



Vikas Dhiman¹, Shurjo Banerjee¹, Jeffrey M. Siskind², Jason J Corso¹
¹EECS, University of Michigan, Ann Arbor, MI. ²ECE, Purdue University, West Lafayette, IN

Goal-conditioned reinforcement learning (GCRL) addresses tasks where the desired goal can change for every trial. State-of-the-art algorithms model these problems such that the reward formulation depends on the goal rewards, to associate goals with high reward.

We propose a reformulation of goal-conditioned value functions for GCRL that yields a similar algorithm, while removing the dependence of reward functions on the goal. Our formulation thus obviates the requirement of reward-resampling that is needed by HER and its extensions.

Figure 1 displays eight different robotic tasks, arranged in two rows of four. The top row shows tasks performed by the Fetch robot, and the bottom row shows tasks performed by the Hand robot. The tasks are:

- FetchReach: A Fetch robot arm reaching for a red ball on a white block.
- FetchPush: A Fetch robot arm pushing a red ball on a white block.
- FetchSlide: A Fetch robot arm sliding a red ball on a white block.
- FetchPickAndPlace: A Fetch robot arm picking up a red ball from a white block.
- HandReach: A Hand robot arm reaching for a red ball on a white block.
- HandManipulateBlockRotateXYZ: A Hand robot arm manipulating a red block with a blue 'I' on it.
- HandManipulateEggFull: A Hand robot arm manipulating a colorful egg.
- HandManipulatePenRotate: A Hand robot arm manipulating a blue pen.

1. DDPG Loss
$$\mathcal{L}_{DDPG}(\theta_Q, \theta_\pi) = \mathbb{E}[(Q_m(s_t, a_t; \theta_Q) - y_t)^2]$$

$$y_t = R(s_t, a_t) + \gamma Q_{\text{tgt}}(s_{t+1}, \pi_{\text{tgt}}(s_{t+1}, \theta_\pi); \theta_{Q_{\text{tgt}}})$$
2. Step Loss
$$\mathcal{L}_{\text{step}}(\theta_Q) = (Q_*^P(s_{l-1}, a_{l1}, g_l; \theta_Q) - R(s_{l-1}, a_{l-1}))$$
3. FWRL Loss upper bound
$$\mathcal{L}_{\text{lo}} = \text{ReLU}[Q_{\text{tgt}}(s_t, a_t, g_w) + Q_{\text{tgt}}(s_w, \pi_t(s_w, g_{t+f}; \theta_\pi), g_{t+f}) - Q_m(s_t, a_t, g_{t+f})]^2$$
4. FWRL Loss lower bound
$$\mathcal{L}_{\text{up}} = \text{ReLU}[Q_m(s_t, a_t, g_w) + Q_{\text{tgt}}(s_w, \pi_t(s_w, g_{t+f}; \theta_\pi), g_{t+f}) - Q_{\text{tgt}}(s_t, a_t, g_{t+f})]^2$$

1. Hindsight Experience Replay (HER)[1]:	\mathcal{L}_{DDPG}
2. Path Reward Reinforcement Learning (Ours):	$\mathcal{L}_{DDPG} + \mathcal{L}_{\text{step}}$
3. Floyd-Warshall RL (FWRL [2]):	$\mathcal{L}_{DDPG} + \mathcal{L}_{\text{up}} + \mathcal{L}_{\text{lo}}$

1. With Goal Rewards: $R(s, a, g) = (0 \text{ if } s == g \text{ else } -1)$
2. Without Goal Rewards: $R(s, a, g) = -1$

Figure 1 displays the performance of three algorithms (Ours, HER, and FWRL) across four tasks: Fetch Push, Fetch Pull, Fetch Risk And Place, and Fetch Size. The figure is organized into a 4x2 grid of plots. The left column shows the 'Distance from goal (m)' and the right column shows the 'Success Rate (test)'. The x-axis for all plots is 'Epochs' (0 to 200) and 'Reward Computations' (0 to 10^6). The y-axis for the left column is 'Distance from goal (m)' (0 to 0.2 or 0.4) and for the right column is 'Success Rate (test)' (0 to 1). The legend indicates: Ours (red line), HER (blue line), and FWRL (grey line).

- Fetch Push:** Ours reaches the goal distance of 0m faster than HER and FWRL. In the success rate plot, Ours reaches a success rate of 1.0 much faster than the other two algorithms.
- Fetch Pull:** Ours reaches the goal distance of 0m faster than HER and FWRL. In the success rate plot, Ours reaches a success rate of 1.0 faster than HER and FWRL.
- Fetch Risk And Place:** Ours reaches the goal distance of 0m faster than HER and FWRL. In the success rate plot, Ours reaches a success rate of 1.0 faster than HER and FWRL.
- Fetch Size:** Ours reaches the goal distance of 0m faster than HER and FWRL. In the success rate plot, Ours reaches a success rate of 1.0 faster than HER and FWRL.

Figure 2: For the Fetch tasks, we compare our method (red) against HER (blue)(Andrychowicz et al., 2016) and FWRL (green) (Kaelbling, 1993) on the distance-from-goal and success rate metrics. Both metrics are plotted against two progress measures: the number of training epochs and the number of reward computations. Except for the Fetch Slide task, we achieve comparable or better performance across the metrics and progress measures.

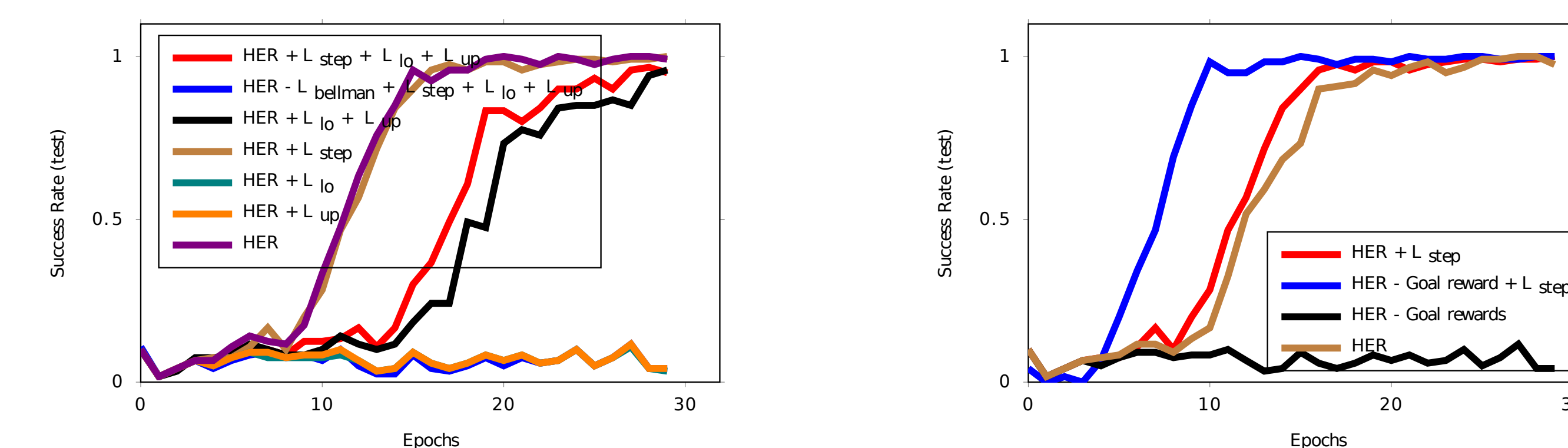
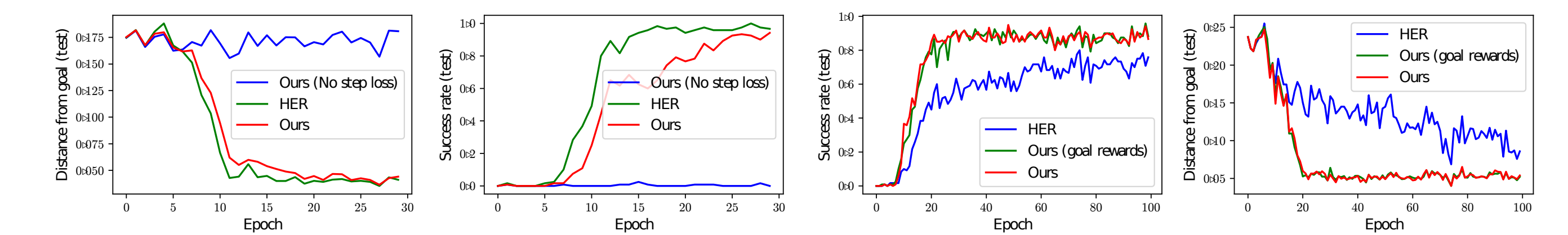


Figure 5 (left) Ablation on loss functions for Fetch Push task. The Floyd-Warshall inspired loss functions L lo and L up do not help much. L step helps a little but only in conjunction with HER [1]. (right) Even when the Goal rewards are removed from HER [1] training, the HER is able to learn only if the L step is added again.



(a) Do we really need the step-loss? (b) Effect of goal-rewards

Figure 3: (a) Effects of removing the step-loss from our methods. Results show that it is a critical component to learning in the absence of goal-rewards. (b) Adding goal-rewards to our algorithm that does have an effect further displaying how they are avoidable.

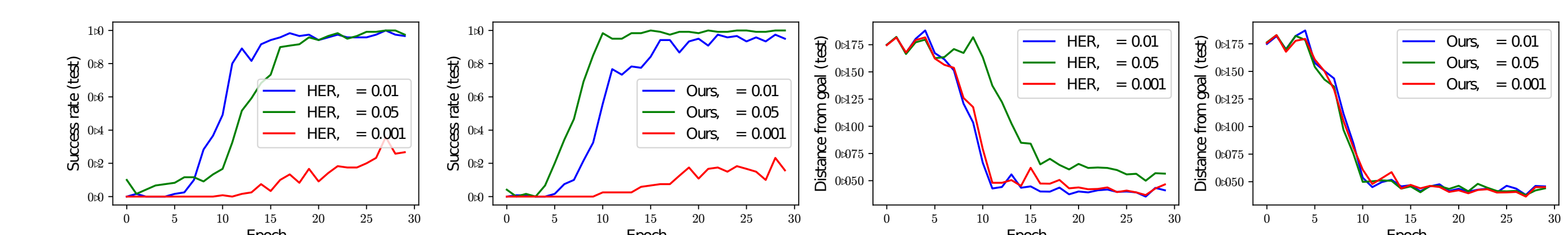


Figure 4: We measure the sensitive of HER and our method to the distance-threshold (δ) with respect to the success-rate and distance-from-goal metrics. Both algorithms success-rate is sensitive the threshold while only HER's distance-from-goal is affected by it.

In this work we pose a reinterpretation of goal-conditioned value functions and show that under this paradigm learning is possible in the absence of goal reward. This is a surprising result that runs counter to intuitions that underlie most reinforcement learning algorithms. In future work, we will augment our method to incorporate the distance-threshold information to make the task easier to learn when the threshold is high. We hope that the experiments and results presented in this paper lead to a broader discussion about the assumptions actually required for learning multi-goal tasks.

[1] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In *Advances in Neural Information Processing Systems*, pp. 5048–5058, 2017.

[2] Leslie Pack Kaelbling. Learning to achieve goals. In *IJCAI*, pp. 1094–1099. Citeseer, 1993.