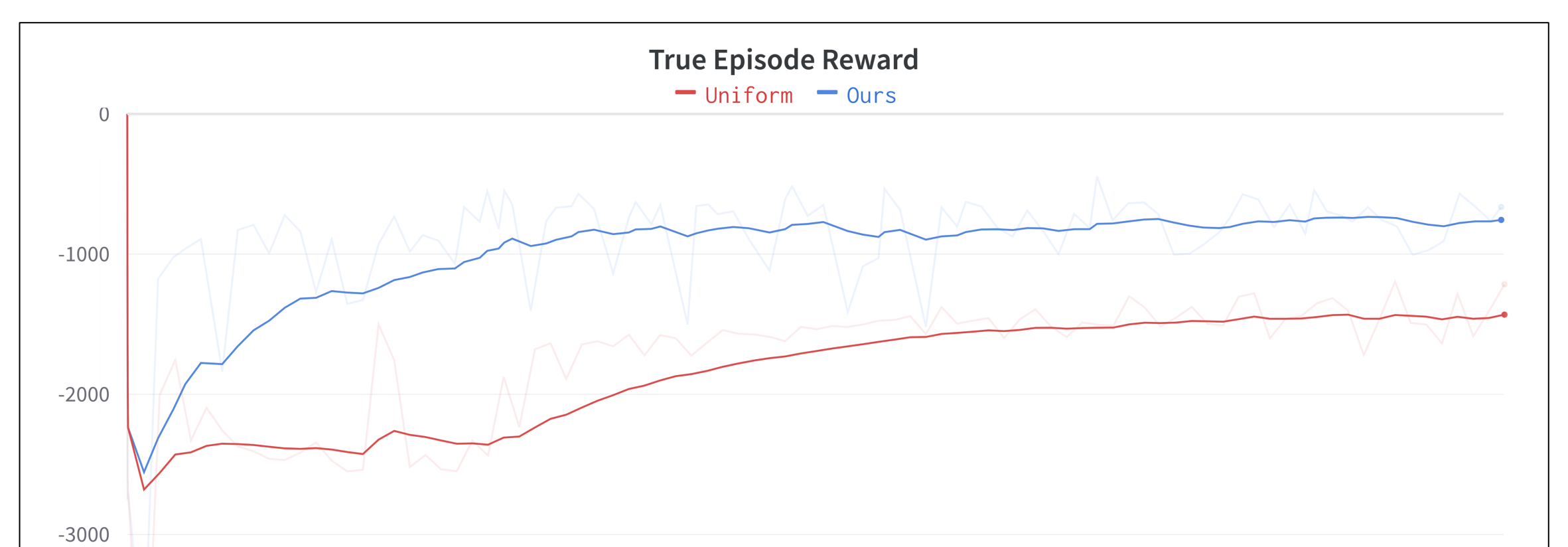
> We compare against baselines : Ones, Zeros, Xavier-Uniform, Orthonormal

(b) Reward values for 7x7 Gridworld states after data-driven Reward Initialization

(c) Learning curves for 7x7 gridworld across various initialization schemes.

(d) Learning curves for 15x15 gridworld comparing Ours v/s Uniform.

## Key Takeaways

1. Reward learning and Policy learning *is sensitive to Reward Initialization.*
2. Empirically, we show that several initialization schemes for the reward model results in high variance in agent performance.
3. We propose a pre-training data-driven reward initialization step that forces the reward-model co-domain to be uniform, minimizing any prior bias over trajectories.