

Density-based Curriculum for Multi-goal Reinforcement Learning with Sparse Rewards

Deyu Yang, Hanbo Zhang, Xuguang Lan*, Jishiyu Ding

Abstract—Multi-goal reinforcement learning (RL) aims to qualify the agent to accomplish multi-goal tasks, which is of great importance in learning scalable robotic manipulation skills. However, reward engineering always requires strenuous efforts in multi-goal RL. Moreover, it will introduce inevitable bias causing the suboptimality of the final policy. The sparse reward provides a simple yet efficient way to overcome such limits. Nevertheless, it harms the exploration efficiency and even hinders the policy from convergence. In this paper, we propose a density-based curriculum learning method for efficient exploration with sparse rewards and better generalization to desired goal distribution. Intuitively, our method encourages the robot to gradually broaden the frontier of its ability along the directions to cover the entire desired goal space as much and quickly as possible. To further improve data efficiency and generality, we augment the goals and transitions within the allowed region during training. Finally, We evaluate our method on diversified variants of benchmark manipulation tasks that are challenging for existing methods. Empirical results show that our method outperforms the state-of-the-art baselines in terms of both data efficiency and success rate.

I. INTRODUCTION

Reinforcement learning (RL) has attracted widespread attention in recent years with inspiring results in video games [1], [2] and complex robot manipulation tasks [3]–[5]. Despite its notable success, RL always needs arduous reward engineering, which might bias the learning and limit the scalability of the learned policies. Therefore, it’s more natural and convenient to use sparse reward settings in manipulation tasks, i.e., the robot will be awarded only when it reaches the final goal. However, the sparse reward makes it extremely hard to explore in the initial stage, resulting in few useful learning signals, and even hindering the agent from convergence.

Andrychowicz et al. [6] proposed a Hindsight Experience Replay (HER) to tackle sparse-reward RL in multi-goal robot manipulation tasks. By replacing desired goals with the randomly selected achieved goals, namely hindsight goals, of the past experiences, an agent can get much more positive reward signals and learn to master skills of reaching handy goals. However, this mechanism does not provide the guarantee of accomplishment of original tasks. The agent is able to generalize to finishing original tasks only when the distribution of hindsight goals is close enough to those originally desired goals. When there’s little or no

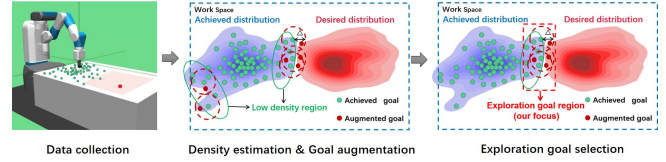


Fig. 1: Illustration of goal selection. We utilize density estimation and goal augmentation to select goals from the exploration goal region.

overlap between desired and hindsight goal distributions, the originally desired goals are too hard for HER agent to achieve with random hindsight goal selection.

In this paper, we propose to use a curriculum for the efficient exploration, which gradually expands the frontier of the reachable range along the directions to cover the entire originally desired goal space. To be specific, as shown in Fig. 1, we firstly evaluate achieved goal distribution using Kernel Density Estimation (KDE) [7] and select a set of goals in the low-density region of the distribution. Then we rank the selected goals by measuring the element-wise entropy computed by the estimated desired goal distribution. Exploring and exploiting these goals (called exploration goals) makes the agent push the frontier of achieved goals forward and simultaneously generalize better to the desired goal space. Moreover, we propose goal and transition augmentation to achieve higher data efficiency and better generalization.

To summarize, our main contributions of the paper are :

- We propose a novel density-based curriculum learning method to select hindsight goals for hindsight reinforcement learning, which helps guide the agent efficiently explore towards the desired goal distribution.
- We introduce two simple but effective data augmentation methods: goal augmentation and transition augmentation to help further improve the data efficiency and generality.
- We evaluate our method on several challenging variants of robotic manipulation tasks and demonstrate that our method achieves state-of-the-art performance compared to other baselines.

II. RELATED WORK

Our work focuses on the problem of how to design a curriculum to efficiently guide the agent to explore the informative space during reinforcement learning with sparse rewards. Therefore, in this section, we will review the key related works on these topics.

Corresponding Author: Xuguang Lan. Xuguang Lan is with the Institute of Artificial Intelligence and Robotics, Xi’an Jiaotong University, No.28 Xianning Road, Xi’an, Shaanxi, China. xglan@mail.xjtu.edu.cn

This work was supported in part by NSFC under grant No.6212500145, No.62088102, No.61973246, No.91748208, Shaanxi Project under grant No.2018ZDCXLYG0607, and the program of the Ministry of Education.

A. Exploration in Sparse-reward RL

Informative and effective explorations are essential for solving RL tasks. However, for sparse-reward RL, traditional exploration strategies, including ϵ -greed [8], Ornstein–Uhlenbeck noise [9], maximum entropy explore in action space [10], and noise adding in parameter space [11], [12] cannot work well since they cannot help the agent to receive enough positive rewards.

Recent works impel agents to discover novel states by intrinsic motivation. They can be mainly divided into three categories: count-based methods, curiosity-driven methods, and entropy-based methods. Count-based methods estimate visit counts and turn them into intrinsic bonuses to help agents reduce uncertainty [13]–[16]. Such methods are not suitable for tasks with large state space like robotic manipulation. The curiosity-driven method often utilizes prediction errors in the learned feature space to provide intrinsic reward [17]–[22]. However, these shaped rewards may cause the learning process detoured and unstable [23]. The entropy-based method aims to maximize information gain [24]–[27], state entropy [28]–[30], or mutual information between the agent’s action and the next states [20], [31] to encourage agents to discover novel states or explore unexplored regions of the environment.

In our work, we focus on setting curriculum-based goals step by step to guide the agent explore the environment efficiently. The most related work to ours is OMEGA [24], which selects achieved goals in sparsely explored areas to maximize entropy gain to push the frontier of achieved goals forward. However, OMEGA usually turns to unsafe goals (e.g. below the table) pursuit at the early exploration because of the instability of the mixture distribution, which may hinder the learning process or even cause irrevocable damages to the environment. By contrast, our work utilizes estimated desired goal distribution to constraint novel goal states and use goal augmentation to quickly expand the boundary of achieved goals, thus making the exploration more efficient and safe.

B. Curriculum Learning

Our work is also highly related to curriculum learning (CL) [32], which has been introduced to deep RL to guide data collection and exploitation [33]. During data collection, CL defines a set of auxiliary tasks to be used as stepping stones towards the main task. For example, CL could be used to generate a sequence of sub-goals [34]–[37] to guide the agent to finish the original task step by step. In multi-goal reinforcement learning, CL could be also applied to guide the goal selection [38]–[43]. Particularly, hindsight goal generation (HGG) [42] is a state-of-the-art method in the area of exploration algorithms in sparse-reward RL, which optimizes desired goals by solving the Wasserstein Barycenter problem. CL can also be used in data exploitation by performing transition prioritization and transition modification [44]–[46]. For example, CHER [46] adaptively selects the failed experiences for replay according to the proximity to original goals and the curiosity of exploration. However,

the metric to select goals is hand-designed and influenced by hyperparameters. Besides, the selection mechanism [47] of CHER is quite time-consuming. By contrast, our work utilizes a more generalizable KDE algorithm without increasing computational complexity.

III. BACKGROUND

A. Goal-conditioned Reinforcement Learning

Compared with the standard RL which is formulated by Markov Decision Process (MDP) represented as a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma, \rho_0)$, Goal-conditioned Reinforcement Learning (GCRL) arguments it with an additional set of goals \mathcal{G} . Formally, the Goal-conditioned MDP is denoted as a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{G}, \mathcal{P}, r_g, \gamma, \rho_0)$, where \mathcal{S} and \mathcal{A} are set of states and actions, respectively, $\mathcal{P} = p(s_{t+1}|s_t, a_t)$ defines the transition dynamics of the environment. $\mathcal{G} \in \mathbb{R}^{d_g}$ is the goal space and $r_g : \mathcal{S} \times \mathcal{A} \times \mathcal{G} \rightarrow \mathbb{R}$ defines the reward function that describes how close the agent is to the desired goals. $\gamma \in [0, 1]$ is a discount factor and $\rho_0 = p(s_0)$ is the initial state distribution. Specifically, in the sparse reward setting, the reward function has the following form:

$$r_g(s_t, a_t, g) = \begin{cases} 0, & \|\phi(s_t) - g\| \leq \delta \\ -1, & \text{otherwise} \end{cases}, \quad (1)$$

where $\phi : \mathcal{S} \rightarrow \mathcal{G}$ is a given function that maps state space to goal space. δ is a known distant threshold that represents the tolerance of reached goal. $\|\bullet\|$ is a distant metric.

B. Deep Deterministic Policy Gradient

Deep Deterministic Policy Gradient (DDPG) [9] is an off-policy model-free reinforcement learning algorithm, which utilizes an actor-critic [48] framework and is proposed to deal with continuous control tasks. The actor is for approximating policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ and the critic aims to learn action-value function $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$.

As for multi-goal continuous control tasks, DDPG can be extended with Universal Value Function Approximators (UVFA) [49], whose idea is to augment the Q-function and policy with multiple goal states $g \in \mathcal{G}$. Thus the actor and critic network additionally depend on goal state: $\pi = \pi(s, g)$, $Q = Q(s, a, g)$.

C. Hindsight Experience Replay

Hindsight experience replay (HER) [6] was proposed to solve the sparse reward environments in multi-goal RL. HER learns from failed experience using goal relabeling technique. Given a trajectory $\tau = \{(s_t, a_t, g, r_t, s_{t+1})\}_{t=1}^T$ of length T , HER alternates g and r in the i -th transition $(s_i, a_i, g, r_i, s_{i+1})$ with a future achieved goal in the same trajectory $g' = \phi(s_{i+k})$, $0 < k \leq T - i$ and $r' = r(s_i, a_i, g')$ computed by Eq. (1). Through relabeling, transitions from a failed trajectory can gain non-negative reward, thus helping the agent’s learning in sparse reward environment. HER can be combined with any off-policy algorithms such as DDPG [9], TD3 [50], SAC [10]. In this paper, we adopt DDPG+HER framework following [6].

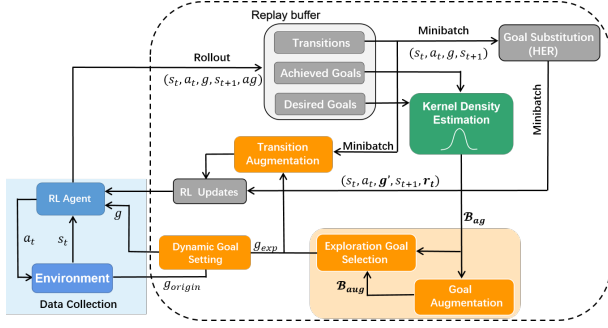


Fig. 2: Framework of our proposed method. The dash box contains the main modules of our method.

IV. METHOD

A. Overview

The overview of our proposed method is shown in Fig. 2. We mark the key components in orange and green colors. Roughly, we follow the standard off-policy RL skeleton, in which the past experiences from the interaction between the agent and the environment will be stored in the replay buffer and used in policy updates. The key contribution of our work lies in the dashed box defining the sampling and post-process of past experiences. To be specific, firstly, we store the achieved and desired goals into the corresponding replay buffer after trajectory data collection. Then we fit density models of achieved and desired distribution using the KDE module, and select exploration goals according to Section IV-B and IV-C. These exploration goals have two aspects of usefulness. One (Dynamic exploration goal setting module, see Section IV-D) is to replace the original goal provided by the environment with probability ϵ at the beginning of trajectory collection, the other (Transition augmentation module, see Section IV-E) is to relabel transition goals like HER substitution mechanism. Finally, we mix augmented transitions with relabeled ones that use HER *future* strategy for policy updates.

By setting exploration goals at the beginning of every training epoch, the agent is capable to get more guidance to explore along the reachable boundary. Moreover, by cooperating with goal and transition augmentation, the agent could achieve better generalization and performance.

B. Curriculum via Achieved Goal Density

We hope that the agent could set its own exploration goals to avoid pursuing unobtainable goals. But which goals are more valuable and should be selected for exploration? One notable fact in HER is that achieved goals easily get stuck in the area around the starting point during early exploration. However, learning these repeated goal states is useless to improve the agent’s capability, while goal states with low reachability are more informative to extend the reachable boundary. Therefore, we focus on the low-density region of achieved goals, which is similar to the curiosity learning.

To do so, during exploration, we collect achieved goals and store them in the replay buffer. Then we introduce the KDE algorithm to estimate their density, and thus obtain the

density model of the achieved goals M_{ag} . For an achieved goal ag_i in the replay buffer, its density can be computed as

$$\rho_i = M_{ag}(ag_i), \quad (2)$$

then the density is normalized using the following equation:

$$\hat{\rho}_i = \frac{\rho_i - \rho_{\min}}{\sum_{n=1}^N (\rho_n - \rho_{\min})}, \quad (3)$$

where ρ_{\min} represents the minimum value of density, and N is the size of the achieved goal buffer. After normalization, all density values are limited within the range $0 < \hat{\rho} < 1$. To ensure that low-density achieved goals have higher prioritization, we use $1 - \hat{\rho}_i$ and softmax function to calculate sample probability of ag_i and store the value in the goal replay buffer.

$$p(ag_i) = \frac{e^{1-\hat{\rho}_i}}{\sum_{n=1}^N e^{1-\hat{\rho}_n}} \quad (4)$$

Finally, we will sample a batch of achieved goals $\mathcal{B}_{ag} = \{g_i\}_{i=1}^n$ according to sample probability defined in Eq. (4).

C. Curriculum via Desired Goal Density

Although goals in \mathcal{B}_{ag} could be more informative to learn, not all these goals are safe or valuable for exploration. Actually, some goals are exactly opposite to the desired distribution or out of the safe workspace (e.g., below the table), which are useless to learn and may cause damages to the environment. To constrain the selected goals, we propose to utilize desired goal density model to further rank the selected goals. Formally, for each selected achieved goal g_i , its density probability in desired goal distribution can be calculated as follows:

$$p(g_i) = M_{dg}(g_i), \quad (5)$$

where M_{dg} is the desired goal density learned in a similar way to M_{ag} .

Goals with higher desired density will be naturally sampled more frequently and hence dominate the learning process. However, we found that such domination will harm the final performance since the agent tends to overfit such high-density goals. To achieve better performance, the agent should not only learn to reach high-density goals, but also be capable of relatively low-density goals. Therefore, learning to reach goals that lie in the uncertain areas of the desired distribution is valuable to improve the final performance. Here, we choose element-wise entropy to measure the uncertainty of g_i in desired goal distribution following $e(g_i) = -p(g_i) \log p(g_i)$, and the normalized entropy is

$$\hat{e}(g_i) = \frac{e(g_i)}{\sum_{j=1}^n e(g_j)}. \quad (6)$$

We rank the selected achieved goals \mathcal{B}_{ag} by the normalized entropy and then only preserve the top-ranking ones to form the training batch. Intuitively, such a mechanism will impose high priority on goals at the boundary of the desired distribution, and ignore the ones with extremely low density, which are not valuable to learn, and high density, which

have already been enough learned. It helps to achieve better generalization to the whole distribution instead of only the high-density area.

D. Dynamic Exploration Goal Setting

Our method is proposed to help the agent efficiently explore the environment especially at the beginning stage. However, such goal relabeling techniques will inevitably introduce learning bias, i.e., it encourages the agent to finish the relabeled tasks, and sometimes prevents the agent from learning to finish the original tasks. Therefore, we hope that the agent could explore in high efficiency with our method at the beginning, and gradually turn to learn the original tasks after enough exploration.

To achieve this, at the beginning of each episodes, we replace the original goal provided by the environment with the selected goals with probability ϵ , which is defined as a function that is related to test success rate p_s :

$$\epsilon = \alpha e^{-2p_s}, \quad (7)$$

where $\alpha \in (0, 1)$ is a hyperparameter. Intuitively, when the test success rate is low, reaching the original goal is beyond the agent's ability, so we encourage agents to explore with the curriculum-based exploration goals. As the test success rate is increasing, the agent learns to reach more and more original goals. Under this circumstance, we gradually reduce the portion of the relabeled data to alleviate the learning bias.

E. Goal and Transition Augmentation

We use augmentation techniques to further expand the boundary of achieved goals and increase exploration efficiency.

In sparse reward environments with multiple goals, a transition gets a non-negative reward when the achieved goal is within a small radial threshold (a ball) around the desired goal. As Fig. 1 shows, goals that within the circle with radius Δ (red dot) also satisfy the success condition for hindsight replacement. These goals form a set $\mathcal{B}_{near}^{g'} = \{g \in \mathcal{G} | d(g, g') < \Delta, \Delta \leq \delta\}$, where g' is a hindsight goal, $d(\bullet)$ is the distance metric used in reward function (Eq. (1)) and δ is the maximum tolerance to get a non-negative reward. In our practice, after we select a set of goals \mathcal{B}_{ag} , we further randomly sample a batch of goals in \mathcal{B}_{ag} . For each sampled goal g' , we randomly generate an additional goal from its corresponding $\mathcal{B}_{near}^{g'}$, thus forming a set of augmented goals \mathcal{B}_{aug} . Then we utilize the set $\mathcal{B}_{ag} \cup \mathcal{B}_{aug}$ to further sample exploration goals following Eq. (6). The augmented goals could improve the robustness against the small perturbations.

For goal-conditioned RL, regardless of the goal being pursued, (s, a, s') transitions are unbiased samples from the environment dynamics. As a result, an agent is free to pair transitions with any goal and corresponding reward. To provide more instructive samples, we propose to utilize the selected goals $\mathcal{B}_{ag} \cup \mathcal{B}_{aug}$ to augment sampled transitions. The transition augmentation technique has two significant advantages: 1) increasing the diversity of learning goals at

early exploration; and 2) learning more information of better generalization to desired goals.

Accompanied by the above two augmentation techniques, the agent is capable to achieve more efficient exploration and better generalization.

V. EXPERIMENTS AND RESULTS

In this section, we first introduce the simulation environments of robotic manipulation tasks. Then we present performance results of our method and other baselines as well as analyzing the results. Finally, we conduct ablation experiments to confirm the effectiveness of our method.

A. Simulation Environments

The experimental environments are all goal-conditioned tasks with sparse reward settings and simulated with MuJoCo physics engine [51]. The robot we use in experiments is a 7-DOF (degrees of freedom) Fetch robotic arm with a two-fingered parallel gripper. The action space is 4-dimensional, with 3-dimensional end-effector velocities in x, y, z direction, and 1-dimensional grasping indicator. The goal space is 3-dimensional indicating the Cartesian coordinate in working space. The state is 25-dimensional for all environments. We evaluate our method on several variants of the standard benchmarks [52] proposed by Open AI gym [53].

- **Sliding** (Fig. 4a) The robot arm is required to contact the object with enough force and slide the object to a goal position on the table.
- **Throwing** (Fig. 4b) Similar to the *Sliding* task, but the robot arm needs to throw a rubber ball to a goal position outside the robot's workspace. The rubber ball is initialized between the two-fingered parallel gripper.
- **Pushing** (Fig. 4c) The robot arm aims to push an object to a position which is far away from initial position. We add two brick-like obstacles to separate the target position and initial position of the object based on the standard benchmark task *FetchPush-v1*.
- **Pick and Place** We extend the *FetchPickAndPlace-v1* task by increasing the height of desired goals and adding an obstacle. In *FetchPnP-InAir* (Fig. 4d) task, the heights of desired goals are between $0.2m$ and $0.45m$ above the table, which is the same setting with [24]. In *FetchPnP-Obstacle* task (Fig. 4e), the obstacle is added to separate the object initial position and desired goal distribution.

B. Training Setting

We evaluate our method with five baseline algorithms: HER [6], CHER [46], MEP [44], HGG [42] and OMEGA [24]. Moreover, to verify the effectiveness of learning uncertain areas of desired goals, we compare our method that uses Eq. (6) to select top-ranking goals following entropy prioritization (called Ours_E) with the one that uses Eq. (5) (called Ours_D). All the algorithms are combined with DDPG [9] framework. We implement HER, CHER, MEP, HGG with minor modifications based on the official

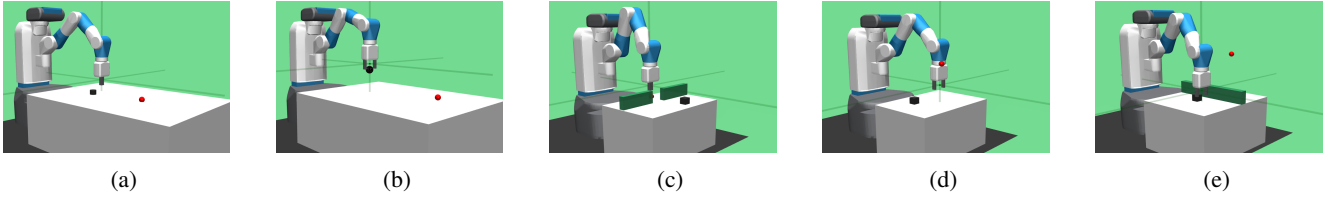


Fig. 3: Five challenging manipulation tasks with sparse reward settings. (a) *FetchSlide-v1*; (b) *FetchThrowBall-v1*; (c) *FetchPush-Obstacle*; (d) *FetchPnP-InAir*; (e) *FetchPnP-Obstacle*.

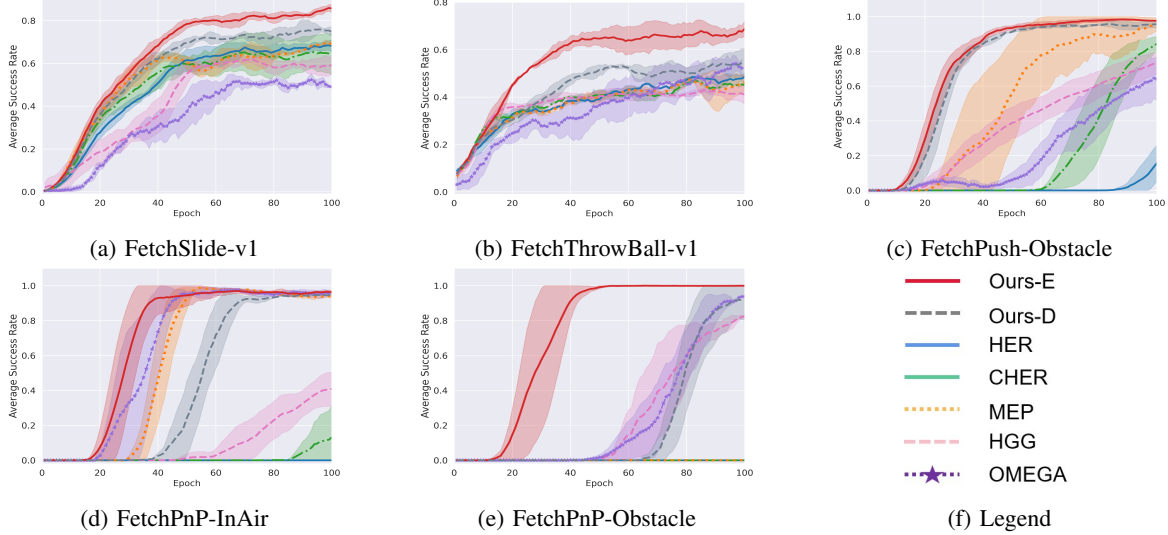


Fig. 4: Learning curves of our method compared with baselines on five challenging robotic manipulation tasks. The solid line indicates the average success rate and the shaded area is the standard error over 5 different random seeds.

codebase and unify these baseline algorithms under the same framework. For OMEGA, we use the author’s code.

In all experiments, we use 16 CPU cores to train the agent and each core generates experiences by using 2 parallel rollouts with MPI for synchronization. For OMEGA, we create 16 independent environments and keep the same timesteps and training frequency as other algorithms. We use the same set of training hyperparameters for HER and our method across all environments except for some specific variables. During training, each epoch generates $16 \times 2 \times 50 = 1600$ episodes. After each epoch, the test success rate is computed over $16 \times 10 = 160$ episodes. We observed that the performance of the proposed approach is robust by choosing explore ratio $\alpha = 0.2$ for *FetchSlide-v1* and $\alpha = 0.5$ for the other four environments. The size of the selected goal set \mathcal{B}_{ag} and augmented goal set \mathcal{B}_{aug} are 256 and 128 respectively. As for exploration goals, we keep the same size as the number of episodes in each epoch and select the top 100 from the set $\mathcal{B}_{ag} \cup \mathcal{B}_{aug}$ according to the entropy prioritization. Besides, we choose radius $\Delta = 3cm$ for goal augmentation and keep the augmented transition batch size 128 across all environments. The desired and achieved goal buffer size is 1000 and the transition buffer size is 10^6 .

C. Performance and Analysis

As shown in Fig. 4, our method (Ours-E) achieves state-of-the-art performances in all environments. It can be seen that

TABLE I: Average training time (hour) for different methods

	Ours	HER	CHER	MEP	HGG	OMEGA
Time	3.5h	3.25h	13.2h	5.2h	4.5h	8.3h

learning the goals that lie in uncertain areas of desired goal distribution contributes to better performances in all environments compared to choosing high-density goals (Ours-D).

Besides, we notice that OMEGA has lower sample efficiency than other baselines in *FetchSlide-v1* and *FetchThrowBall-v1*, indicating that maximizes the entropy of achieved goals may hinder the learning process in those environments that easy to generate unsafe goals (e.g., below the table). However, it’s undeniable that entropy-based methods (OMEGA and MEP) contribute to the faster exploration and better generalization than HER. Our method shows significant advantage in exploration efficiency especially in the most difficult task *FetchPnP-Obstacle*, where HER, CHER, and MEP are failed to solve. Comparing with HGG and OMEGA, our approach achieves over $3 \times$ speedup in exploration efficiency. As for training time (see Table I), our method only takes 0.25h longer than HER and has much lower computational complexity than CHER and OMEGA.

To gain a more intuitive understanding of the learned behavior across training, we visualize the achieved goals and exploration goals after the 5th, 15th, and 25th epoch (see

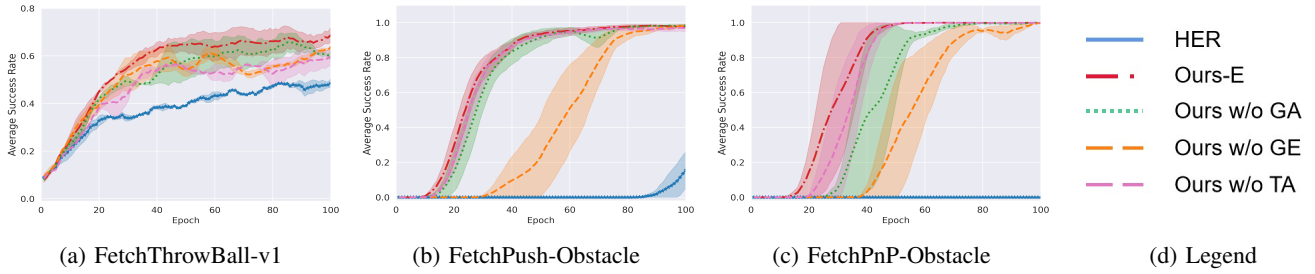


Fig. 5: Ablation study on different components of our method. Goal exploration (GE), goal augmentation (GA), and transition augmentation (TA). In Fig. 5d, w/o represents without.

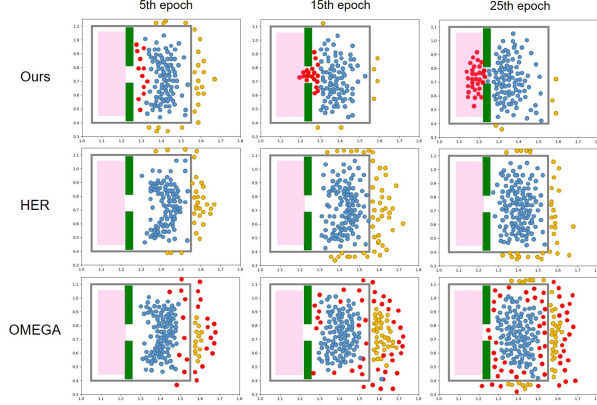


Fig. 6: Visualization of hindsight goals on *FetchPush-Obstacle*. The purple area is the region of desired goals. The gray rectangular frame represents the table workspace. The blue dot marks hindsight goal, the red dot marks the exploration goal of our method and OMEGA, and the orange dot represents the goal below the table. HER does not have red dots because of no exploration goal setting.

Fig. 6) during *FetchPush-Obstacle* task training. At the early stage of training, the exploration goals are mainly distributed on the boundary of the target distribution. Exploring those goals enables agents approach real targets faster and easier. However, HER always provides the agent completely unreachable goals that are out of the agent’s ability and learns a lot of invalid goals which fall off the table. OMEGA achieves more diverse goals than HER but plenty of goals are in the opposite direction to the desired goals, thus distracting agents from the original purpose. Our method encourages the agent to explore the boundary of its ability while getting closer to the real target. Therefore, our method makes agent learning in a curriculum learning manner.

D. Ablation Experiments

To measure the effect of each component in our method, we conduct ablation studies on three challenging environments. From Fig. 5 we can observe that the learning performance suffers a varying degree of decline without either component of the proposed method, but still achieves better than HER. It can be seen that transition augmentation is most important for *FetchThrowBall-v1*, indicating that learning goals within the uncertain area of desired distribution can further improve the performance. While in the rest of the hard exploration environments, exploration efficiency is severely

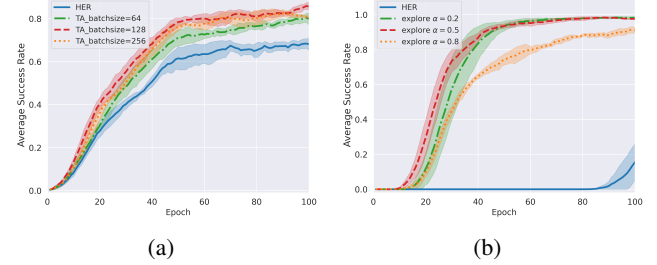


Fig. 7: Parameters study. (a) Different batch size of transition augmentation on *FetchSlide-v1*; (b) Different explore ratio α on *FetchPush-Obstacle*.

reduced without setting exploration goals. In addition, goal augmentation brings up an obvious advance in exploration when cooperating with the other two techniques. The synergy between the three components significantly contributes to better exploration efficiency and converged performance.

We also conduct parameter studies (see Fig. 7) on the batch size of augmented transitions and explore ratio α . It can be seen from Fig. 7a that less augmented transitions hinder the performance slightly in *FetchSlide-v1*, but still much better than HER. As shown in Fig. 7b, explore ratio α influences the exploration efficiency and performance slightly. A larger explore ratio ($\alpha = 0.8$) distracts the agent’s attention from learning the original task and decreases the performance. A proper explore ratio ($\alpha = 0.5$) balances the learning of local information and original tasks, thus achieving better exploration and performance.

VI. CONCLUSION

In this work, we propose a density-based exploration goal selection mechanism to guide goal-conditioned agents to explore efficiently towards the frontier of desired distribution. Cooperating with two simple but effective techniques: goal augmentation and transition augmentation, the agent is capable to generalize hindsight goal distribution better. We evaluate our method in five challenging sparse reward robotic environments and demonstrate substantial improvements in performance and sample efficiency compared to HER and other standard RL algorithms. Through further analysis, we demonstrate that each component of our method contributes a positive effect on exploration efficiency. Future research could concentrate on extension (e.g., image-based tasks), improvement, and real-world deployment of our method.

REFERENCES

- [1] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, *et al.*, “Human-level control through deep reinforcement learning,” *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [2] N. Justesen, P. Bontrager, J. Togelius, and S. Risi, “Deep learning for video game playing,” *IEEE Transactions on Games*, vol. 12, no. 1, pp. 1–20, 2019.
- [3] S. Levine, C. Finn, T. Darrell, and P. Abbeel, “End-to-end training of deep visuomotor policies,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 1334–1373, 2016.
- [4] M. Vecerik, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. Heess, T. Rothörl, T. Lampe, and M. Riedmiller, “Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards,” *arXiv preprint arXiv:1707.08817*, 2017.
- [5] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine, “Learning complex dexterous manipulation with deep reinforcement learning and demonstrations,” *arXiv preprint arXiv:1709.10087*, 2017.
- [6] M. Andrychowicz, D. Crow, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba, “Hindsight experience replay,” in *NIPS*, 2017.
- [7] M. Rosenblatt, “Remarks on some nonparametric estimates of a density function,” *Annals of Mathematical Statistics*, vol. 27, pp. 832–837, 1956.
- [8] C. Watkins and P. Dayan, “Q-learning,” *Machine Learning*, vol. 8, pp. 279–292, 2004.
- [9] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. A. Riedmiller, “Deterministic policy gradient algorithms,” in *ICML*, 2014.
- [10] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *ICML*, 2018.
- [11] M. Plappert, R. Houthoofd, P. Dhariwal, S. Sidor, R. Y. Chen, X. Chen, T. Asfour, P. Abbeel, and M. Andrychowicz, “Parameter space noise for exploration,” *ArXiv*, vol. abs/1706.01905, 2018.
- [12] M. Fortunato, M. G. Azar, B. Piot, J. Menick, I. Osband, A. Graves, V. Mnih, R. Munos, D. Hassabis, O. Pietquin, C. Blundell, and S. Legg, “Noisy networks for exploration,” *ArXiv*, vol. abs/1706.10295, 2018.
- [13] A. L. Strehl and M. L. Littman, “An analysis of model-based interval estimation for markov decision processes,” *Journal of Computer and System Sciences*, vol. 74, no. 8, pp. 1309–1331, 2008.
- [14] A. Strehl and M. Littman, “A theoretical analysis of model-based interval estimation,” *Proceedings of the 22nd international conference on Machine learning*, 2005.
- [15] G. Ostrovski, M. G. Bellemare, A. Oord, and R. Munos, “Count-based exploration with neural density models,” in *International conference on machine learning*, pp. 2721–2730, PMLR, 2017.
- [16] H. Tang, R. Houthoofd, D. Foote, A. Stooke, X. Chen, Y. Duan, J. Schulman, F. De Turck, and P. Abbeel, “# exploration: A study of count-based exploration for deep reinforcement learning,” in *31st Conference on Neural Information Processing Systems (NIPS)*, vol. 30, pp. 1–18, 2017.
- [17] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, “Curiosity-driven exploration by self-supervised prediction,” in *International conference on machine learning*, pp. 2778–2787, PMLR, 2017.
- [18] B. C. Stadie, S. Levine, and P. Abbeel, “Incentivizing exploration in reinforcement learning with deep predictive models,” *arXiv preprint arXiv:1507.00814*, 2015.
- [19] Y. Burda, H. Edwards, A. Storkey, and O. Klimov, “Exploration by random network distillation,” *arXiv preprint arXiv:1810.12894*, 2018.
- [20] H. Kim, J. Kim, Y. Jeong, S. Levine, and H. O. Song, “Emi: Exploration with mutual information,” *arXiv preprint arXiv:1810.01176*, 2018.
- [21] T. Zhang, P. Rashidinejad, J. Jiao, Y. Tian, J. Gonzalez, and S. Russell, “Made: Exploration via maximizing deviation from explored regions,” *arXiv preprint arXiv:2106.10268*, 2021.
- [22] J. Fu, J. D. Co-Reyes, and S. Levine, “Ex2: Exploration with exemplar models for deep reinforcement learning,” *arXiv preprint arXiv:1703.01260*, 2017.
- [23] J. Jiang and Z. Lu, “Generative exploration and exploitation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 4337–4344, 2020.
- [24] S. Pitis, H. Chan, S. Zhao, B. C. Stadie, and J. Ba, “Maximum entropy gain exploration for long horizon multi-goal reinforcement learning,” in *ICML*, 2020.
- [25] R. Houthoofd, X. Chen, Y. Duan, J. Schulman, F. De Turck, and P. Abbeel, “Vime: Variational information maximizing exploration,” *arXiv preprint arXiv:1605.09674*, 2016.
- [26] J. Achiam and S. Sastry, “Surprise-based intrinsic motivation for deep reinforcement learning,” *arXiv preprint arXiv:1703.01732*, 2017.
- [27] P. Shyam, W. Jaśkowski, and F. Gomez, “Model-based active exploration,” in *International conference on machine learning*, pp. 5779–5788, PMLR, 2019.
- [28] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey, *et al.*, “Maximum entropy inverse reinforcement learning,” in *Aaai*, vol. 8, pp. 1433–1438, Chicago, IL, USA, 2008.
- [29] L. Lee, B. Eysenbach, E. Parisotto, E. Xing, S. Levine, and R. Salakhutdinov, “Efficient exploration via state marginal matching,” *arXiv preprint arXiv:1906.05274*, 2019.
- [30] Y. Seo, L. Chen, J. Shin, H. Lee, P. Abbeel, and K. Lee, “State entropy maximization with random encoders for efficient exploration,” in *ICML*, 2021.
- [31] R. Zhao, Y. Gao, P. Abbeel, V. Tresp, and W. Xu, “Mutual information state intrinsic control,” *arXiv preprint arXiv:2103.08107*, 2021.
- [32] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, “Curriculum learning,” in *Proceedings of the 26th annual international conference on machine learning*, pp. 41–48, 2009.
- [33] R. Portelas, C. Colas, L. Weng, K. Hofmann, and P.-Y. Oudeyer, “Automatic curriculum learning for deep rl: A short survey,” *arXiv preprint arXiv:2003.04664*, 2020.
- [34] Y. Lai, W. Wang, Y. Yang, J. Zhu, and M. Kuang, “Hindsight planner,” in *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 690–698, 2020.
- [35] E. Chane-Sane, C. Schmid, and I. Laptev, “Goal-conditioned reinforcement learning with imagined subgoals,” in *ICML*, 2021.
- [36] O. Nachum, S. Gu, H. Lee, and S. Levine, “Data-efficient hierarchical reinforcement learning,” in *NeurIPS*, 2018.
- [37] O. Kilinc and G. Montana, “Follow the object: Curriculum learning for manipulation tasks with imagined goals,” *ArXiv*, vol. abs/2008.02066, 2020.
- [38] C. Florensa, D. Held, X. Geng, and P. Abbeel, “Automatic goal generation for reinforcement learning agents,” in *International conference on machine learning*, pp. 1515–1528, PMLR, 2018.
- [39] H. Charlesworth and G. Montana, “Plangan: Model-based planning with sparse rewards and multiple goals,” *arXiv preprint arXiv:2006.00900*, 2020.
- [40] M. Tomar, A. Sathuluri, and B. Ravindran, “Mamic: macro and micro curriculum for robotic reinforcement learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 10053–10054, 2019.
- [41] V. H. Pong, M. Dalal, S. Lin, A. Nair, S. Bahl, and S. Levine, “Skew-fit: State-covering self-supervised reinforcement learning,” *ArXiv*, vol. abs/1903.03698, 2020.
- [42] Z. Ren, K. Dong, Y. Zhou, Q. Liu, and J. Peng, “Exploration via hindsight goal generation,” in *NeurIPS*, 2019.
- [43] Y. Zhang, P. Abbeel, and L. Pinto, “Automatic curriculum learning through value disagreement,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [44] R. Zhao, X. Sun, and V. Tresp, “Maximum entropy-regularized multi-goal reinforcement learning,” in *ICML*, 2019.
- [45] R. Zhao and V. Tresp, “Energy-based hindsight experience prioritization,” in *Conference on Robot Learning*, pp. 113–122, PMLR, 2018.
- [46] M. Fang, T. Zhou, Y. Du, L. Han, and Z. Zhang, “Curriculum-guided hindsight experience replay,” in *NeurIPS*, 2019.
- [47] B. Mirzasoleiman, A. Badanidiyuru, A. Karbasi, J. Vondrák, and A. Krause, “Lazier than lazy greedy,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 29, 2015.
- [48] R. Sutton and A. Barto, “Reinforcement learning: An introduction,” *IEEE Transactions on Neural Networks*, vol. 16, pp. 285–286, 2005.
- [49] T. Schaul, D. Horgan, K. Gregor, and D. Silver, “Universal value function approximators,” in *ICML*, 2015.
- [50] S. Fujimoto, H. Hoof, and D. Meger, “Addressing function approximation error in actor-critic methods,” in *International Conference on Machine Learning*, pp. 1587–1596, PMLR, 2018.
- [51] E. Todorov, T. Erez, and Y. Tassa, “Mujoco: A physics engine for model-based control,” in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033, IEEE, 2012.

- [52] M. Plappert, M. Andrychowicz, A. Ray, B. McGrew, B. Baker, G. Powell, J. Schneider, J. Tobin, M. Chociej, P. Welinder, V. Kumar, and W. Zaremba, “Multi-goal reinforcement learning: Challenging robotics environments and request for research,” *ArXiv*, vol. abs/1802.09464, 2018.
- [53] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, “Openai gym,” *arXiv preprint arXiv:1606.01540*, 2016.