# Goal-Conditioned On-Policy Reinforcement Learning

Xudong Gong[1,2]  Dawei Feng[1,2]  Kele Xu[1,2]  Bo Ding[1,2]  Huaimin Wang[1,2]

[1]College of Computer, National University of Defense Technology, Changsha, Hunan, China

[2]State Key Laboratory of Complex & Critical Software Environment, Changsha, Hunan, China

## Multi-Goal Problems

➢ Controlling robotic arms to grasp objects at any location on a table

➢ Operating fixed-wing UAVs to navigate towards any specified velocity vector

➢ ......

## Goal-Conditioned Reinforcement Learning

Learns goal-conditioned behaviors that can achieve and generalize across a range of different goals
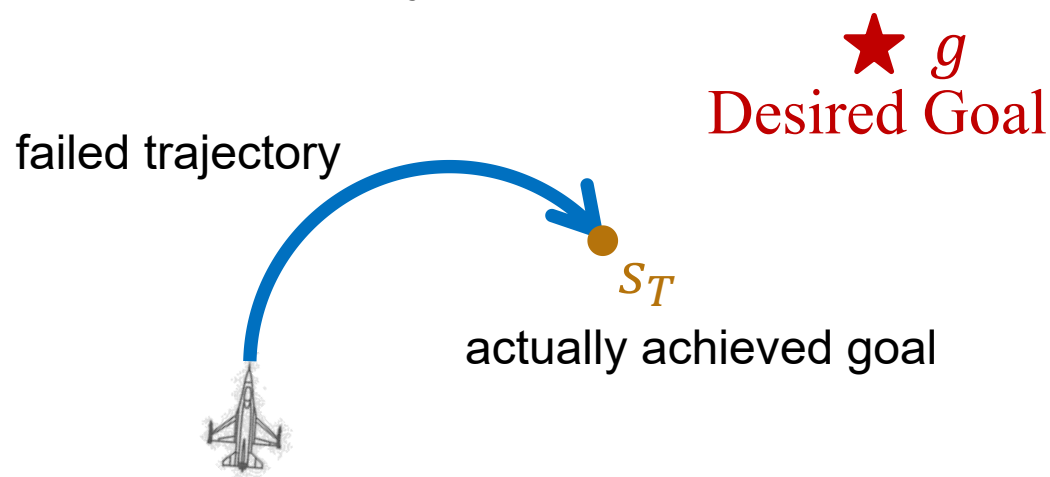
## Challenge

➢ The additional goal space intensifies the complexity of exploration

➢ The goals used in sampling data and the order of these goals affects GCRL's training efficiency and effectiveness

## Mainstream methods

Based on HER[1] (Hindsight Experience Replay)

1. Relabel failed trajectories[1]



★ $g$
Desired Goal

failed trajectory

$s_T$

actually achieved goal

Relabel desired goal and recompute reward.

| step | failed traj | relabeled traj |
|------|-------------|----------------|
| 0 | $(s_0, a_0, r(s_0, g))$ | $(s_0, a_0, r(s_0, s_N))$ |
| 1 | $(s_1, a_1, r(s_1, g))$ | $(s_1, a_1, r(s_1, s_N))$ |
| ⋮ | ⋮ | ⋮ |
| N | $(s_N, a_N, r(s_N, g))$ | $(s_N, a_N, r(s_N, s_N))$ |

➢ Failed trajectories contribute almost nothing to policy optimization.
➢ If the failed but actually achieved goal state $s_T$ is considered as the desired goal, then this failed trajectory becomes a successful one, which can help optimize the policy.
➢ Besides the last state, any state in the trajectory can be mapped to a potential hindsight goal.

**Assigning valuable use to failure experiences, alleviating the exploration challenge!**

1. Andrychowicz M, Wolski F, Ray A, et al. Hindsight experience replay[J]. Advances in neural information processing systems, 2017, 30.

## Mainstream methods

Based on HER (Hindsight Experience Replay)

2. Arrange behavioral goals[1]

$\ldots$   ★$g_N$

★$g_2$

★$g_1$

Learning to achieve desired goals in some orders

Initial state

➤ Evaluate policy's ability of achieving goals based on the data in the ER, and sample behavioral goals of appropriate difficulty for sampling training data

**Further improve training efficiency by arranging behavioral goals!**

1. Pitis S, Chan H, Zhao S, et al. Maximum entropy gain exploration for long horizon multi-goal reinforcement learning[C]//International Conference on Machine Learning. PMLR, 2020: 7750-7761.

## Mainstream methods

Based on HER (Hindsight Experience Replay)

➤ Relabel failure experience $\rightarrow$ alleviate the exploration challenge

➤ Evaluate policy's ability of achieving goals based on the data in the ER, and select behavioral goals of appropriate difficulty for sampling training data $\rightarrow$ further improve training efficiency by arranging behavioral goals

**Algorithm 1** Unified Framework for Multi-goal Agents

**function** TRAIN(*$*args$*):
  Alternate between collecting experience using ROLLOUT and optimizing the parameters using OPTIMIZE.

**function** ROLLOUT (policy $\pi_{explore}$, buffer $\mathcal{B}$, *$*args$*):
  $g \leftarrow$ SELECT($*args$)
  $s_0 \leftarrow$ initial state
  **for** $t$ in $0 \dots T-1$ **do**
    $a_t, s_{t+1} \leftarrow$ execute $\pi_{explore}(s_t, g)$ in environment
    $r_t \leftarrow$ REWARD($s_t, a_t, s_{t+1}, g$)
    Store ($s_t, a_t, s_{t+1}, r_t, g$) in replay buffer $\mathcal{B}$

**function** OPTIMIZE (buffer $\mathcal{B}$, algorithm $\mathcal{A}$, parameters $\theta$):
  Sample mini-batch $B = \{(s, a, s', r, g)_i\}_{i=1}^N \sim \mathcal{B}$
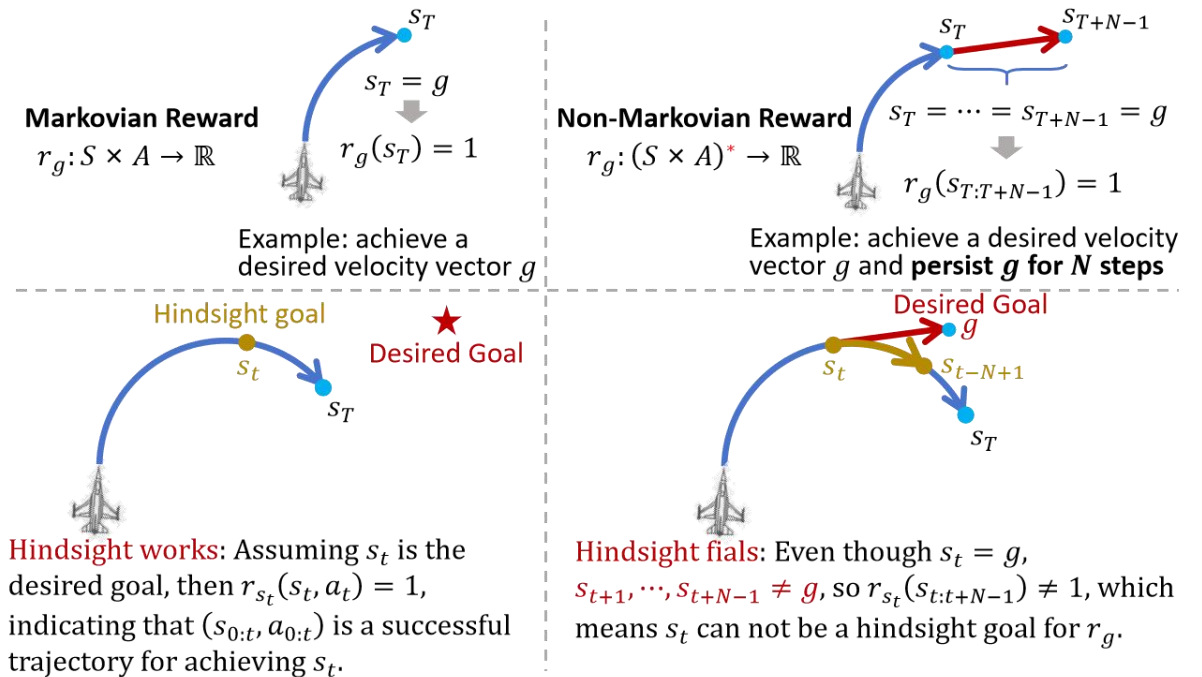  $B' \leftarrow$ RELABEL($B, *args$)
  Optimize $\theta$ using $\mathcal{A}$ (e.g., DDPG) and relabeled $B'$

A general GCRL training framework[1]

1. Pitis S, Chan H, Zhao S, et al. Maximum entropy gain exploration for long horizon multi-goal reinforcement learning[C]//International Conference on Machine Learning. PMLR, 2020: 7750-7761.

## Challenge of HER-based methods

➤ HER-based methods make an implicit assumption that the rewards **depends only on the current state** (Markovian Reward, MR). Under this assumption, any single state can potentially become a hindsight goal.

➤ However, when the computation of rewards **depends on multi-step states** (Non-Markovian Reward, NMR), a failed trajectory may not contain any state sequence that satisfies the reward.



**Markovian Reward**
$r_g: S \times A \to \mathbb{R}$

$s_T$
$s_T = g$
$r_g(s_T) = 1$

Example: achieve a desired velocity vector $g$

Hindsight goal
$s_t$
$s_T$
★ Desired Goal

Hindsight works: Assuming $s_t$ is the desired goal, then $r_{s_t}(s_t, a_t) = 1$, indicating that $(s_{0:t}, a_{0:t})$ is a successful trajectory for achieving $s_t$.

**Non-Markovian Reward**
$r_g: (S \times A)^* \to \mathbb{R}$

$s_T$ $s_{T+N-1}$
$s_T = \cdots = s_{T+N-1} = g$
$r_g(s_{T:T+N-1}) = 1$

Example: achieve a desired velocity vector $g$ and **persist $g$ for $N$ steps**

Desired Goal $g$
$s_t$ $s_{t-N+1}$
$s_T$

Hindsight fials: Even though $s_t = g$, $s_{t+1}, \cdots, s_{t+N-1} \neq g$, so $r_{s_t}(s_{t:t+N-1}) \neq 1$, which means $s_t$ can not be a hindsight goal for $r_g$.

**Can a GCRL framework be proposed that does not rely on HER and can simultaneously address both MR and NMR problems?**

# Goal-conditioned on-policy reinforcement learning (GCPO)

Insight: refer to the two successful designs of HER-based methods

**HER-based methods**

**GCPO**

Relabel failure experience

← alleviate the exploration challenge →

Pre-training from demonstrations to provide a behavioral prior for the policy

(1) Evaluate policy's goal-achieving ability **with the help of ER**
(2) select appropriate behavioral goals
(3) optimize policy with off-policy RL

← further improve training efficiency by arranging behavioral goals →

(1) evaluate policy periodically
(2) estimate policy's goal achieving ability with **Off-Policy Evaluation (OPE)** method
(3) select appropriate behavioral goals
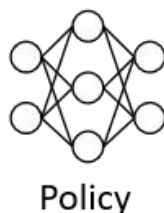(4) optimize polify with KL-Regularized on-policy RL

# General framework



**1.Pre-Training from Demonstrations**
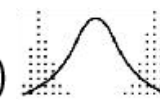
pre-train

Policy

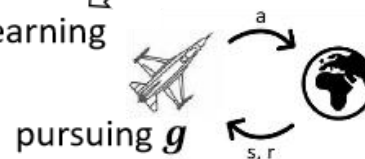Demonstrations

**2.Online Self-Curriculum Learning**

① estimate $p_{ag}$

② sample $g \sim f(p_{ag})$

③ online learning

pursuing $g$

1.Pre-training from demonstrations

pre-training provides the policy with an initial ability to achieve some of the desired goals, enhancing informative rewards during online learning.

2.Online self-curriculum learning

design an online self-curriculum learning mechanism that autonomously constructs a curriculum, generating behavioral goals that are incrementally more difficult than those the policy is currently capable of achieving.

# A practical implementation

**1.Pre-training from demonstrations**

utilize Behavioral Cloning (BC) to pre-train policies

$$\mathcal{L}(\theta) = -\mathbb{E}_{(s,a)\sim\mathcal{D}_E}\left[\log \pi_\theta(a|s)\right]$$

**2.Online self-curriculum learning**

*2.1 estimating the current policy's goal-achieving ability, $p_{ag}$*

employ GMM to estimate $p_{ag}$ with historical evaluation data

*2.2 sampling progressively challenging behavioral goals*

utilize inverse probability weighting, $[f_{MEGA}(p_{ag}, p_{dg})](g) = \dfrac{\frac{1}{p_{ag}(g)}}{\sum_{p'}\frac{1}{p_{ag}(g')}}$

*2.3 Conducting online RL learning with behavioral goals*

optimize policy with KL-regularized RL, $J_{kl}(\pi_\theta) = \mathbb{E}\left[\sum_t \gamma^t\left(r - \lambda log(\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_0}(a_t|s_t)})\right)\right]$

---

**Algorithm 1** Goal-Conditioned Policy Optimization (GCPO) framework

**Require:** demonstrations $\mathcal{D}_E$, distribution of desired goal $p_{dg}$, goal weight discount factor $\kappa$, online evaluation budget $N$, probability transform function $f$
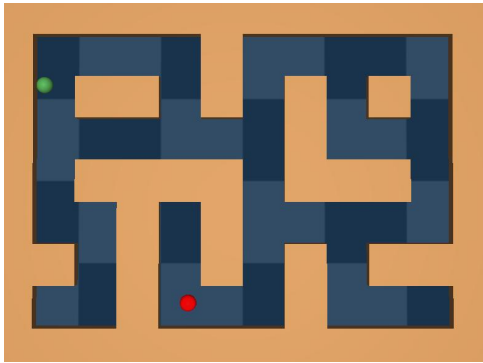
**Ensure:** $\pi_\theta(\cdot|s,g)$

1: Initialize goal-conditioned policy $\pi_\theta(\cdot|s,g)$, goal buffer $B_g$ which stores tuples of achieved goals in evaluation and their corresponding weight $(g, w_g)$
2: pre-train $\pi_\theta(\cdot|s,g)$ by Eq. 1       ▷ pre-train policy
3: **while** Not converge **do**
4:     **for all** $(g, w_g)$ in $B_g$ **do**
5:       $w_g \leftarrow \kappa \cdot w_g$       ▷ decay weight of historically achieved goals
6:     **end for**
7:     sample $N$ goals, $g_1, g_2, \ldots, g_N$ uniformly from $p_{dg}$       ▷ online policy evaluation
8:     **for all** $g$ in $g_1, g_2, \ldots, g_N$ **do**
9:       **if** $\pi_\theta$ finishes $g$ successfully **then**
10:         add $(g, 1.0)$ to $B_g$
11:       **end if**
12:     **end for**
13:     estimate $p_{ag}$ with GMM on $B_g$       ▷ estimate $p_{ag}$
14:     $\mathcal{D} \leftarrow \varnothing$       ▷ roll-out samples
15:     **while** Not collect enough online samples **do**
16:       sample a goal $g$ from Eq. 2
17:       sample a trajectory $\tau$ by $\pi_\theta$ on $g$
18:       append $\tau$ to $\mathcal{D}$
19:     **end while**
20:     update $\pi_\theta$ by Eq. 3 on $\mathcal{D}$       ▷ update policy
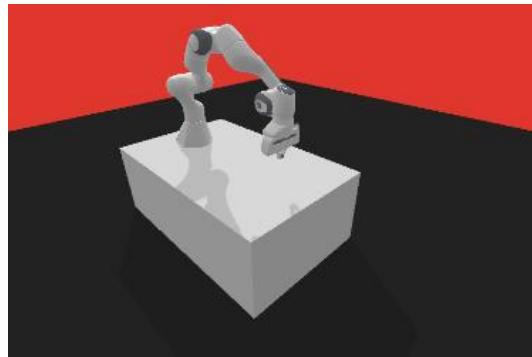21: **end while**

# Outline

## Settings

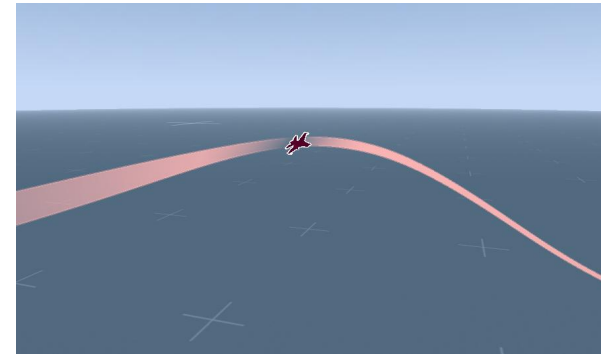➤ Tasks



PointMaze



Reach



Fixed-wing velocity vector control

## Main Results

On PointMaze and Reach

Table 7: Comparison between GCPO and baselines on Reach and PointMaze. The mean and variance of % success rates are presented over 5 random seeds. Optimal values are highlighted in bold, and sub-optimal values are underlined.

| Task | Reward | SAC + HER + MEGA | BC | GCPO |
|---|---|---|---|---|
| Reach | MR | **100.0±0.0** | 70.63±2.99 | **100.0±0.0** |
| | NMR | 0.72±1.34 | 10.52±11.70 | **80.26±17.01** |
| PointMaze | MR | **100.0±0.0** | 75.96±5.34 | 93.33±3.06 |
| | NMR | 4.17±0.93 | 22.8±3.71 | **47.50±8.06** |

➢ under MR settings, GCPO exhibits similar performance to HER-based methods.
➢ under NMR settings, GCPO shows significantly superior performance than HER-based methods.

# Main Results

On fixed-wing velocity vector control

Table 2: Comparison between GCPO and baselines on NMR. The mean and variance of % success rates are presented over 5 random seeds. Optimal values are highlighted in bold, and sub-optimal values are underlined.

| | Demonstration | | SAC + HER + MEGA | BC | GCPO w/o pre-training | GCPO w/o self-curriculum | GCPO |
|---|---|---|---|---|---|---|---|
| notation | #traj | traj length | | | | | |
| $\mathcal{D}_E^0$ | 10264 | 281.83±149.48 | | 17.08±0.57 | | 31.28±8.97 | **45.87±3.09** |
| $\mathcal{D}_E^1$ | 27021 | 119.64±47.55 | | 36.54±1.97 | | 43.49±3.85 | **49.12±1.67** |
| $\mathcal{D}_E^2$ | 34952 | 115.76±45.65 | 8.32±1.86 | 41.79±0.44 | 0.04±0.03 | 51.28±2.07 | **57.45±2.49** |
| $\mathcal{D}_E^3$ | 39835 | 116.56±47.62 | | 42.77±1.35 | | 53.51±3.18 | **59.90±1.78** |



(a) Success rate on MR and NMR   (b) Histogram of achieved goals   (c) Distribution of goals from self-curriculum during learning
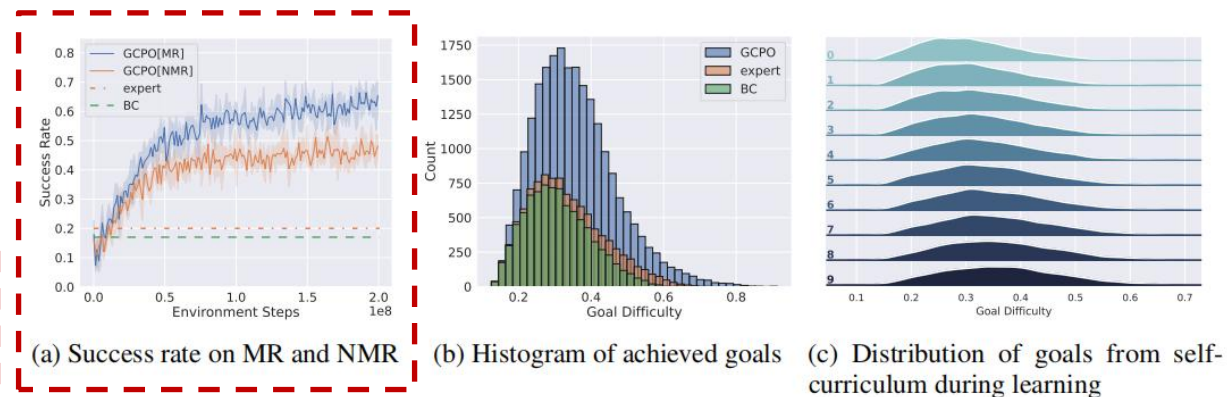
Figure 3: Main results of GCPO. 'expert' refers to the demonstrator that generates demonstrations. 'BC' refers to the pre-trained policy. Results are derived from experiments across 5 random seeds. For sub-figure (a), expert and BC are both evaluated in the NMR setting. For sub-figure (c), the vertical axis represents the training progress, where $0, 1, \cdots, 9$ correspond to $10\%, 20\%, \cdots, 100\%$ of the training progress, respectively.

➢ **GCPO is applicable to both MR and NMR problems**

- GCPO outperforms all baselines on NMR

- The learning progression of GCPO for both NMR and MR shows that GCPO is effective in solving both types of problems

## Main Results

On fixed-wing velocity vector control

Table 2: Comparison between GCPO and baselines on NMR. The mean and variance of % success rates are presented over 5 random seeds. Optimal values are highlighted in bold, and sub-optimal values are underlined.

| | Demonstration | | SAC + HER + MEGA | BC | GCPO w/o pre-training | GCPO w/o self-curriculum | GCPO |
|---|---|---|---|---|---|---|---|
| notation | #traj | traj length | | | | | |
| $\mathcal{D}_E^0$ | 10264 | 281.83±149.48 | | 17.08±0.57 | | 31.28±8.97 | **45.87±3.09** |
| $\mathcal{D}_E^1$ | 27021 | 119.64±47.55 | 8.32±1.86 | 36.54±1.97 | 0.04±0.03 | 43.49±3.85 | **49.12±1.67** |
| $\mathcal{D}_E^2$ | 34952 | 115.76±45.65 | | 41.79±0.44 | | 51.28±2.07 | **57.45±2.49** |
| $\mathcal{D}_E^3$ | 39835 | 116.56±47.62 | | 42.77±1.35 | | 53.51±3.18 | **59.90±1.78** |



(a) Success rate on MR and NMR  (b) Histogram of achieved goals  (c) Distribution of goals from self-curriculum during learning
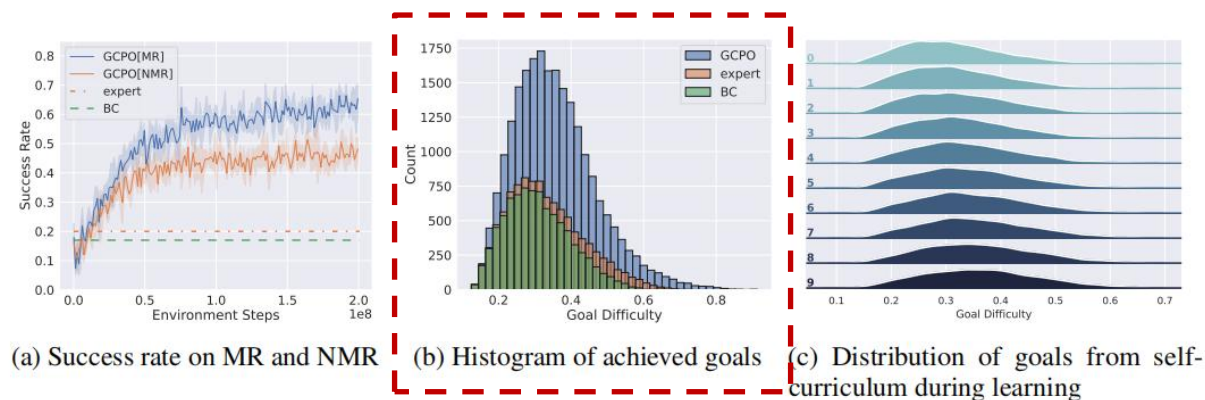
Figure 3: Main results of GCPO. 'expert' refers to the demonstrator that generates demonstrations. 'BC' refers to the pre-trained policy. Results are derived from experiments across 5 random seeds. For sub-figure (a), expert and BC are both evaluated in the NMR setting. For sub-figure (c), the vertical axis represents the training progress, where $0, 1, \cdots, 9$ correspond to $10\%, 20\%, \cdots, 100\%$ of the training progress, respectively.

➢ **Pre-training is crucial for the success of GCPO**

- without pre-training, GCPO struggles to learn meaningful skills

- even with a pre-trained policy that initially exhibits inferior performance compared to the demonstrator, GCPO's online self-curriculum learning facilitates significant improvement in the policy's performance, surpassing that of the demonstrator

## Main Results

On fixed-wing velocity vector control

Table 2: Comparison between GCPO and baselines on NMR. The mean and variance of % success rates are presented over 5 random seeds. Optimal values are highlighted in bold, and sub-optimal values are underlined.

| | Demonstration | | SAC + HER + MEGA | BC | GCPO w/o pre-training | GCPO w/o self-curriculum | GCPO |
|---|---|---|---|---|---|---|---|
| notation | #traj | traj length | | | | | |
| $\mathcal{D}_E^0$ | 10264 | 281.83±149.48 | | 17.08±0.57 | | 31.28±8.97 | **45.87±3.09** |
| $\mathcal{D}_E^1$ | 27021 | 119.64±47.55 | 8.32±1.86 | 36.54±1.97 | 0.04±0.03 | 43.49±3.85 | **49.12±1.67** |
| $\mathcal{D}_E^2$ | 34952 | 115.76±45.65 | | 41.79±0.44 | | 51.28±2.07 | **57.45±2.49** |
| $\mathcal{D}_E^3$ | 39835 | 116.56±47.62 | | 42.77±1.35 | | 53.51±3.18 | **59.90±1.78** |



(a) Success rate on MR and NMR  (b) Histogram of achieved goals  (c) Distribution of goals from self-curriculum during learning
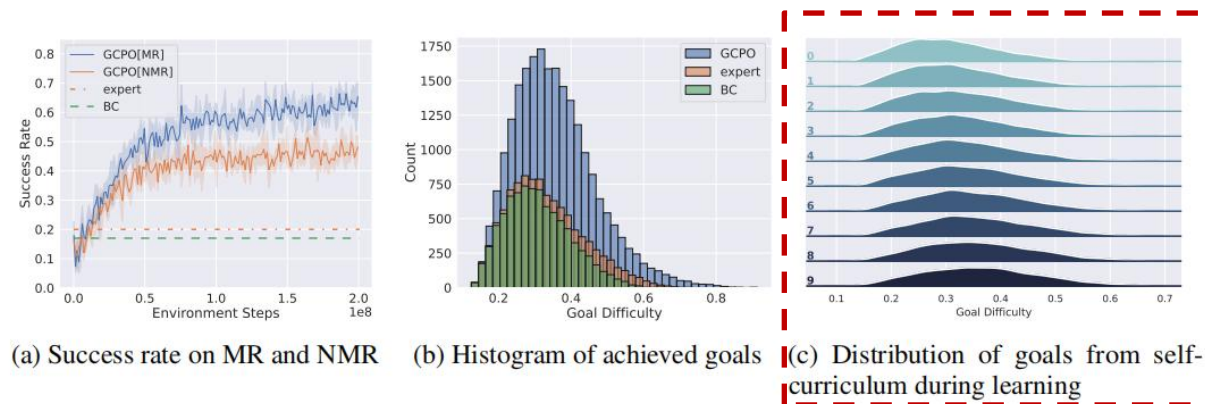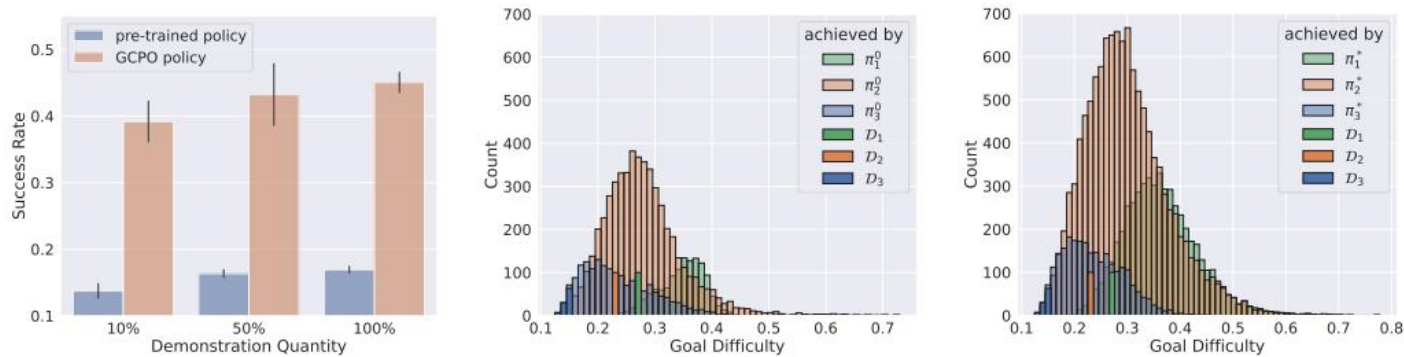
Figure 3: Main results of GCPO. 'expert' refers to the demonstrator that generates demonstrations. 'BC' refers to the pre-trained policy. Results are derived from experiments across 5 random seeds. For sub-figure (a), expert and BC are both evaluated in the NMR setting. For sub-figure (c), the vertical axis represents the training progress, where $0, 1, \cdots, 9$ correspond to $10\%, 20\%, \cdots, 100\%$ of the training progress, respectively.

➤ **Online self-curriculum facilitates the mastery of challenging goals**

- the application of self-curriculum within GCPO leads to an average 8.2% increase in policy performance compared to its absence

- online self curriculum mechanism systematically introduces more difficult goals into the learning progression as the policy gains proficiency

## Ablation Studies

1. Ablation on Quantity of Demonstrations

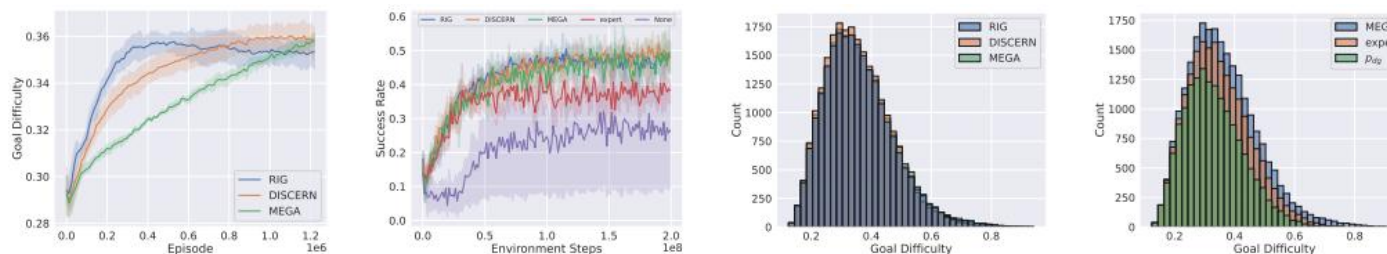2. Ablation on Goal Distribution of Demonstrations



(a) Success rate of GCPO with dif-(b) Histogram of achieved goals of (c) Histogram of achieved goals of ferent demonstration quantity the pre-trained policy GCPO policy

Figure 4: The influence of demonstration quantity and the distribution of goals covered by demonstrations on GCPO. $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3$ represent sets of demonstrations that are difficult, medium, and easy, respectively. The pre-trained policies obtained from $\mathcal{D}_1, \mathcal{D}_2$, and $\mathcal{D}_3$ are denoted as $\pi_1^0, \pi_2^0$, and $\pi_3^0$, respectively. The corresponding GCPO policies are denoted as $\pi_1^*, \pi_2^*$, and $\pi_3^*$, respectively. Results are derived from experiments across 5 random seeds.

- An increase in the quantity of demonstrations can enhance the performance of GCPO, yet the marginal gains diminish as the quantity of demonstrations grows
- when preparing demonstrations for GCPO, it is preferable to sample goals and generate demonstrations as closely as possible to the desired goal distribution $p_{dg}$

## Ablation Studies

3. Comparison on Different Self-Curriculum Methods



(a) Difficulty of goals sampled by self-curriculum methods

(b) Success rate of self-curriculum and non-curriculum methods

(c) Histogram of achieved goals of different self-curriculum methods

(d) Histogram of achieved goals of MEGA and non-curriculum methods

Figure 5: Analysis of the influence of different self-curriculum methods on the learning progression of GCPO, as well as a comparison between self-curriculum and non-curriculum methods. 'expert' and 'None' are two non-curriculum methods, where 'expert' refers to sampling goals from those that the demonstrator can achieve, and 'None' signifies directly sampling from $p_{dg}$. Results are derived from experiments across 5 random seeds.

- Fig 5(a) suggests that different self-curriculum methods exhibit distinctly different learning progressions
- Figs 5(b) and 5(c) show that there is no significant difference in performance between different self-curriculum methods, whether in the learning progression or in the final policy
- Figs 5(b) and 5(d) show that self-curriculum methods outperform non-curriculum methods in both the learning progression and the final policy performance

1. Motivation

2. Method

3. Experiments

4. Discussion

# Discussion

## Our Contributions

➢ We propose an on-policy goal-conditioned reinforcement learning framework, GCPO, designed to address the limitations of existing methods in solving non-Markovian reward (NMR) problems.

➢ We demonstrate the effectiveness of GCPO in handling both Markovian reward (MR) and NMR problems through experimental evaluation.

## Limitations

➢ In the implementation of the two components within GCPO, we employ relatively simple methods, such as behavioral cloning and Gaussian mixture model. Whether the use of alternative methods could lead to more efficient learning and better-performing policies is yet to be further validated

➢ Under the sparse reward setting, the successful training of GCPO relies on the pre-trained policy possessing a certain level of goal-achieving capability. Otherwise, if the policy achieves nothing, it becomes ineffective in establishing a self-curriculum.

➢ The specific implementation of GCPO has not explicitly incorporated components that are specifically designed to handle NMR problems. It is not clear whether integrating the most advanced methods for handling NMR problems within GCPO would lead to a more effective resolution

# Thanks for watching!

➢ Code is available at https://github.com/GongXudong/GCPO

➢ Happy to answer any questions by email:

gongxudong_cs@aliyun.com          davyfeng.c@qq.com