

Autonomous Learning From Success and Failure: Goal-Conditioned Supervised Learning with Negative Feedback

Zeqiang Zhang Fabian Wurzberger Gerrit Schmid Sebastian Gottwald
Daniel A. Braun

Abstract

Reinforcement learning faces significant challenges when applied to tasks characterized by sparse reward structures. Although imitation learning, within the domain of supervised learning, offers faster convergence, it relies heavily on human-generated demonstrations. Recently, Goal-Conditioned Supervised Learning (GCSL) has emerged as a potential solution by enabling self-imitation learning for autonomous systems. By strategically relabelling goals, agents can derive policy insights from their own experiences. Despite the successes of this framework, it presents two notable limitations: (1) Learning exclusively from self-generated experiences can exacerbate the agents’ inherent biases; (2) The relabelling strategy allows agents to focus solely on successful outcomes, precluding them from learning from their mistakes. To address these issues, we propose a novel model that integrates contrastive learning principles into the GCSL framework to learn from both success and failure. Through empirical evaluations, we demonstrate that our algorithm overcomes limitations imposed by agents’ initial biases and thereby enables more exploratory behavior. This facilitates the identification and adoption of effective policies, leading to superior performance across a variety of challenging environments.

1 Introduction

Reinforcement learning (RL) has experienced rapid progress over the past decades. With the advent of deep learning, RL has been able to tackle a wide array of complex tasks previously deemed intractable [33, 43, 50]. Nonetheless, as a framework that optimizes an agent’s behaviors through the sole guidance of a reward signal, RL algorithms encounter significant challenges in tasks characterized by sparse rewards [24, 48, 19]. This limitation is particularly pronounced in *goal-conditioned reinforcement learning*, where the agent must learn a policy to solve diverse tasks [22, 2]. In such scenarios, the agent might not receive a reward signal for extended periods due to continually changing goals, hampering its ability to improve.

To increase the efficiency of goal-conditioned reinforcement learning approaches in the face of sparse reward functions, a number of goal-relabelling strategies like Hindsight Experience Replay (HER) [2] and goal-conditioned supervised learning (GCSL) [14] have been suggested. Inspired by imitation learning [3, 28, 40], these strategies aim to utilize the agent’s actual experiences to learn from unintended effects. By posthoc relabelling achieved states as goals, they generate a substantial volume of successful training data, facilitating the learning process for the agent. HER recalculates the reward signal according to the state attained and subsequently improves the policy using off-policy RL algorithms. Conversely, GCSL adopts a direct approach to policy learning through self-imitation learning, employing relabelled goals. Although hindsight relabelling with actual achieved states enhances learning efficiency, it may also introduce biases or lead to convergence to suboptimal policies. Some methods, such as weighted GCSL [53] and normalized outcome-conditioned behavioral cloning (OCBC) [11], have sought to refine performance by weighting relabelled trajectories. However, all of these supervised learning methods disregard any information carrying negative feedback when failing to achieve an intended goal. Since supervised learning agents exclusively focuses on re-enacting past experiences with relabelled goals, they remain completely unaware of their performance relative to the originally intended goals.

To maintain the simplicity of GCSL while mitigating the aforementioned biases, we introduce a novel approach dubbed *Goal-conditioned Supervised Learning with Negative Feedback (GCSL-NF)* that utilizes both relabelled successful experience and failures by merging contrastive learning methods. In this approach, trajectories are collected based on the current policy and specified goals. Subsequently, by relabelling the achieved states as goals, these trajectories serve as positive samples for the agent’s imitation, similar to previous GCSL approaches. To obtain negative samples, the agent needs some kind of criterion to judge the quality of achieved states with respect to an intended goal in an autonomous fashion. To this end, a learned distance function assesses the quality of the sampled trajectories, comparing the achieved states against the original goals. Contrary to traditional distance functions, which typically convert states

All authors are with Institute for Neural Information Processing, Ulm University, Germany.

into a latent space to compute the Euclidean distance [45, 49], our distance leverages contrastive learning principles, enhancing its generalizability across various environments. This mechanism for identifying discrepancies proves particularly valuable in instances where agents fail to achieve their intended goals, prompting exploration and avoiding behavioral stagnation.

The main contribution of this paper is GCSL-NF, a novel method that integrates experience from original failures into the GCSL framework. To establish this method, we introduce an integrated framework for policy and distance learning. We demonstrate that our method effectively overcomes the limitations associated with agent biases and thereby promotes exploratory behavior. This, in turn, aids in the discovery and implementation of effective policies. Our experiments reveal that our method surpasses other leading GCSL-based and HER-based methods in multiple challenging environments.

2 Preliminaries

2.1 Goal-conditioned reinforcement learning

The goal-conditioned reinforcement learning (GCRL) [21] problem is usually formalized as a Markov Decision Process (MDP), represented by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{G}, \mathcal{T}, r_g, \gamma, (T))$, where \mathcal{S}, \mathcal{A} and \mathcal{G} are the state space, action space and goal space, respectively; \mathcal{T} describes the transition dynamics as $\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$; r_g is the reward function, typically defined as a binary reward $r_g(s) = \mathbb{1}(s = g)$ for goal reaching problems; γ represents the discount factor; and T specifies the horizon length. The objective in goal-conditioned RL is to devise a policy $\pi(a|s, g)$, a mapping from the current state s and the goal g to a distribution over actions a , aiming to maximize the expected cumulative return $J(\pi) = \mathbb{E}_{g \sim p(g), \tau \sim \pi_g} [\sum_{t=1}^T \gamma^{t-1} r_g(s_t)]$, where p is a distribution over goals and τ signifies the trajectory produced by the policy $\pi_g(a|s) := \pi(a|s, g)$. In hindsight experience replay, the policy is not only updated based on the original goals, but also on relabelled goals $g = s_T$.

2.2 Goal-conditioned supervised learning

In contrast to the GCRL problem defined in 2.1, where the agent seeks to learn a policy to achieve goals based on reward signals, GCSL [14] enables the agent to directly learn a policy either through behavioral cloning from expert demonstrations or via self-imitation learning with hindsight relabelling. Specifically, the dataset for training is represented as $\mathcal{D} = \{(s_t, a_t, g')\}$, with g' being either the goal from expert demonstrations or a relabelled goal defined as $g' = s_i$ for $i \geq t$. The objective is to optimize the policy by maximizing the likelihood of the recorded actions in \mathcal{D} relative to the specified goal: $\pi = \arg \max_{\pi} \mathbb{E}_{\tau \sim \mathcal{D}} [\log \pi(a_t|s_t, g')]$.

2.3 Criteria for goal attainment

Recognizing the attainment of a goal is essential for evaluating goal-conditioned approaches if the agent is to learn from negative feedback whenever it did not reach the intended goal. In discrete state spaces, goal attainment is often encoded by a simple indicator function, otherwise similarity measures or distance functions are commonly used to measure the disparity between states. However, such measures often make assumptions about the state space, for instance underlying Euclidean properties. Data-driven methods try to estimate the distance directly from observed data. In an ideal scenario, a global distance function between two states would be proportional to the number of steps necessitated by the optimal policy to transition between those states. However, in general this optimal policy is not known beforehand. Consequently, alternative criteria are required.

3 Goal-Conditioned Supervised Learning with Negative Feedback

In this section, we present the details of our proposed method. We start with an illustrative example to explain the motivation behind our approach. Following that, we detail the GCSL-NF framework as well as the approach to learning distance functions.

3.1 A motivating example

To highlight the limitations of GCSL, let's consider a scenario in a two-dimensional space where the goal is to navigate from a starting point to a target location. The agent begins at some *Point A* with the objective of reaching *Point B*. However, due to the shortcomings of its current policy, the agent ends up at *Point C* instead (see Figure 1).

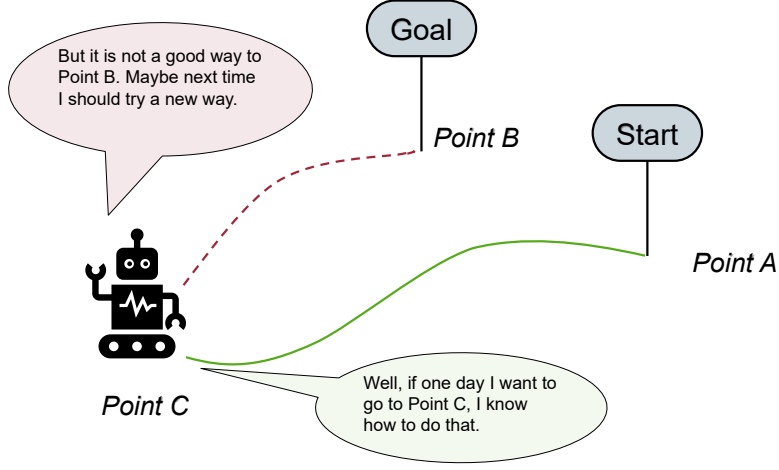


Figure 1: A Motivating Example: the agent starts from *Point A*, initiates its journey with the objective of reaching *point B*, but ultimately arrives at *point C*. The green line indicates the actual path, while red dotted lines represent the distance between the achieved state and the original goal.

GCSL hinges on the notion that a trajectory failing to reach its intended goal effectively constitutes a successful path towards the state it actually attains. Thus, suboptimal trajectories can be reinterpreted as optimal by reassigning their goals. As illustrated in Figure 1, this process involves redefining *Point C* as the target, thereby transforming the trajectory represented by the green line into a form of expert guidance for achieving this new goal.

Crucially, this reinterpretation is not the sole lesson to be drawn here. In fact, there are two insights to be gained: while the green line denotes an optimal route to the newly assigned goal *Point C*, it does not represent the most efficient path to the original goal *Point B*. GCSL capitalizes on the former perspective to refine the policy but ignores the latter, which could provide valuable cues for future exploration. We turn GCSL into GCSL-NF by incorporating this additional insight into the learning framework. Specifically, for each trajectory $\tau = \{s_1, a_1, s_2, a_2, \dots, s_T, a_T\}$ generated under the current policy $\pi(a|s, g)$ and a specified goal g , we analyze it from the perspectives of both the relabelled goal g' and the original goal g . For any given timestep t , if action a_t facilitates reaching the relabelled goal $g' = s_{t+i}$, we enhance the likelihood of selecting a_t at state s_t with goal g' through optimization of $\pi(a_t|s_t, g')$. Concurrently, the suitability of a_t in achieving the original goal g is assessed based on the proximity of the final state s_T to goal g , utilizing a distance function. This dual approach allows the agent to either refine its policy or embark on further exploration by adjusting its strategy to maximize or minimize $\pi(a_t|s_t, g)$ accordingly. The comparison of these different feedbacks is illustrated in Section 6.1.

3.2 Proposed algorithm

Our approach relies on learning a parametrized function Q_θ of actions a_t , states s_t , and goals g , which approximates the probability $p(g|s_t, a_t)$ of reaching goal g from state s_t by acting according to a_t (Step 5 below). Note that Q_θ deviates from the conventional interpretation of Q -value functions in reinforcement learning, which are derived from the Bellman equation [51]. We do not explicitly model the success probability $p(g|s_t, a_t)$, but instead, we learn an approximate (inverse) distance function p_φ between arbitrary states s, s' , which allows us to use $p_\varphi(s_T, g)$ as targets for $Q_\theta(a_t, s_t, g)$, where s_T denotes a state reached at time T . Figure 2 outlines our approach.

In the following, we list the steps in more detail:

1. Sample a goal g from the goal space \mathcal{G} using a fixed distribution p on \mathcal{G} .
2. Generate a trajectory $\tau_g = (s_0, a_0, s_1, a_1, \dots, s_T, a_T)$ from Q_θ by choosing actions according to

$$a_t = \arg \max_{a \in \mathcal{A}} Q_\theta(a, s_t, g)$$

for T steps.

3. Add the trajectory τ_g to the replay buffer \mathcal{R} .
4. Transform trajectories from \mathcal{R} into expert tuples by relabelling, which are used as positive examples for imitation learning,

$$\mathcal{T}_+ := \{(s_t, a_t, g' = s_{t+i}) : t \geq 0, i > 0, t + i \leq T\}, \quad (1)$$

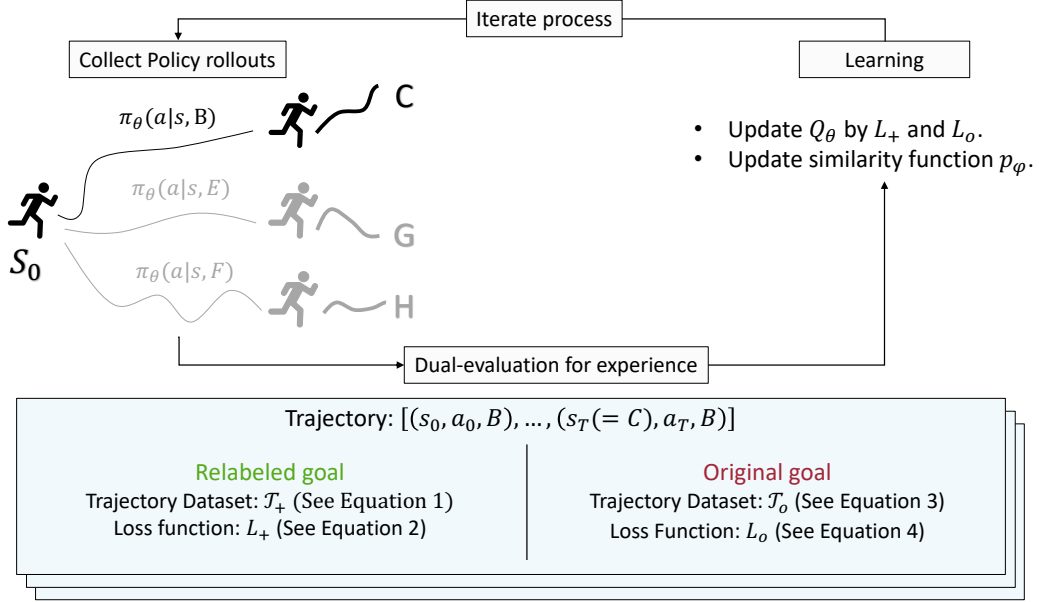


Figure 2: Diagram of GCSL-NF: The agent learns how to reach goals by sampling trajectories, reconsidering the trajectories with both relabelled goals and original goals, and then updating the policy and similarity function.

with states s_t , corresponding actions a_t , and future states s_{t+i} as relabelled goals g' . Notice, that not only the last state s_T is considered as a goal, but also any state s_{t+i} subsequent to the state s_t is valid for hindsight relabelling.

5. Imitate the expert tuples through supervised learning by minimizing the loss function

$$L_+(\theta) = \mathbb{E}_{(s_t, a_t, g') \sim \mathcal{T}_+} \left[\underbrace{H(Q_\theta(a_t, s_t, g'), 1)}_{\text{Imitation learning}} + \alpha \underbrace{\sum_{a \in \mathcal{A}} H(Q_\theta(a, s_t, g'), 0)}_{\text{Regularization}} \right], \quad (2)$$

where $H(x, y) = -y \log x - (1 - y) \log(1 - x)$ denotes the binary cross entropy between x and y (c.f. the discussion in Section 4) and α is the regularization coefficient.

6. Sample trajectories from the replay buffer \mathcal{R} to create the dataset of tuples \mathcal{T}_o with the original goals according to

$$\mathcal{T}_o = \{(s_t, a_t, s_T, g) : t > 0\}, \quad (3)$$

with states s_t , corresponding actions a_t , the final achieved state s_T , the original goals g .

7. Evaluate the tuples \mathcal{T}_o by the similarity function $p_\varphi(s_T, g)$ (detailed further in Section 3.3).
8. Compute the loss L_o over the tuples \mathcal{T}_o with the estimated similarity:

$$L_o(\theta) = \mathbb{E}_{\mathcal{T}_o} [\gamma^{T-t} H(Q_\theta(a_t, s_t, g), p_\varphi(s_T, g))], \quad (4)$$

where γ^{T-t} exponentially discounts the loss according to the discount factor γ and the remaining steps $T - t$.

9. Update the policy π_θ by minimizing the combined loss:

$$\theta_{new} \leftarrow \arg \min_{\theta} (\beta_1 L_+(\theta) + \beta_2 L_o(\theta)), \quad (5)$$

where β_1 and β_2 are the hyperparameters to balance the two losses. We set them to 1.

10. Update the parameters φ of the distance function p_φ according to Section 3.3.

3.3 Distance function p_φ

In this Section, we introduce a method of learning distance functions between states, using Monte Carlo methods [30] and contrastive learning. Broadly speaking, $p_\varphi(s, s')$ is trained to approximate the probability that the distance between two arbitrary states s, s' is smaller than a given threshold (see Figure 3).

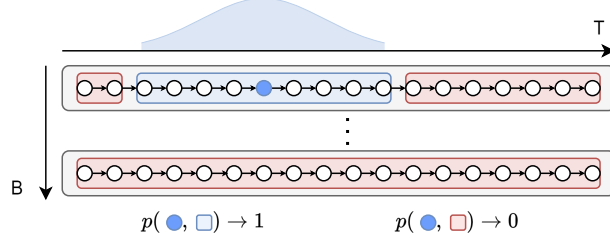


Figure 3: Illustration of the learning process for p_φ from trajectories. The figure shows how the neighbourhood relationship is defined on trajectories. The vertical axis indicates different trajectories, and the horizontal axis indicates time steps within the same trajectory. The blue region corresponds to positive samples from a unimodal distribution around the reference state, and the red region corresponds to negative samples.

Given the unknown structure and transition dynamic of the environment beforehand, it is impractical to sample states based on their actual distances. Therefore, we devise an approximation strategy which involves sampling states from collected trajectories. State pairs are sampled under three distinct types of constraints:

1. Positive pairs in \mathcal{D}_+ are extracted from the same trajectory, utilizing a unimodal distribution centered around the reference state:

$$\mathcal{D}_+ = \{(s_i, s_j) : s_i, s_j \in \tau, \tau \in \mathcal{R}, j = \text{round}(j'), j' \sim \Delta(i, n)\}, \quad (6)$$

where i and j represent time steps, and $\Delta(i, n)$ denotes a symmetric triangle distribution with mode i and threshold n . The triangular distribution is selected for its inherent boundary at the edge.

2. Negative pairs in \mathcal{D}_-^1 , drawn from the same trajectory, are defined as:

$$\mathcal{D}_-^1 = \{(s_i, s_j) : s_i, s_j \in \tau, \tau \in \mathcal{R}, |i - j| > n\}, \quad (7)$$

where i and j represent time steps. To qualify as negative pairs, the states must be significantly separated from each other, achieved by imposing a threshold n that specifies the minimum step distance between them.

3. Negative pairs in \mathcal{D}_-^2 , selected from distinct trajectories, are defined as:

$$\mathcal{D}_-^2 = \{(s_i, s_j) : s_i \in \tau, s_j \in \tilde{\tau}, \tau \neq \tilde{\tau} \in \mathcal{R}\}, \quad (8)$$

where τ and $\tilde{\tau}$ represent different trajectories, and i and j denote the respective time steps.

The function p_φ is trained to represent the probability of a given state pair to be close by minimizing the NCE loss [16, 29, 12]

$$L(\varphi) = -\left[\mathbb{E}_{(s, s') \sim \mathcal{D}_+} [\log p_\varphi(s, s')] + \sum_{i=1,2} \mathbb{E}_{(s, s') \sim \mathcal{D}_-^i} [\log(1 - p_\varphi(s, s'))]\right]. \quad (9)$$

4 Related work

Our research addresses goal-conditioned tasks within the reinforcement learning (RL) domain, characterized by multi-goal settings and sparse rewards. Kaelbling [22] firstly proposed the idea of hindsight relabelling, which augments the dataset with successful outcomes by reinterpreting failed trajectories. Hindsight relabelling approaches enhance data efficiency [2, 39, 10, 15] since trajectories generated by sub-optimal agents are still considered useful for training goal-conditioned policies, by simply pretending the actually reached states were the goals. This approach integrates well with other RL enhancements, such as curriculum learning [13].

GCSL [14] adopts hindsight relabelling in conjunction with self-imitation learning, offering stability and compatibility with offline datasets. Different from self-imitation learning which usually chooses a subset of optimal trajectories to imitate [52, 36, 17] or learn a separate value function [35, 1, 37], GCSL maximizes data reuse by training on every collected trajectory. Advancements of GCSL aim to improve performance by assigning weights to relabelled trajectories [53, 11]. Drawing inspiration from

these methods, our approach combines GCSL with negative feedback. Unlike HER-based methods that retain original trajectories, our strategy benefits from the precision and stability of supervised learning, obviating the need for explicit reward functions and thus broadening its applicability to autonomous learning across various tasks.

Our approach also utilizes distance function learning within RL, where most strategies map state representations into a latent space and calculate the l^2 -norm or cosine similarity between them, and then either directly predict distances or assess if the state distance falls below a certain threshold [45, 49, 55]. Value functions or Q-functions [44, 5, 34], successor representations [23, 44], and similarity-based approaches such as SimCLR-TT [41, 6] can also be considered as kinds of distance functions in reinforcement learning. Diverging from these approaches, we employ Monte-Carlo based contrastive learning methods [30] to estimate state pair co-occurrence frequencies, a technique agnostic to observation structure, thereby ensuring wide applicability. Compared to these approaches, which primarily capture global structural relationships within the environment, our method excels at modeling local spatial structures. Since the agent can leverage hindsight relabeling to improve policy learning, providing feedback from a local distance function can be more effective in guiding exploration—see Section 6.2 for a detailed comparison.

In our work, we make use of contrastive learning, a method predominantly applied in unsupervised or self-supervised representation learning, where positive and negative examples are used to extract discriminative features from unlabeled data [6, 7, 18, 4]. Contrastive learning has shown success not only in computer vision but also in other fields, such as natural language processing [31, 8], recommender systems [54, 20], and reinforcement learning [25, 46, 41, 27, 12]. Instead of simply using contrastive learning to learn latent representations of states as a drop-in replacement for reconstruction or other representation learning methods, Contrastive Goal-Conditioned Reinforcement Learning (Contrastive GCRL) [12] frames the goal-conditioned reinforcement learning problem *as* contrastive learning. This is possible because, under the assumption that the immediate reward is defined as the probability of reaching the goal at the next time step, the optimal goal-conditioned Q-value function is given by the (discounted) state occupancy measure evaluated at the goal state. Since the state occupancy measure can be approximated contrastively from reached states as positive examples and other states as negative examples, Contrastive GCRL learns the Q-value function by contrasting experienced goals within the trajectory with goals from other trajectories. While both Contrastive GCRL and our approach are grounded in contrastive learning principles, our method differs in utilizing the original goal to refine policy learning through negative feedback, rather than solely relying on negative examples from different trajectories.

Table 1: Comparison between different algorithms for goal-condition problems

	Vanilla RL	RL with HER	GCSL	Contrastive GCRL	GCSL-NF
Relabelled goal		✓	✓	✓	✓
Original goal	✓	✓			✓
Contrastive goal				✓	✓
Contrastive action		✓	✓		✓
External reward function	✓	✓			

We compare and summarize the differences between widely used algorithms for goal-conditioned problems and our proposed method in Table 1. With the exception of standard goal-conditioned reinforcement learning, all listed algorithms incorporate goal relabeling to improve learning efficiency. However, only our method, vanilla RL and RL with HER utilize information from the original goal. Notably, our approach does not rely on an external reward function, distinguishing it from vanilla RL and HER. Furthermore, our method employs both random actions as negative examples during policy training (contrastive action) and random states as negative examples during distance function training (contrastive goal). These mechanisms enhance learning stability and robustness, setting our approach apart from existing goal-conditioned methods.

5 Experiments

In this first part of our experimental evaluations, we assess the efficacy of our method across a variety of goal-conditioned learning tasks characterized by discrete action spaces. Our investigation is particularly focused on addressing the following research questions:

1. Can GCSL-NF effectively navigate out of policy biases to improve agent performance?
2. How does GCSL-NF compare to other algorithms in scenarios where spatial vicinity does not imply temporal vicinity, for example in the presence of obstacles?

3. What is GCSL-NF’s performance in contexts where similarities in observation space carry no information about temporal vicinity?

5.1 Environments

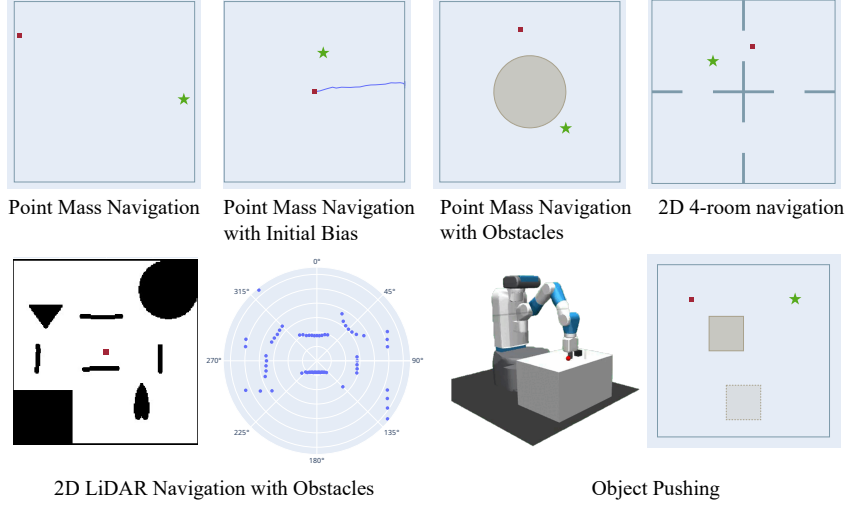


Figure 4: Environments: Red points indicate the position of the agent, and green stars denote the goal position.

As shown in Fig 4, we consider six different environments to evaluate the performance of our method: point mass navigation, point mass navigation with initial bias, point mass navigation with obstacles, 2D 4-room navigation, and 2D LiDAR navigation with obstacles and object pushing.

All tasks are 2D navigation or manipulation tasks with varying observation space and difficulty. The point mass navigation environment represents a basic goal-conditioned task where an agent is tasked with moving from an initial location to a designated target. The variant, point mass navigation with initial bias, introduces a preliminary movement bias that hinders the agent’s progress. The point mass navigation with obstacles environment and 2D 4-room navigation introduce obstacles that the agent must circumvent to reach the goal, making the optimal path distorted. For the 2D LiDAR navigation with obstacles environment, observations are derived from a LiDAR sensor. Therefore the distances in observation space have no direct correlation with the actual distances in physical space. The object-pushing environment demands that the agent manipulates an object to reach a target location, necessitating the achievement of an intermediate goal before reaching the final objective. This task is characterized by complex optimal paths and potentially deceptive $L2$ distance measurements in observation space. Details of these environments are listed below:

5.1.1 Point Mass Navigation

This scenario tasks the agent with reaching a goal location within a two-dimensional space. The observation space comprises two dimensions that represent the agent’s current coordinates, with the value from -1 to 1 , corresponding to the size of the map. The agent’s action space consists of five discrete options: move up, move down, move left, move right, or remain stationary. The step size is 0.05 , with a random normal noise of 0.01 . Both the initial state and the goal are uniformly selected from the available state space. The trajectory length is 50 .

5.1.2 Point Mass Navigation with Initial Bias

This setup mirrors the point mass navigation environment but introduces additional complexity by imposing an initial bias in the agent’s movement. Specifically, the agent is predisposed to move in a particular direction (to the right) by following gradient descent of the action labeled “right” for 10 times before the experiment.

5.1.3 Point Mass Navigation with Obstacles

In this variation of the point mass navigation task, the environment includes an additional challenge: a circular obstacle positioned centrally. The obstacle is a circle with a radius of 0.4 and a center at $(0, 0)$.

Attempts by the agent to traverse this obstacle will result in no movement, keeping the agent at its current location. The trajectory length is 70.

5.1.4 2D 4-room Navigation

In this environment, the agent is tasked with reaching a goal position within a two-dimensional space featuring a layout of four interconnected rooms, with the width of 0.4 for each passageway between rooms. The observation space comprises two dimensions that represent the agent’s current coordinates, with the value from -1.2 to 1.2 , corresponding to the size of the map. The action space is identical to that in the point mass navigation scenario. Attempts to pass through walls will result in the agent maintaining its current position. The trajectory length is 70.

5.1.5 2D LiDAR Navigation with Obstacles

This scenario challenges the agent to navigate to a goal within a two-dimensional space filled with obstacles. The agent is equipped with a 360° LiDAR sensor for detecting distances to edges or nearby obstacles. The size of the map is 200×200 . The observation space is a 64-dimensional vector, mirroring the LiDAR’s resolution. The action space includes five discrete options: move forward, move backward, turn right, turn left, or remain stationary. The step size is 10 for moving forward or backward, and $\frac{\pi}{2}$ for turning right or left. Goals are defined by the LiDAR sensor’s observation at specific positions (x, y, a) , with x, y indicating coordinates and a the rotation angle, all uniformly sampled. The agent’s position remains unchanged if it attempts to traverse through obstacles. The trajectory length is 50.

5.1.6 Object Pushing

This environment tasks the robot arm with moving a puck to a designated goal location within a two-dimensional space. The puck is a square with a side length of 0.2. The observation space encompasses four dimensions, representing both the end-effector’s and the puck’s coordinates. Values for the observation vary from -0.5 to 0.5 , corresponding to the size of the map. The action space mirrors that of the point mass navigation setup. Goals are defined by both the puck’s and the end-effector’s position, matching the observation space’s dimensions. Initial states and goals are uniformly selected from the state space. When the robot arm attempts to push the puck, the end-effector moves in the direction dictated by its action. Success in this environment requires the agent to first reach the puck, then push it to the goal, and finally move to its target location. The trajectory length is 50. This environment is a customized discrete version of Fetch Push, a standard manipulation task described in [38]. The continuous version is discussed in section 5.4.

5.2 Comparisons

To better understand the performance of our method, we conducted a comparative analysis against different reinforcement learning variants with HER and different supervised learning methods based on GCSL. Specifically, our study encompasses: (1) HER with Deep Q-Network (HER DQN) [2], which applies a value-based approach incorporating hindsight experience replay; (2) HER with Advantage Actor-Critic (HER A2C) [2, 32], employing a policy-based strategy through an actor-critic framework; (3) GCSL [14], utilizing self-imitation learning augmented by hindsight relabelling; (4) Weighted GCSL [53], an advanced iteration of GCSL that optimizes performance by applying weights to relabelled trajectories; (5) and Contrastive GCRL [12], utilizing contrastive learning to learn representations and correlations between state-action pairs and future states. Note that for the applicability of HER DQN and HER A2C we need to define an external reward function. Here we define rewards based on the Euclidean distances between the achieved states and the goal. Details of the algorithms are described in Appendix A.

5.3 Experimental results

The results, as depicted in Fig 5 and Table 3 (in Appendix C), reveal that our method either matches or surpasses the comparative methods in learning efficiency and ultimate performance across all evaluated tasks.

For the point mass navigation task, all algorithms achieve proficiency due to their simplicity. Notably, in the scenario of point mass navigation with an initial bias, only our method and Contrastive GCRL successfully overcome the bias to achieve the goal. This finding underscores GCSL-NF’s efficacy not just in recovering from suboptimal states but also in adapting to evolving environments by discarding obsolete experiences and promptly exploring new strategies.

In environments featuring barriers, such as point mass navigation with obstacles, 2D 4-room navigation, and 2D LiDAR navigation, our method exhibited superior performance. This suggests that our

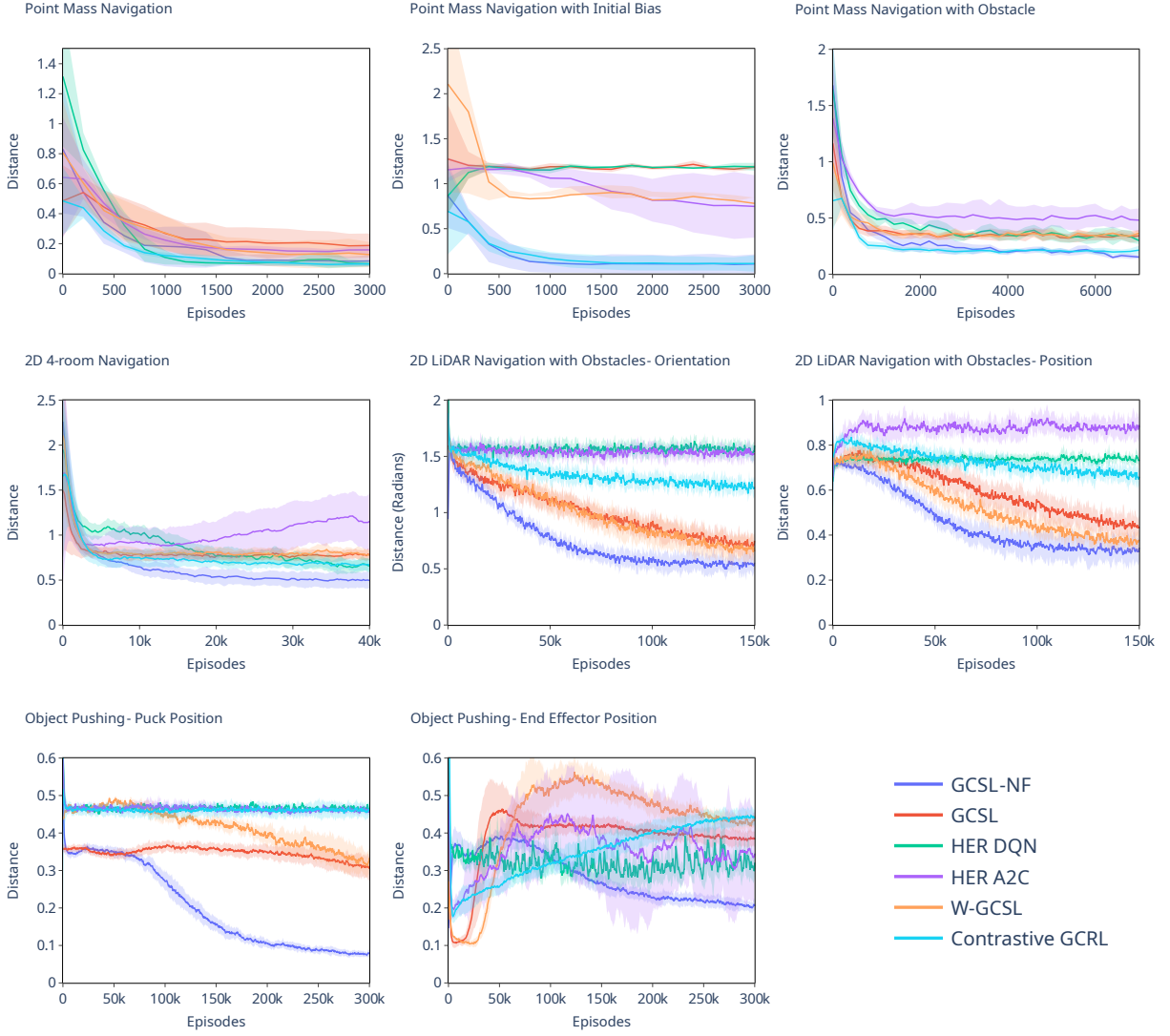


Figure 5: Results: GCSL-NF draws or outperforms other baselines on the test tasks. The results are averaged over 5 runs with different random seeds. The shaded area represents the standard deviation.

method employs a more effective exploration strategy, enabling it to discover optimal paths and avoid the pitfalls of local optima, in contrast to other methods. HER-based approaches, in particular, struggle in the 2D LiDAR navigation scenario. A contributing factor may be these methods' reliance on a reward function defined by the $L2$ distance within the observation space, which, in the context of LiDAR navigation, poorly correlates with the actual distance in the state space, thereby complicating policy updates.

In the object-pushing environment, our method stands out by successfully learning to manipulate the puck. This task implicitly requires reaching an intermediary goal—the puck—prior to achieving the final destination. GCSL-NF's similarity function consistently delivers negative feedback, motivating the agent to explore and ultimately master puck manipulation.

In conclusion, the empirical evidence highlights our method as a promising methodology for goal-conditioned learning tasks, based on a simple motivation. As there is no need to predefine a reward function, we anticipate that the method should generalize well across different tasks. Moreover, by learning from negative feedback, our method not only facilitates escape from suboptimal states but also significantly enhances its adaptability to dynamic environmental alterations.

5.4 Environments with continuous action spaces

So far, we have focused on environments with discrete action spaces. However, our method can be extended to environments with continuous action spaces.

5.4.1 Algorithm

In order to manage continuous action spaces, we utilize a similar approach as Deep Deterministic Policy Gradient (DDPG) [42, 26]. Like the discrete case, the approach for continuous action spaces also includes a parametrized function Q_{θ_1} of action, state, and goal, such that $Q_{\theta_1}(a, s, g)$ approximates the probability $p(g|s, a)$ of reaching goal g from state s by acting according to a . The function Q_{θ_1} is trained to imitate the expert tuples by minimizing the loss function

$$L_+(\theta_1) = \mathbb{E}_{(s_t, a_t, g') \sim \mathcal{T}_+} \left[\underbrace{H(Q_{\theta_1}(a_t, s_t, g'), 1)}_{\text{Imitation learning}} + \underbrace{\alpha H(Q_{\theta_1}(\pi_{\theta_2}(s_t, g'), s_t, g'), 0)}_{\text{Regularization}} \right], \quad (10)$$

where α denotes the regularization coefficient, and $H(x, y)$ denotes the binary cross entropy. Note that, while in the discrete case the regularization term consists of the sum over all possible actions (2), here only the action generated by the policy π_{θ_2} is used. The loss L_o is defined as the same as in the discrete case, and the function Q_{θ_1} is updated by minimizing the combined loss $\beta_1 L_+(\theta_1) + \beta_2 L_o(\theta_1)$.

The parameterized function $Q_{\theta_1}(a_t, s_t, g)$ is used to train the deterministic policy function $\pi_{\theta_2}(s_t, g)$. The loss for the actor is defined by evaluating the proposed actions with the function $Q_{\theta_1}(s, a, g)$:

$$L_a(\theta_2) = \mathbb{E}_{(s, s') \sim \mathcal{R}} [-Q_{\theta_1}(\pi_{\theta_2}(s, s'), s, s')]. \quad (11)$$

The distance function p_φ is trained in the same way as in the discrete case, according to Section 3.3.

5.4.2 Experiments

To evaluate the performance of the model for continuous action spaces, we consider two variants of the point mass navigation task with continuous actions, and two standard manipulation tasks including Fetch Reach and Fetch Push from [38, 9]. For both point mass navigation tasks, the observation space comprises four dimensions, representing the agent’s current coordinates and current velocity in x and y direction. The action space is continuous with four dimensions, representing the gas usage and braking force in the x and y directions. The goal is defined by the agent’s desired position and velocity (always zero). The first environment is the simple point mass navigation task without obstacles, while the second environment is the 4-room navigation task. Note that compared to the original version with discrete action space, the environments here are more challenging, as the agent cannot directly control the movement in position, but only the acceleration. For Fetch environments, we evaluate with state-based observations for the default settings of the tasks. Note that for these tasks, goals have different dimensions as the states.

We compare the performance of our method with the continuous GCSL method, the Contrastive GCRL, as well as the HER method with DDPG [2]. The original GCSL paper is limited to discrete action spaces [14], so here for comparison, we extend the GCSL method to the continuous action space by performing a regression to imitate the relabelled trajectories. The results are shown in Figure 6 and Table 4 (in Appendix C).

For the car racing point mass environment, all of the methods perform well. For the car racing 4-room environment, only HER with DDPG, Contrastive GCRL and GCSL-NF can improve the performance during the training. The continuous GCSL method is not able to learn the policy in this environment, as it collapses to a non-optimal policy and cannot get out of it due to the lack of exploration. DDPG gets its exploration from the noise added to the action, while our method utilizes a deterministic policy, where exploration is implicitly encouraged by negative feedback. For the fetch reach task, all of the algorithms show strong learning curves, although GCSL converges the fastest. For the fetch push task, GCSL-NF and Contrastive GCRL show a similar performance, while GCSL and HER DDPG fail to show any learning.

Another notable finding is that while Contrastive GCRL performs well across tasks requiring the agent to reach a specific goal position, it tends to disregard local information that varies similarly in most trajectories. In car racing tasks, for instance, it successfully reaches the target position but neglects velocity. This limitation arises from the way Contrastive GCRL selects negative examples for contrastive learning. Since its primary objective is to learn representations that capture the correlation between state-action pairs and future states, it designates future states as positive examples and randomly sampled states from other trajectories as negative examples. Since position information serves as a key distinguishing feature across different trajectories, Contrastive GCRL prioritizes it while local features are less useful to distinguish trajectories and therefore might be considered noise. Consequently, posture and velocity information are often disregarded in the learned state representations. While this property enables Contrastive GCRL to learn representations more efficiently in certain environments, it also leads to an incomplete representation of the goal.

In contrast, our method learns a local distance function to provide local feedback. While this approach may result in slower learning or fail to prioritize the most salient features of the goal space in some cases, it ensures that all features within the goal space, rather than just position, are optimized.

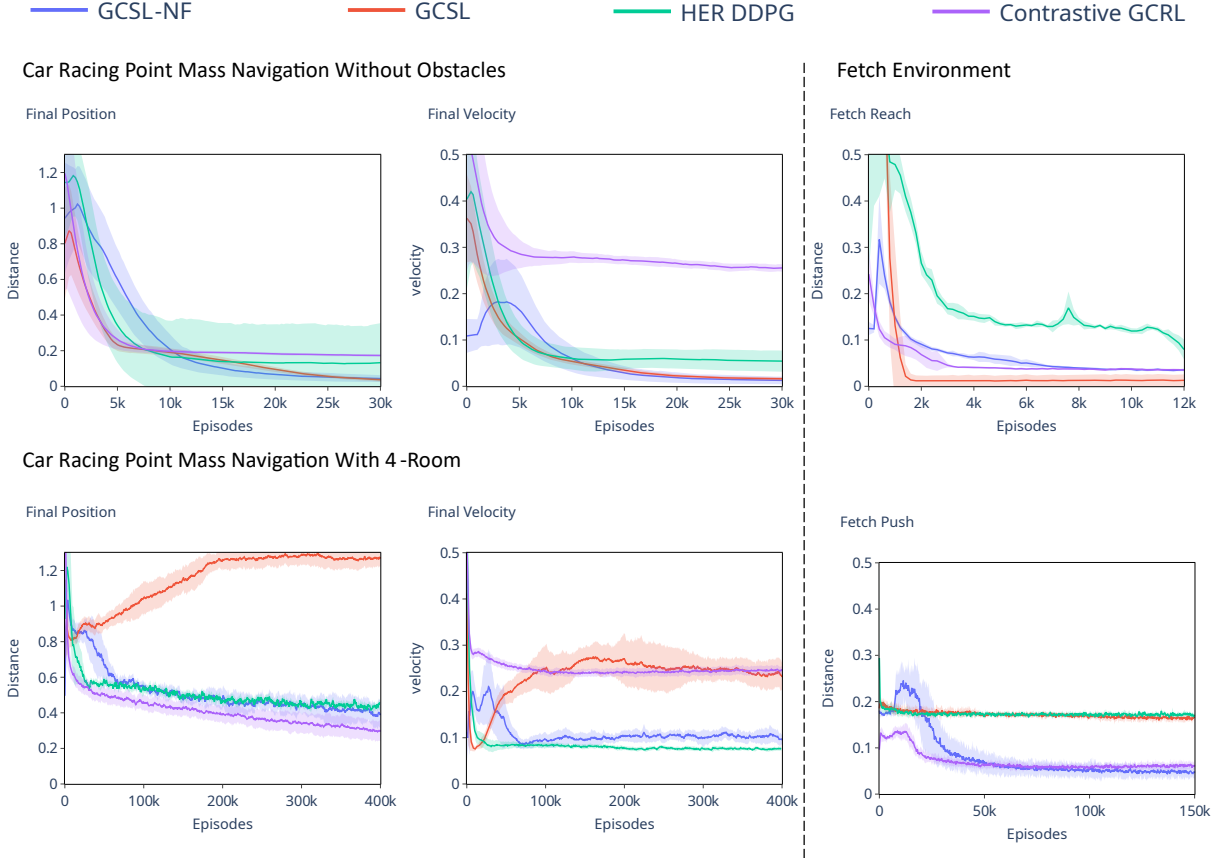


Figure 6: Results: Comparison between continuous GCSL-NF, continuous GCSL, HER with DDPG and Contrastive GCRL. The results are averaged over 5 runs with different random seeds. The shaded area represents the standard deviation.

6 Discussion

6.1 Comparison between different kinds of feedback

In this study, we investigate the impact of various feedback types on the learning process by comparing the performance of different methods: (1) GCSL-NF trained using trajectories that incorporate relabelled and original goals; (2) a policy trained exclusively with relabelled goals; (3) a policy trained solely with original goals. All of the algorithms start without an initial collection of random trajectories. We compare the different feedback types in the *Point Mass Navigation with Obstacles Environment*. The results are depicted in Figure 7.

The findings demonstrate that GCSL-NF, when trained with a mix of relabelled and original goals, surpasses the performance of the other two approaches. This suggests that integrating both feedback types significantly enhances the learning process. Training with only relabelled goals often leads to impasses, indicating the limitations of relying solely on this type of feedback. While the initial collection of random trajectories offers some improvements, it too can result in convergence to suboptimal solutions. The increase in the relative importance of the loss L_o given by the originally intended goals, as depicted in the bottom panel in figure 7, also supports this observation. With the progression of training, the proportion of loss attributable to the original goals steadily rises. This trend indicates that the agent progressively learns to imitate the current trajectory with respect to the relabelled goals. Meanwhile, the losses originating from the original goals serve to incentivize the agent’s exploration. Although training exclusively with original goals seems to achieve a satisfactory result eventually, it tends to slow down policy convergence due to the scarcity of successful trajectories, a challenge that becomes more pronounced in complex environments.

6.2 Comparison between different kinds of distance measures

To better understand the impact of different distance measures on the learning process, we compare the performance of GCSL-NF using different (inverse) distance functions: (1) distance p_φ proposed by this paper; (2) similarity measure learned by SimCLR-TT [41, 6] as the distance function; and (3) successor

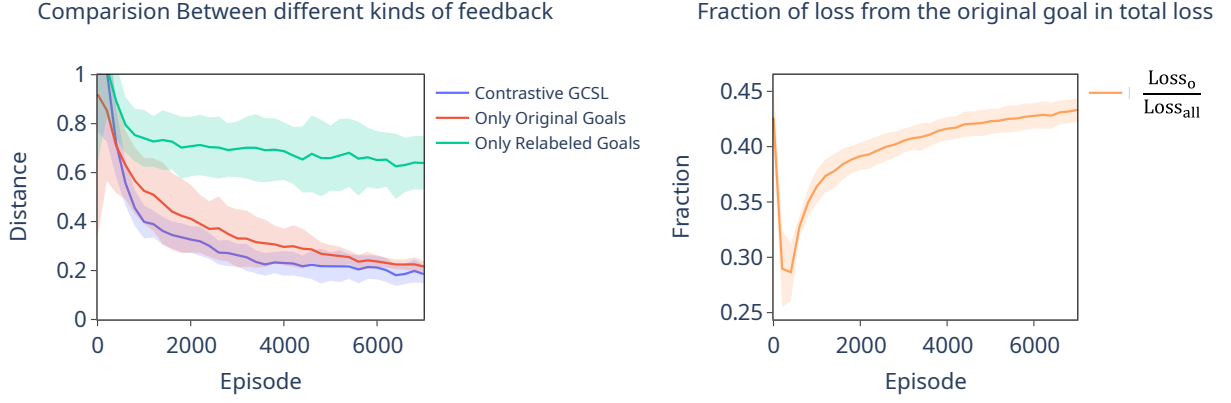


Figure 7: Results: Comparison between different kinds of feedback. The left figure shows the performance of the algorithms trained with three different kinds of feedback (positive, negative, or both) during the training. The right figure shows the fraction $\frac{L_o}{L_o+L_+}$ over time, that is the relative importance of the loss L_o under the originally intended goal. The results are averaged over 5 runs with different random seeds. The shaded area represents the standard deviation.

representation [23, 44] as the distance function.

SimCLR-TT is an effective contrastive learning method for representation learning with temporal neighbourhoods as positive pairs. The loss of SimCLR-TT is defined as:

$$L(\theta) = -\log \frac{\exp(\text{sim}(z_\theta(s), z_\theta(s^+)))}{\sum_{s^-} [\exp(\text{sim}(z_\theta(s), z_\theta(s^-)))]}, \quad (12)$$

where $z_\theta(s)$ denotes the latent representation of state s , $z_\theta(s^+)$ and $z_\theta(s^-)$ are the latent representations of positive and negative samples, respectively, and $\text{sim}(\cdot, \cdot)$ denotes the cosine similarity.

The successor representation is a method to predict the expected future state occupancy from any given state. Given a state s and a future state s' , the successor representation is defined as $M(s, s') = \mathbb{E} [\sum_{t=0}^{\infty} \gamma^t \mathbb{1}[s_t = s']]$. With the inspiration from the Bellman equation, the successor representation can be expressed in a recursive form and learned by TD-learning:

$$M(s, s') = \mathbb{1}[s_t = s'] + \gamma \mathbb{E}[M(s_{t+1}, s')]. \quad (13)$$

Table 2: Definitions of distance measures

Distance Measure	Definition
GCSL-NF- p_φ	$p_\varphi(s, s')$ (see Section 3.3)
SimCLR-TT	$p_{TT}(s, s') = (1 + \text{sim}(z_\theta(s), z_\theta(s')))/2$
Successor Representation	$p_{SR}(s, s') = \min(M(s, s')/\hat{M}, 1)$ \hat{M} is the normalization parameter defined as 0.95-quantiles of a batch

The definitions of the three distance functions can be found in Table 2. We compare the performance of GCSL-NF using different distance functions in the 4-room environment and the LiDAR environment. Relevant distance maps are shown in Figure 8 and the comparative results in Figure 9. These distance maps suggest that our distance function is good at capturing the local spatial structure of the environment, while others are more global. For the 4 room environment, our distance function can learn the existence of walls, which is crucial for the navigation task. The distance function learned by SimCLR-TT can also learn the position of the walls, but still has responses to the states that are far away from the reference state. The distance function learned by the successor representation can capture the room structure of the environment, but when the reference state is close to the wall, it becomes less effective.

While in the 4-room environment the observation is given by the position, whose difference is directly related to the distance of the states, in the LiDAR environment the observation is the LiDAR sensor reading, where the L_2 distance in the observation space does not directly correspond to the distance in the state space. Our distance function can still capture the local structure of the environment, while the distance function learned by SimCLR-TT and the successor representation are less effective locally. This suggests that our distance function is more robust to the observation space and can be more effective in environments where the observation space does not directly correspond to the state space.

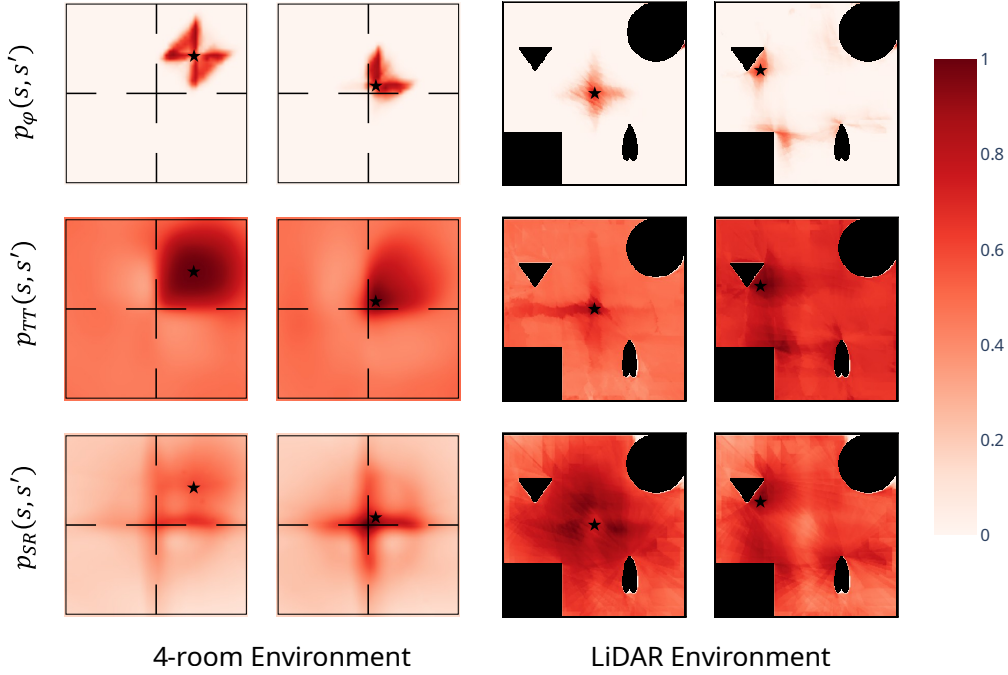


Figure 8: Heatmap: Comparison between different kinds of distance measures. The heatmap shows the values of the distance functions with respect to a fixed state which is indicated by a black star. The first row shows the distance function proposed by this paper, the second row shows the distance function learned by SimCLR-TT, and the third row shows the distance function learned by the successor representation. The first two columns are the results for the 4-room environment, and the last two columns are the results for the LiDAR environment. Note that darker colour indicates closer proximity, because we have inverse distance measures.

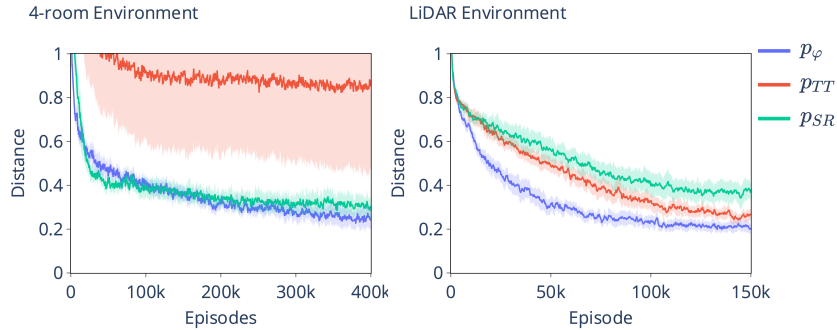


Figure 9: Results: Comparison between different kinds of distance measures. The results are averaged over 5 runs with different random seeds. The shaded area represents the standard deviation.

One potential shortcoming of our distance function is that it may be overgeneralizing similar observations across otherwise unrelated states, which may lead to a loss of the fine-grained structure of the environment. This is particularly evident in the LiDAR environment, where similar observations may be far away and represent totally different states. This problem occurs with all of the methods presented here. One potential solution is to use memory or observation traces as input to specify the corresponding states.

6.3 Policy-dependence of distance function

We also investigate the policy-dependence of the proposed distance function. We compare the distance function learned together with the policy, the distance function learned by the trajectories generated by a random policy, and the distance function learned by trajectories generated by the optimal policy in the 4-room environment. We also show the influence of the threshold n on the distance function and its policy-dependence. The heatmaps are shown in Figure 10.

The results suggest that the distance function learned together with the policy is more sensitive to the

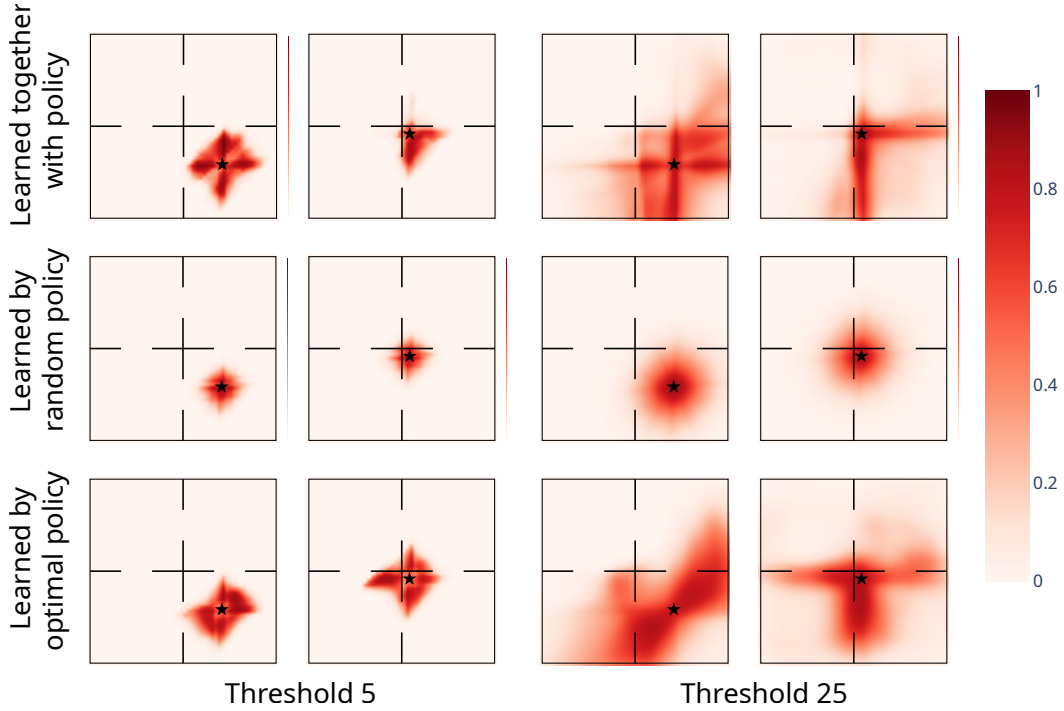


Figure 10: Heatmap: Comparison between distance function learned together with policy, distance function learned by the trajectories generated by a random policy, and distance function learned by trajectories generated by an optimal policy in the 4-room environment. The heatmap shows the values of the distance functions with respect to a fixed state which is indicated by a black star. The first two columns are the results for distance functions with threshold $n = 5$, and the last two columns are the results for distance functions with threshold $n = 25$. Note that darker colour indicates closer proximity, because we have inverse distance measures.

structure given by the walls in the environment. The distance function learned with a random policy tends to cover a smaller region, while the opposite is true for the distance function learned with the optimal policy. This can be clearly seen in the bottom half of Figure 10, where the temporal neighbourhood is defined very broadly and we use policy extremes (completely random vs. optimal) for illustrative purposes. Despite the fact that distance functions learned with different policies can lead to different properties, these differences are negligible when small thresholds are used to define temporal neighbourhood, as we did throughout the paper (see top panels in Figure 10). Therefore, our method should be fairly robust with respect to policy variations, such that the distance function is locally independent of the policy approximately. This finding is consistent with the idea that a local distance function is much easier to learn and more trustworthy.

7 Conclusion

In this study, we introduce GCSL-NF, a novel approach designed to overcome some of the inherent challenges of goal-conditioned reinforcement learning. In particular, we aim to improve upon GCSL by incorporating negative feedback, utilizing information about the desired goal within each trajectory that is otherwise disregarded. As a consequence, our approach reduces the impact of biases that most approaches trained on self-generated data suffer from. Additionally, we suggest a novel method for learning distance functions between states, drawing on Monte-Carlo based contrastive learning methods. Our empirical findings demonstrate that our approach surpasses existing baselines across a range of goal-conditioned learning tasks. A promising future direction would be to develop sophisticated goal-sampling strategies, moving beyond the current reliance on randomly selected goals from the goal space to enable novelty-driven exploration and hence more autonomy in the learning process.

References

- [1] Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin Riedmiller. Maximum a posteriori policy optimisation. *arXiv preprint arXiv:1806.06920*, 2018. 4

- [2] 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. *Advances in neural information processing systems*, 30, 2017. 1, 4, 5.2, 5.4.2, A, A
- [3] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016. 1
- [4] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020. 4
- [5] Elliot Chane-Sane, Cordelia Schmid, and Ivan Laptev. Goal-conditioned reinforcement learning with imagined subgoals. In *International conference on machine learning*, pages 1430–1440. PMLR, 2021. 4
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 4, 6.2
- [7] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020. 4
- [8] Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao, Heyan Huang, and Ming Zhou. InfoXLM: An information-theoretic framework for cross-lingual language model pre-training. *arXiv preprint arXiv:2007.07834*, 2020. 4
- [9] Rodrigo de Lazcano, Kallinteris Andreas, Jun Jet Tai, Seungjae Ryan Lee, and Jordan Terry. Gymnasium robotics, 2024. 5.4.2
- [10] Yiming Ding, Carlos Florensa, Pieter Abbeel, and Mariano Phielipp. Goal-conditioned imitation learning. *Advances in neural information processing systems*, 32, 2019. 4
- [11] Benjamin Eysenbach, Soumith Udatha, Russ R Salakhutdinov, and Sergey Levine. Imitating past successes can be very suboptimal. *Advances in Neural Information Processing Systems*, 35:6047–6059, 2022. 1, 4
- [12] Benjamin Eysenbach, Tianjun Zhang, Sergey Levine, and Russ R Salakhutdinov. Contrastive learning as goal-conditioned reinforcement learning. *Advances in Neural Information Processing Systems*, 35:35603–35620, 2022. 3.3, 4, 5.2, A
- [13] Meng Fang, Tianyi Zhou, Yali Du, Lei Han, and Zhengyou Zhang. Curriculum-guided hindsight experience replay. *Advances in neural information processing systems*, 32, 2019. 4
- [14] Dibya Ghosh, Abhishek Gupta, Ashwin Reddy, Justin Fu, Coline Manon Devin, Benjamin Eysenbach, and Sergey Levine. Learning to reach goals via iterated supervised learning. In *International Conference on Learning Representations*, 2021. 1, 2.2, 4, 5.2, 5.4.2, A
- [15] Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. *arXiv preprint arXiv:1910.11956*, 2019. 4
- [16] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In Yee Whye Teh and Mike Titterton, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 297–304, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR. 3.3
- [17] Xiaotian Hao, Weixun Wang, Jianye Hao, and Yaodong Yang. Independent generative adversarial self-imitation learning in cooperative multiagent systems. *arXiv preprint arXiv:1909.11468*, 2019. 4
- [18] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 4
- [19] Xiangkun He and Chen Lv. Robotic control in adversarial and sparse reward environments: A robust goal-conditioned reinforcement learning approach. *IEEE Transactions on Artificial Intelligence*, 5(1):244–253, 2024. 1

- [20] Chao Huang, Xiang Wang, Xiangnan He, and Dawei Yin. Self-supervised learning for recommender system. In *Proceedings of the 45th international ACM SIGIR conference on research and development in information retrieval*, pages 3440–3443, 2022. 4
- [21] L. Kaelbling. Learning to Achieve Goals. In *International Joint Conference on Artificial Intelligence*, 1993. 2.1
- [22] Leslie Pack Kaelbling. Learning to achieve goals. In *IJCAI*, volume 2, pages 1094–8. Citeseer, 1993. 1, 4
- [23] Tejas D Kulkarni, Ardavan Saeedi, Simanta Gautam, and Samuel J Gershman. Deep successor reinforcement learning. *arXiv preprint arXiv:1606.02396*, 2016. 4, 6.2
- [24] Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. Exploration in deep reinforcement learning: A survey. *Information Fusion*, 85:1–22, 2022. 1
- [25] Michael Laskin, Aravind Srinivas, and Pieter Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. In *International conference on machine learning*, pages 5639–5650. PMLR, 2020. 4
- [26] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015. 5.4.1
- [27] Minsong Liu, Luntong Li, Shuai Hao, Yuanheng Zhu, and Dongbin Zhao. Soft contrastive learning with q-irrelevance abstraction for reinforcement learning. *IEEE Transactions on Cognitive and Developmental Systems*, 15(3):1463–1473, 2023. 4
- [28] Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and Pierre Sermanet. Learning latent plans from play. In *Conference on robot learning*, pages 1113–1132. PMLR, 2020. 1
- [29] Zhuang Ma and Michael Collins. Noise contrastive estimation and negative sampling for conditional models: Consistency and statistical efficiency. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3698–3707, 2018. 3.3
- [30] Nicholas Metropolis and Stanislaw Ulam. The monte carlo method. *Journal of the American statistical association*, 44(247):335–341, 1949. 3.3, 4
- [31] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013. 4
- [32] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016. 5.2, A
- [33] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015. 1
- [34] Soroush Nasiriany, Vitchyr Pong, Steven Lin, and Sergey Levine. Planning with goal-conditioned policies. *Advances in neural information processing systems*, 32, 2019. 4
- [35] Gerhard Neumann and Jan Peters. Fitted q-iteration by advantage weighted regression. *Advances in neural information processing systems*, 21, 2008. 4
- [36] Junhyuk Oh, Yijie Guo, Satinder Singh, and Honglak Lee. Self-imitation learning. In *International conference on machine learning*, pages 3878–3887. PMLR, 2018. 4
- [37] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *CoRR*, abs/1910.00177, 2019. 4
- [38] Matthias Plappert, Marcin Andrychowicz, Alex Ray, Bob McGrew, Bowen Baker, Glenn Powell, Jonas Schneider, Josh Tobin, Maciek Chociej, Peter Welinder, et al. Multi-goal reinforcement learning: Challenging robotics environments and request for research. *arXiv preprint arXiv:1802.09464*, 2018. 5.1.6, 5.4.2

- [39] Paulo Rauber, Avinash Ummadisingu, Filipe Mutz, and Jürgen Schmidhuber. Hindsight policy gradients. *arXiv preprint arXiv:1711.06006*, 2017. 4
- [40] Joe Saunders, Chrystopher L. Nehaniv, and Kerstin Dautenhahn. Using self-imitation to direct learning. In *ROMAN 2006 - The 15th IEEE International Symposium on Robot and Human Interactive Communication*, pages 244–250, 2006. 1
- [41] Felix Schneider, Xia Xu, Markus R Ernst, Zhengyang Yu, and Jochen Triesch. Contrastive learning through time. In *SVRHM 2021 Workshop@ NeurIPS*, 2021. 4, 6.2
- [42] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. In *International conference on machine learning*, pages 387–395. Pmlr, 2014. 5.4.1
- [43] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017. 1
- [44] Kimberly L Stachenfeld, Matthew M Botvinick, and Samuel J Gershman. The hippocampus as a predictive map. *Nature neuroscience*, 20(11):1643–1653, 2017. 4, 6.2
- [45] Lorenzo Steccanella and Anders Jonsson. State representation learning for goal-conditioned reinforcement learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 84–99. Springer, 2022. 1, 4
- [46] Adam Stooke, Kimin Lee, Pieter Abbeel, and Michael Laskin. Decoupling representation learning from reinforcement learning. In *International Conference on Machine Learning*, pages 9870–9879. PMLR, 2021. 4
- [47] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016. A
- [48] Mel Vecerik, Todd Hester, Jonathan Scholz, Fumin Wang, Olivier Pietquin, Bilal Piot, Nicolas Heess, Thomas Rothörl, Thomas Lampe, and Martin Riedmiller. Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards. *arXiv preprint arXiv:1707.08817*, 2017. 1
- [49] Srinivas Venkattaramanujam, Eric Crawford, Thang Doan, and Doina Precup. Self-supervised learning of distance functions for goal-conditioned reinforcement learning. *arXiv preprint arXiv:1907.02998*, 2019. 1, 4
- [50] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019. 1
- [51] C. J. C. H. Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College (University of Cambridge), 1989. 3.2
- [52] Mao Xu and Qian Zhao. Improving deep reinforcement learning with intrinsic rewards via self-imitation learning. In *2024 39th Youth Academic Annual Conference of Chinese Association of Automation (YAC)*, pages 2017–2022, 2024. 4
- [53] Rui Yang, Yiming Lu, Wenzhe Li, Hao Sun, Meng Fang, Yali Du, Xiu Li, Lei Han, and Chongjie Zhang. Rethinking goal-conditioned supervised learning and its connection to offline rl. *arXiv preprint arXiv:2202.04478*, 2022. 1, 4, 5.2, A
- [54] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Jundong Li, and Zi Huang. Self-supervised learning for recommender systems: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 36(1):335–355, 2023. 4
- [55] Hongyu Zang, Xin Li, and Mingzhong Wang. Simsir: Simple distance-based state representations for deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 8997–9005, 2022. 4

A Comparison of Algorithms

We provide comparisons to HER DQN, HER A2C, GCSL, Weighted GCSL, and Contrastive GCRL.

- **HER DQN** [2] uses Q-learning to learn a policy that maximizes the expected return with hindsight experience replay [2]. The Q function is learned by minimizing the TD-error

$$L_Q(\varphi) = \mathbb{E}_{\tau \sim \mathcal{D}}[(r_t + \gamma \max_{a'} Q_\varphi(s_{t+1}, a', g) - Q_\varphi(s_t, a_t, g))^2]. \quad (14)$$

We put the trajectories both with the relabelled goals g' and the original goals g into the replay buffer \mathcal{D} . The reward is defined based on the $L2$ distances between the current states and the goal $r_t = -\|s_t - g\|_2$. To stabilize the learning process, we use a target network to calculate the target value [47]. The target network is updated by the main network after every trajectory. To encourage exploration, we use an ϵ -greedy policy where ϵ is set to 0.001. The discount factor γ is set as 0.99.

- **HER A2C** [2, 32] is an actor-critic method that optimizes a policy and a value function simultaneously [32]. The actor is trained to maximize the expected return by policy gradient with respect to an advantage function,

$$L_\pi(\theta) = \mathbb{E}_{\tau \sim \mathcal{D}}[-\log \pi_\theta(a_t | s_t, g) \cdot \underbrace{(r_t + \gamma V_\varphi(s_{t+1}, g) - V_\varphi(s_t, g))}_{\text{Advantage function}}], \quad (15)$$

and the critic is trained to minimize the TD-error

$$L_V(\varphi) = \mathbb{E}_{\tau \sim \mathcal{D}}[(r_t + \gamma V_\varphi(s_{t+1}, g) - V_\varphi(s_t, g))^2]. \quad (16)$$

We put the trajectories both with the relabelled goals g' and the original goals g into the replay buffer \mathcal{D} . The reward only depends on the Euclidean distances between the current states and the goal: $r_t = 1/(1 + \|s_t - g\|_2)$. The discount factor γ is set as 0.99. To encourage exploration, we sample actions from a softmax-transformed policy.

- **GCSL** [14] considers the total experience \mathcal{D} as successful trajectories by considering the actually achieved states as goals. When executing in all of the environments except *Point Mass Navigation with Initial Bias*, the first 200 trajectories are collected by executing a random policy to accumulate random trajectories, followed by a greedy policy that selects actions according to $\arg \max_a \pi(a | s, g)$.
- **Weighted GCSL** [53] is an advanced iteration of GCSL that optimizes performance by applying weights to relabelled trajectories. The loss function is

$$L_\pi(\theta) = \mathbb{E}_{\tau \sim \mathcal{D}}[-w_t \log \pi_\theta(a_t | s_t, g)], \quad (17)$$

where w_t is the weight of the trajectory defined as

$$w_t = \gamma^h \cdot \exp_{clip}(A(s_t, a_t, g)) \cdot \epsilon(A(s_t, a_t, g)), \quad (18)$$

where $A(s_t, a_t, g) = r_t + \gamma V_\varphi(s_{t+1}, g) - V_\varphi(s_t, g)$ defines the advantage function, \exp_{clip} is the exponential advantage weight with a clipped range $(0, M]$, and ϵ denotes the best-advantage weight to filter the advantage by a threshold. The value function is learned by minimizing the TD-error

$$L_{V_\varphi} = \mathbb{E}_{\tau \sim \mathcal{D}}[(r_t + \gamma V_\varphi(s_{t+1}, g) - V_\varphi(s_t, g))^2]. \quad (19)$$

For the experiments we set the discount factor γ to be 0.99, and the clipped upper range M to be 10. The first 200 trajectories are collected by executing a random policy to accumulate random trajectories, then the greedy policy $a = \arg \max_a \pi(a | s, g)$ is used.

- **Contrastive GCRL** [12] uses contrastive learning to train a Q-function relating the current state-action pair and the future states. The loss for the critic function is defined as:

$$L = \log \sigma(f(u_{s,a}, v_{s+})) + \log(1 - \sigma(f(u_{s,a}, v_{s-}))), \quad (20)$$

where $u_{s,a} = \phi(s_t, a_t)$ is the representation of the state-action pair, and $v_s = \varphi(s_f)$ is the representation of the future state. The positive pairs are sampled from the same trajectory, whereas the negative pairs are randomly sampled from different trajectories. The policy is learned to choose the action that can maximize the likelihood that the goal states occur in the future:

$$\pi(a | s, g) = \arg \max_{\pi(a | s, g)} \mathbb{E}[f(u_{s,a}, v_g)]. \quad (21)$$

B Implementation Details

We employ a consistent architecture for the neural networks representing the Q-value GCSL-NF. This architecture comprises a fully connected neural network with two hidden layers, featuring [400, 300] neurons and utilizing the SiLU activation function. The input layer merges state and goal information, while the output layer is tailored to match the action space’s dimensionality. The network utilizes a logistic activation function for its outputs.

Similarly, for the similarity function in GCSL-NF, we adopt an identical neural network architecture. This structure mirrors the previously described architecture, albeit with the output layer reduced to a single neuron. logistic activation is applied to the output of this network.

The optimizer is set as Adam optimizer with a learning rate of 0.001, the regularization coefficient α is set to 0.2, and the threshold n for the similarity function is set to 5. Our method follows the deterministic greedy policy $\arg \max_a \pi(a|s, g)$ from the beginning, while exploration is encouraged by feedback from the similarity function p_φ .

C Experiment Results

Table 3: Results for baseline algorithms and GCSL-NF

	GCSL-NF	GCSL	HER DQN	HER A2C	W-GCSL	Contrastive GCRL
Point Mass Navigation	0.085 \pm 0.037	0.187 \pm 0.080	0.066 \pm 0.019	0.158 \pm 0.029	0.127 \pm 0.084	0.065 \pm 0.014
Point Mass Navigation with Initial Bias	0.108 \pm 0.106	1.184 \pm 0.022	1.187 \pm 0.044	0.748 \pm 0.347	0.756 \pm 0.030	0.113 \pm 0.091
Point Mass Navigation with Obstacles	0.152 \pm 0.015	0.341 \pm 0.050	0.299 \pm 0.031	0.481 \pm 0.105	0.364 \pm 0.034	0.213 \pm 0.038
2D 4-Room Navigation	0.495 \pm 0.089	0.781 \pm 0.066	0.658 \pm 0.075	1.153 \pm 0.312	0.775 \pm 0.052	0.668 \pm 0.070
2D LiDAR Navigation with Obstacles - Orientation	0.555 \pm 0.070	0.749 \pm 0.087	1.597 \pm 0.086	1.520 \pm 0.099	0.654 \pm 0.147	1.227 \pm 0.096
2D LiDAR Navigation with Obstacles - Position	0.344 \pm 0.060	0.435 \pm 0.040	0.734 \pm 0.030	0.879 \pm 0.058	0.379 \pm 0.043	0.659 \pm 0.033
Object Pushing Puck Position	0.078 \pm 0.009	0.306 \pm 0.035	0.463 \pm 0.005	0.458 \pm 0.008	0.334 \pm 0.031	0.465 \pm 0.019
Object Pushing Endeffector Position	0.205 \pm 0.0158	0.384 \pm 0.019	0.321 \pm 0.037	0.345 \pm 0.112	0.445 \pm 0.016	0.440 \pm 0.022

Table 4: Results for baseline algorithms and GCSL-NF for continuous environments

		GCSL-NF	GCSL	HER DDPG	Contrastive GCRL
Car Racing Point Mass Navigation without Obstacles	Position	0.043 \pm 0.019	0.037 \pm 0.010	0.133 \pm 0.222	0.173 \pm 0.007
	Velocity	0.013 \pm 0.008	0.016 \pm 0.004	0.054 \pm 0.023	0.256 \pm 0.008
Car Racing Point Mass Navigation with 4-Room	Position	0.387 \pm 0.077	1.267 \pm 0.045	0.457 \pm 0.021	0.296 \pm 0.055
	Velocity	0.095 \pm 0.011	0.233 \pm 0.030	0.076 \pm 0.004	0.246 \pm 0.009
Fetch Reach		0.036 \pm 0.003	0.013 \pm 0.013	0.079 \pm 0.021	0.035 \pm 0.002
Fetch Push		0.049 \pm 0.014	0.165 \pm 0.009	0.171 \pm 0.003	0.060 \pm 0.007