# Discrete Factorial Representations as an Abstraction for Goal-Conditioned Reinforcement Learning

Anonymous Author(s) Affiliation Address email

# Abstract

Goal-conditioned reinforcement learning (RL) is a promising direction for train-1 ing agents that are capable of solving multiple tasks and reach a diverse set of 2 objectives. How to specify and ground these goals in such a way that we can both 3 reliably reach goals during training as well as generalize to new goals during eval-4 5 uation remains an open area of research. Defining goals in the space of noisy and 6 high-dimensional sensory inputs poses a challenge for training goal-conditioned agents, or even for generalization to novel goals. We propose to address this by 7 learning factorial representations of goals and processing the resulting representa-8 tion via a discretization bottleneck, for coarser specification of goals, through an 9 approach we call DGRL. We show that applying a discretizing bottleneck can im-10 11 prove performance in goal-conditioned RL setups, by experimentally evaluating this method on tasks ranging from maze environments to complex robotic navi-12 gation and manipulation tasks. Additionally, we prove a theorem lower-bounding 13 the expected return on out-of-distribution goals, while still allowing for specifying 14 goals with expressive combinatorial structure. 15

# **16 1 Introduction**

17 Reinforcement Learning is a popular and highly general framework [25, 57] focusing on how to select actions for an agent to yield high long-term sum of rewards. An important question is how 18 to control the desired behavior of an RL agent during both training and evaluation [24]. One way 19 to control this behavior is by specifying a reward signal [49, 53]. While this approach is very 20 general, the reward signal can be hard to design and may not be the most informative form of 21 feedback. The credit assignment problem in RL can become difficult when the reward signal is 22 sparse [59, 60, 33, 34, 54], such as policy gradients becoming nearly flat in regions where reward is 23 24 almost never achieved. Generalization can also suffer if the agent only learns one way to achieve a high reward rather than learning a diverse set of skills for coping with novel challenges [19]. 25

26 One potential way to flexibly specify and ground the desired behavior of RL agents is by training agents that receive a reward when they reach a goal specified explicitly to them [22]. In this ap-27 28 proach, called *Goal-Conditioned RL*, a single agent is trained to reach a diverse set of goals, and is given a reward only when it reaches the goal that it was instructed to reach [58, 45, 41]. This 29 provides a richer signal for the agent than simply collecting more samples oriented around a single 30 goal, as reaching multiple goals requires the agent to learn a more diverse and robust set of skills. It 31 also allows for more flexible and tightly constrained control over the desired behavior of a learned 32 agent [11, 4, 27]. Moreover, the diversity of goals seen during training should help improve both 33 credit assignment and generalization [45, 40]. 34

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

While this framework is promising, it introduces two new challenges: *goal grounding* [7, 1] and *goal specification* [3]. Goal grounding refers to defining the goal space and goal specification refers to selecting what goal the agent should try to reach in a given context. The agent is only rewarded in goal-conditioned RL when it reaches the reward which it was specified to reach, whereas in goalfree RL a reward is provided regardless of any such specification, which makes the nature of the agent's task fundamentally different.

41 What makes grounding and specifying goals

challenging? Consider trying to train a goal-42 conditioned RL agent to pick up various 43 fruits from a table. For example, we may 44 want it to pick up a red apple or a green 45 pear (illustrated in Figure 1). The number 46 of possible goals of interests may be fairly 47 small, such as the set of all valid combi-48 nations of fruits and their colors, while the 49 number of possible observations of goals is 50 extremely large when working in a rich ob-51 servation space (e.g images from a camera). 52

Goal Grounding refers to this challenge of

relating high-dimensional observations and

53

54



Figure 1: Illustration of learning *discrete* and *factorial* goal representations.

the space of relevant goals. Goal Specification refers to picking a suitable goal for the agent to 55 reach and computing an appropriate reward when it is reached. It also implies specifying goals 56 reachable in the agent's current context [36, 26]. Goal specification could be done either manu-57 ally by a developer or by another RL agent, such as a high-level agent which generates goals that 58 a lower-level agent tries to reach [11, 4, 27, 21]. Goals specified in language are an excellent fit 59 for these desiderata, as language is a compressed discrete representation which is useful for out-of-60 distribution generalization, while being compositional and expressive [18, 10, 21, 16]. At the same 61 time, connecting language feedback for an agent is non-trivial (requiring special assumptions or a 62 labeling framework) [8]. 63

We propose to learn the goal representations with self-supervised learning (either by itself, or trained 64 65 jointly with the downstream RL objective) while forcing the goal representations to be *discrete* and factorial. To perform this discretization, we use Vector-Quantization [61, 44, 30] which discretizes a 66 continuous representation using a codebook of discrete and learnable codes. The approach proposed 67 here (called DGRL) serves two complementary purposes. First, it provides a structured represen-68 tation of the raw visual goals. By representing the visual goals as a *composition* of discrete codes 69 from a learned dictionary, it makes it easier to ground unseen goals (i.e., goals not seen during train-70 71 ing) using (novel) compositions of the discrete codes learned during the training process. We show empirically that this improves generalization performance of goal-reaching policies to unseen goals 72 while remaining expressive enough. Second, the learned discrete codes can be used by another agent 73 (like a higher-level policy in hierarchical RL) to specify sub-goals to an agent (i.e., a lower-level 74 policy) to complete the task (i.e., reach the goal). In this case goal-inference is learned end-to-end. 75 The effectiveness of goal-conditioned HRL relies on the specification of semantically meaningful 76 sub-goals. Using factorial discrete sub-goals allows the higher-level policy to specify semantically 77 meaningful objectives to the lower-level policy. 78

# 79 2 Preliminaries

**Goal-conditioned RL.** We consider a goal-conditioned Markov Decision Process, where the goals  $\mathcal{G}$  lie in the state space  $\mathcal{S}$ , i.e.,  $\mathcal{G} = \mathcal{S}$  (or in observation space  $\mathcal{O}$ ). We denote a goalconditioned policy as  $\pi(a|s,g)$  (either stochastic or deterministic), and its expected total return as  $J(\pi) = \mathbb{E}\left[\sum_{t=0}^{T} R(s_t, g, a)\right]$  where the goal g is either sampled from a distribution  $\rho_g$  or provided by another higher level policy  $\pi_{\theta_h}^h(g \mid s)$ . The value function  $V^{\pi}$  is additionally conditioned on

goals, and is trained to predict the expected sum of future rewards conditioned on states and goals; 85  $V^{\pi}(s,g) = \mathbb{E}\left[\sum_{t=0}^{T} R(s_t,g,a) \mid s_0 = s; \pi\right]$ . As in standard RL, the objective in goal-conditioned RL is to maximize the expected discounted returns induced by the goal-conditioned policy. 86

87

Hierarchical Reinforcement Learning. We consider goal-conditioned settings in which the goals 88 are specified in the observation space. In the hierarchical reinforcement learning (HRL) setup, goals 89 are provided by a higher level policy  $\pi_{\theta_h}^h(g|s_t)$ . The higher level policy operates at a coarser time 90 scale and chooses a goal  $g_t \sim \pi^h_{\theta_h}(g|s_t)$  to reach for the lower level policy every K steps. The 91 lower level policy executes primitive actions  $\pi_{\theta_l}^l(a|s_t, g_t)$  to reach the goals specified by the high-92 level policy and is trained to maximize the intrinsic reward provided by the high-level policy. The 93 higher level policy is trained to maximize the external reward i.e., the reward function specified by 94 the MDP. Both the higher and lower level policies can be trained with any standard RL algorithms, 95 96 such as Deep Q-Learning (DQN) [33] or policy optimization based algorithms [47, 48]. Alternately, one can also consider another setup for goal-conditioned RL, where the goals are provided by the 97 environment  $g_1, \ldots, g_L$  and are part of the state or observation space. At each episode of training, 98 one of the goals is sampled from the distribution of goals  $\rho_q$  and the policy is trained to reach the 99 sampled goal. At test time, the agent can be evaluated either on its ability to reach goals within the 100 distribution  $\rho_{q}$ , or for its out-of-distribution generalization capability to reach new kinds of goals. 101

In this work, we consider both the HRL and goal-conditioned setups, and evaluate the significance 102 of learning a factorial representation of discrete latent goals in a series of complex goal-conditioned 103 tasks. 104

Vector Quantized Representations. VQ-VAE [61, 44, 30] discretizes the bottleneck representation 105 of an auto-encoder by adding a codebook of discrete learnable codes. The input is passed through an 106 107 encoder. The output of the encoder is compared to all the vectors in the codebook, and the codebook vector closest to the continuous encoded representation is fed to the decoder. The decoder is then 108 tasked with reconstructing the input from this quantized vector. 109

Self-supervised learning of representations. Several papers [28, 55, 50, 31] have demonstrated 110 the benefits of using a pre-training stage where the representations of raw observations are learned 111 using self-supervised objectives in a task-agnostic fashion. After the pre-training stage, the represen-112 tations can be used for (and potentially also fine-tuned on) downstream tasks. These self-supervised 113 representations have been shown to improve sample efficiency. 114

#### 3 **Discrete Goal-Conditioned Reinforcement Learning (DGRL)** 115

In this section, we provide technical details on the proposed framework, DGRL, which consists of 116 three parts: (a) learning representations of raw visual observations through self-supervised repre-117 sentation objectives, (b) processing the resulting representations via a learned dictionary of discrete 118 codes, and (c) using the resulting discrete representations for downstream goal-conditioned and HRL 119 tasks. With the learnt discrete goal representations, we describe in Section 5 how they can accelerate 120 learning in complex navigation and manipulation tasks. The goal representation can be learned at 121 the same time as the downstream-RL objective or pre-trained with self-supervised learning and then 122 used as a fixed representation for RL. 123

#### 3.1 Self-Supervised Goal Representation Learning 124

One can use any off-the shelf self-supervised method for learning representations of the raw state and 125 the goal observations. We denote by  $\phi$  the encoder network that takes as input the raw observation 126 and maps it to a continuous embedding:  $z_e = \phi(o_t)$ . Here, we explore two different self-supervised 127 techniques for learning representations. For simpler environments, we use a simple autoencoder with 128 the reconstruction objective. For more complex environments, we use the Deep InfoMax approach 129 [31] which optimizes for a contrastive objective as a proxy to maximize the mutual information 130 between representations of nearby states in the same trajectory. 131



Figure 2: **Summary of Proposed DGRL Model** for improving *goal grounding* and *goal specification* by making goal representations *discrete* and *factorial*. We learn a latent representation for both observations and goals using a self-supervised learning method (sec. 3.1). We convert the learnt latent representation into discrete latents based on a VQ-VAE quantization bottleneck with multiple factor outputs (sec. 3.2). We use the resulting discrete representations for downstream RL tasks: (i) to train a goal-conditioned policy or value function, and (ii) in the context of goal-conditioned hierarchical reinforcement learning (sec. 3.3).

### 132 3.2 Processing continuous representations via a discrete codebook

We learn discrete representations by using the vector-quantization method from the VQ-VAE paper [61] and follow the multi-factor setup used in Discrete-Value Neural Communication [30]. The discretization process for each vector  $z_e \in \mathcal{H} \subset \mathbb{R}^m$  is described as follows. First, vector  $z_e$  is divided into G segments  $c_1, c_2, \ldots, c_G$  with  $z_e = \text{CONCATENATE}(c_1, c_2, \ldots, c_G)$ , where each segment  $c_i \in \mathbb{R}^{m/G}$  (such that m is divisble by G). Each continuous segment  $c_i$  is mapped separately to a discretized latent vector  $e \in \mathbb{R}^{L \times (m/G)}$  where L is the size of the discrete latent space (i.e., an L-way categorical variable):

$$e_{o_i} = \text{DISCRETIZE}(c_i), \quad \text{where } o_i = \operatorname*{arg\,min}_{j \in \{1, \dots, L\}} ||c_i - e_j||.$$

These discrete codes, which we call the factors of the continuous representation  $z_e$ , are concatenated to obtain the final discretized vector  $z_q$ :

$$z_q = \text{CONCATENATE}(\text{DISCRETIZE}(c_1), \text{DISCRETIZE}(c_2), \dots, \text{DISCRETIZE}(c_G)).$$
(1)

The loss for vector quantization is:  $\mathcal{L}_{\text{discretization}} = \frac{\beta}{G} \sum_{i}^{G} ||c_i - \text{sg}(e_{o_i})||_2^2$ .

The training procedure closely follows both [30] and [61]. Here, sg refers to a stop-gradient op-136 eration that blocks gradients from flowing into  $e_{o_i}$ , and  $\beta$  is a hyperparameter which controls how 137 strongly we move the codes toward the encoded values. Unlike [30], we used a moving average to 138 update the code embeddings rather than learning them directly as parameters. We update  $e_{o_i}$  with 139 an exponential moving average to encourage it to become close to the selected output segment  $c_i$ . 140 This update sets the new value of  $e_{o_i}$  to be equal to  $\eta e_{o_i} + (1 - \eta)c_i$ , where the value of  $\eta$  is a fixed 141 hyperparameter controlling how quickly the moving average updates. The term  $\sum_{i=1}^{G} ||c_i - \text{sg}(e_{o_i})||_2^2$ 142 is the commitment loss, which only applies to the target segment  $c_i$  and trains the encoder that out-143 puts  $c_i$  to make  $c_i$  stay close to the chosen discrete latent vector  $e_{o_i}$ . We trained the VQ-quantization 144 process together with other parts of the model by gradient descent. When there were multiple  $z_e$ 145 vectors to discretize in a model, the mean of the codebook and commitment loss across all  $z_e$  vectors 146 was used. 147

**Summary.** The multiple steps described above can be summarized by  $z_q = q(z_e, L, G)$ , where  $q(\cdot)$ 148 is the whole discretization process using the codebook, L is the codebook size, and G is the number 149 of factors per vector. We train the representations for both the state and goal observations with a 150 discretization bottleneck on the continuous representations resulting from the self-supervised pre-151 training. The number of factors G is a hyper-parameter. In our experiments, we try with different 152 number of discrete factors G = 1, 2, 4, 8, 16, and found that G = 16 worked the best. Using 153 discretization with more factors slightly increases computation but reduces the number of model 154 parameters due to the codebook embeddings being reused across the different factors. 155

#### 156 3.3 Using representations for downstream RL

<sup>157</sup> We use the discrete representations for downstream RL tasks: (i) to train a goal-conditioned policy, <sup>158</sup> and (ii) in the context of hierarchical reinforcement learning.

**Goal-conditioned RL.** Defining goals in the space of noisy, high-dimensional sensory inputs poses 159 a challenge for generalization to novel goals because the encoder that maps the goal observations to 160 the low dimensional latent representation may fail to generalize. One way to address this is to embed 161 the continuous latent representation into a discrete representation such that the representation of the 162 novel goal is mapped to the fixed set of latent discrete codes, and hence facilitate generalization 163 to new combinations of these codes while making it easy for downstream learning to figure out 164 the meaning of each discrete code. In this setup instead of feeding the continuous state and goal 165 embedding we used their discretized versions, thus grounding goal representations in the input space. 166

We use the resulting representations for training a goal-conditioned policy  $a_t \sim \pi_{\theta_l}^l(a|s_t,g_t)$  or a 167 goal-conditioned action value function  $Q(s_t, a_t, g_t)$ . At each episode of training, a goal is sampled 168 from the distribution of goals  $\rho_q$  and the agent gets rewarded for reaching the sampled goal. This 169 170 reward can either be an *extrinsic reward* which is defined as part of the environment or an *intrinsic* reward which is defined as part of the algorithm. In DGRL, we define the intrinsic reward as the 171 fraction of discrete factors which match in the respective representations of the goal observation 172 and the state observation. At test time, the agent can either be evaluated to reach goals within the 173 distribution  $\rho_a$ , or for its generalization capability to reach goals not seen during training. 174

Hierarchical RL. The higher level policy  $g_t \sim \pi_{\theta_h}^h(g \mid s_t)$  outputs a continuous representation of goals g by conditioning on the states every K time-steps, it can also output a sub-goal  $s_g$  by conditioning on both states s and environment goals g, i.e.,  $\pi_{\theta_h}^h(s_g \mid s_t, g_t)$ . The effectiveness of goal-conditioned HRL relies on the specification of semantically meaningful sub-goals. Learned codebooks (section 3.2) consisting of a set of discrete codes can be used by a higher level policy to *specify* which goal to reach to a lower level policy. The use of learned codebooks ensures that the goal specified by the higher level policy is grounded in the space of raw-observations.

In section 5, we empirically show the benefits of the proposed approach for training goal-reaching policies or goal-conditioned value functions, as well as in a goal-conditioned hierarchical RL setup.

#### **184 4 Theoretical Analysis**

In this section, the goal discretization is shown to improve generalization to novel goals by enhanc-185 ing the concentration of the goal distribution within each neighborhood of discretized goal values; 186 i.e., by decomposing the goal probability p(g) into  $p(g) = \sum_k p(g|g \in \mathcal{G}_k) p(g \in \mathcal{G}_k)$  with the 187 neighborhood set  $\{\mathcal{G}_k\}_k$ , it improves the overall performance in p(g) by increasing the concentra-188 tion in  $p(g|g \in \mathcal{G}_k)$ . Intuitively, this is because the discretization removes varieties of possible goal 189 values  $g \in \mathcal{G}_k$  for each neighborhood  $\mathcal{G}_k$ . To state our result, we define  $\varphi_{\theta}(g) = \mathbb{E}_{s_0}[V^{\pi}(s_0, g)]$ , 190 where  $\theta \in \mathbb{R}^m$  is the vector containing model parameters learned through n goals observed dur-191 ing training phase,  $g_1, \ldots, g_n$ . We denote the discretization of g by q(g), and the identity function 192 by id as id(g) = g. Let  $\mathcal{Q} = \{q(g) : g \in \mathcal{G}\}$  and d be a distance function. We use  $\mathcal{Q}_i$  to de-193 note the *i*-th element of Q (by ordering elements of Q with an arbitrary ordering). We also define 194  $[n] = \{1, \dots, n\}, \mathcal{G}_k = \{g \in \mathcal{G} : k = \arg\min_{i \in [|\mathcal{Q}|]} \hat{d}(q(g), \mathcal{Q}_i)\}, \mathcal{I}_k = \{i \in [n] : g_i \in \mathcal{G}_k\}, \text{ and } \mathcal{I}_Q = \{k \in [|\mathcal{Q}|] : |\mathcal{I}_k| \ge 1\}. \text{ We denote by } c \text{ a constant in } (n, \theta, \Theta, \delta, S).$ 195 196

- <sup>197</sup> The following theorem (proof in Appendix D) shows that the goal discretization improves the lower
- bound of the expected sum of rewards for unseen goals  $\mathbb{E}_{g \sim \rho_g}[(\varphi_{\theta} \circ \varsigma)(g))]$  by the margin of  $\omega(\theta)$ :

**Theorem 1.** For any  $\delta > 0$ , with probability at least  $1 - \delta$ , the following holds for any  $\theta \in \mathbb{R}^m$  and  $\varsigma \in \{id, q\}$ :

$$\mathbb{E}_{g \sim \rho_g}[(\varphi_{\theta} \circ \varsigma)(g))] \ge \frac{1}{n} \sum_{i=1}^n (\varphi_{\theta} \circ \varsigma)(g_i) - c\sqrt{\frac{2\ln(2/\delta)}{n}} - \mathbb{1}\{\varsigma = \mathrm{id}\}\omega(\theta)$$

where  $\omega(\theta) = \frac{1}{n} \sum_{k \in \mathcal{I}_Q} |\mathcal{I}_k| \left( \frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} \varphi_{\theta}(g_i) - \mathbb{E}_{g \sim \rho_g}[\varphi_{\theta}(g)|g \in \mathcal{G}_k] \right)$ . Moreover, for any compact  $\Theta \subset \mathbb{R}^m$ , if  $\varphi_{\theta}(g)$  is continuous at each  $\theta \in \Theta$  for almost all g and is dominated by a function  $\chi$  as  $|\varphi_{\theta}(g)| \leq \chi(g)$  for all  $\theta \in \Theta$  with  $\mathbb{E}_g[\chi(g)] < \infty$ , then the following holds:

$$\sup_{\theta \in \Theta} |\omega(\theta)| \xrightarrow{P} 0 \quad when \quad n \to \infty.$$

### 201 Proof. Detailed proof provided in the appendix D

Without the goal discretization, we incur an extra cost of  $\omega(\theta)$ , which is expected to be strictly positive since  $\frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} \varphi_{\theta}(g_i)$  is maximized during training while  $\mathbb{E}_{g \sim \rho_g}[\varphi_{\theta}(g)|g \in \mathcal{G}_k]$  is not. Thus, the goal discretization can improve the expected sum of rewards for unseen goals by the degree of  $\omega(\theta)$ , which measures the concentration of the goal distribution in each neighborhood. This extra cost  $\omega(\theta)$  goes to zero when the number of goal observations *n* approaches infinity.

# 207 5 Experiments

The main goal of our experiments is to show that goal discretization can lead to sample efficient 208 learning and generalization to novel goals, in goal-conditioned RL. First, we directly study this 209 by training on environments with a set of goals (such as 8 positions within a gridworld) and then 210 evaluating the agent's ability to reach a position within the gridworld which it was not tasked with 211 reaching during training. Second, we consider hierarchical goal-conditioned RL, in which a higher-212 level agent generates goals that a lower-level agent is tasked with reaching. In this case, the task 213 of reaching novel goals occurs *organically* as the higher-level model selects new goals. This setup 214 also shows the advantages of DGRL for goal specification. A secondary goal of our experiments is 215 to show that using many discrete factors is often critical for optimal performance, which proves the 216 value of *factorization* in grounding goals. 217

We evaluate our proposed method DGRL by integrating it into existing state-of-the-art goal-218 conditioned and hierarchical RL tasks. Experimentally, we analyse DGRL on several challeng-219 ing goal-conditioned testbeds that have previously been used in the goal-conditioned RL commu-220 nity. DGRL in principle can be applied to any existing downstream goal-conditioned RL tasks, we 221 demonstrate improvements on five distinct goal-conditioned RL tasks. We consider maze navigation 222 tasks where images are used as observations and we show improved generalization to novel goals. 223 We integrate DGRL to an existing goal-conditioned baseline for navigating procedurally-generated 224 hard exploration Minigrid environments [9] and find that it outperforms state-of-the-art exploration 225 baselines. We also show improvements with DGRL on continuous control (Ant) navigation and 226 manipulation tasks, where goals come from a high-level controller. Finally, we show that discrete 227 representations also significantly improve sample efficient learning on a challenging vision-based 228 robotic manipulation environment. 229

Learning to Reach Diverse and Novel Goals. We study a gridworld navigation task in which an agent is trained to reach a goal from a small finite set of training goals, and during evaluation is tasked with reaching a novel goal unseen during training. This is a navigation task with a pixel-level observation space showing the position of the agent and the goal in a gridworld. We consider two mazes spiral and single-loop topology. Experiment setup given in Appendix C.1.

For this task, we train a goal-conditioned Deep Q-Learning (DQN) agent, and use a pre-trained representation  $\phi(\cdot)$  where the encoder is trained using data from a random rollout policy. Because the gridworld is small the random rollout policy achieves good coverage of the state space, so we found this was sufficient for learning a good goal representation. At each episode, a specific goal is randomly sampled from a distribution of goals, and the DQN agent is trained to reach the specified goal for that episode. During evaluation, we test the learned agent on goals either from the training distribution, or not seen during training.

Furthermore, for this task, we additionally use an intrinsic reward for exploration of the goal-DQN 242 agent. Since we learn a discrete factorial representation of the goal, we compute an exploration 243 bonus based on the discrete latent codebooks; i.e., we embed the states and goals using the learned 244 codes and then compute an intrinsic exploration bonus based on the fraction of learned factors that 245 246 match. For the baseline goal-DQN agent, we provide an additional reward bonus based on the cosine distance between continuous embeddings of the state observation and goal. Figure 3 shows 247 that DGRL significantly outperform a continuous baseline goal DQN agent, when trained on either 248 four goals or eight goals. We evaluate generalization to 4 novel goals unseen during training (Figure 249 4) and demonstrate improved generalization to novel goals. 250



Figure 3: Loopworld maze environment. We show that for different discrete factors of 4, 8, 16, DGRL outperforms a goal-DQN baseline agent with continuous goal representations. As we increase the number of factors G to 16, the expressivity of the discrete goal representation increases, lowering the odds of the factors being the same. This provides a better intrinsic reward signal for exploration, resulting in faster convergence for DGRL integrated on a goal-DQN agent.



Figure 4: SpiralWorld environment (left). Generalization to test distribution of 4-goals in a Spiral-World environment (left). We show the total number of steps to solve all test set goals, when trained on either an 8-goal or 16-goal training distribution.

In the previous experiment, we evaluated the generalization ability of DGRL by showing that learning discrete factorial representations of goals can improve generalization to novel goals. Now, we consider various setups in which a goal generating agent specifies goals using the learned codebook and the goal-conditioned agent is tasked with reaching the goals specified by the goal generating agent. We test various settings, where the goal generating agents is parameterized as an adversarial teacher [5], or as a higher-level policy in the case of hierarchical RL.

Procedurally Generated MiniGrid Exploration Task. We follow the experimental setup of [5] and [43] and evaluate DGRL on procedurally generated MiniGrid environments [9]. In [5], a goal-generating teacher proposes goals to train a goal-conditioned "student" policy. We integrate DGRL



Model	KCmedium
AMIGO + DGRL, G=16 AMIGO + DGRL, G=8	$.96 \pm .01 \\ .70 \pm .16$
AMIGO RIDE RND ICM	$.93 \pm .06$ $.90 \pm .00$ $.89 \pm .00$ $.42 \pm .21$

Figure 5: Performance comparison of the Amigo baseline [5] with adversarially intrinsic goals, and adding DGRL for discretization of the goals.

Table 1: We added DGRL on top of the Amigo baseline implementation provided by the authors.



Figure 6: In KeyChest, the agent (A) starts from a random stochastic position, picks up the key (K) and then uses the key to open the chest (C). We find DGRL improved sample efficiency over the HRAC baseline.

on top of AMIGO [5] and compare DGRL on a hard exploration task with state-of-the-art exploration baselines. Experimental results are summarized in Table 1 and more details provided in Appendix C.3. Note that unlike RIDE and RND, we do not provide an additional exploration bonus
 to DGRL, and find that DGRL can still solve this hard exploration task more efficiently.

Goal Grounding in KeyChest Maze Navigation Domain. We consider a simple discrete state 264 action KeyChest maze navigation task, following [68], where discrete goals in the state space are 265 provided by a higher level policy. For this task, to integrate DGRL, we learn an embedding  $\phi(\cdot)$  of 266 the goals, then discretize the representation with a learned codebook. We compare with a baseline 267 HRAC [68] agent (details in Appendix C.2). Figure 6 shows an illustration of the KeyChest environ-268 ment and a performance comparison of DGRL with different group factors G. Using fewer factors 269 (G = 4) performs worse than the HRAC baseline, whereas using a larger number of factors (G = 8)270 or G = 16) improves the sample efficiency of the goal reaching agent, providing evidence for the 271 benefits of factorization. 272

Ant Manipulation Control Domains. We employed DGRL on three different continuous control tasks: AntMazeSparse, AntFall and AntPush task. We emphasize that these tasks are the more challenging counterparts of AntGather and AntMaze tasks, typically used in the hierarchical RL community [35, 37]. Figure 7 provides illustration of these tasks. We evaluate goal discretization by integrating DGRL to the state-of-the-art HRAC baseline. Details of the experimental setup are provided in Appendix 5. Figure 7 shows that specifying the goals using the learned codebook helps DGRL achieve a higher success rate compared to the HRAC baseline.

Ant Navigation Maze Tasks. We consider Ant navigation tasks that require extended temporal rea-280 soning, following the setup in Reinforcement learning with Imagined Subgoals [6, RIS]: a U-shaped 281 maze, and an S-shaped maze (the S-shaped maze is shown in Figure 8). The ant navigating in the 282 maze is trained to reach any goal in the environment. The agent is evaluated for generalization in 283 an extended temporal setting with a difficult configuration, we compare the success rate of DGRL 284 integrated on top of RIS with several baselines. We emphasize the difficulty of these tasks, where 285 existing baselines like soft actor critic [19, SAC] and temporal difference models [42, TDM] fail 286 completely. Results in Figure 8 show that DGRL improves the sample efficiency over the RIS base-287 line. Additional experimental setup and environment configurations are provided in Appendix C.5. 288 289



Figure 8: Performance comparison with the success rate of reaching goal positions during evaluation in an extended temporal configuration of the U-shaped and maze-shaped Ant navigation tasks. We find that integrating DGRL with RIS can lead to more sample efficient convergence on these tasks, while baselines such as SAC and TDM (not shown) fail completely on both AntU and AntMaze as reported by [6]. The RIS baseline is based on raw data provided by the authors.

Vision Based Robotic Manipulation. Finally, we assess DGRL on a hard vision-based robotic 290 manipulation task, and use the same setup as in section 5 to integrate DGRL with the state-of-the-291 art RIS baseline on the Sawyer task in Figure 9. This manipulation task is adapted from [39], where 292 the baseline RIS is already shown to be superior to previous goal conditioning methods. The task 293 of the agent is to control a 2-DoF robotic arm from image input and move a puck positioned on 294 the table. The Sawyer task is designed for training and generalization. At test time, it evaluates the 295 agent's success at placing the puck in desired positions in a temporally extended configuration. This 296 is a challenging vision-based complex motor task, since test time generalization requires temporally 297 extended reasoning. Results in Figure 9 show that DGRL improves the sample efficiency over the 298 RIS baseline. Details of the experimental setup is provided in Appendix C.6. 299



Figure 9: Sawyer Robotic Manipulation Task. Integrating DGRL with RIS (with a larger number of factors, G=8 and G=16, while G=4 fails), improves over the RIS baseline

# 300 6 Conclusion

Our work provides direct evidence that performance of goal-conditioned RL can be improved when the representations of the goals are both *discrete* and *factorial*. We show that an instantiation of this idea using multi-factor discretization significantly improves performance on a diverse set of benchmarks. An interesting question that arises from our work is how to theoretically *ground* and *specify* goals, which might be helpful for efficient structured exploration in tasks where goal seeking is crucial.

## 307 **References**

- [1] Ahmed Akakzia, Cédric Colas, Pierre-Yves Oudeyer, Mohamed Chetouani, and Olivier
   Sigaud. Grounding language to autonomously-acquired skills via goal generation. *arXiv preprint arXiv:2006.07185*, 2020.
- [2] Ankesh Anand, Evan Racah, Sherjil Ozair, Yoshua Bengio, Marc-Alexandre Côté, and R. Devon Hjelm. Unsupervised state representation learning in atari. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 8766–8779, 2019.
- [3] Dzmitry Bahdanau, Felix Hill, Jan Leike, Edward Hughes, Arian Hosseini, Pushmeet Kohli,
   and Edward Grefenstette. Learning to understand goal specifications by modelling reward.
   *arXiv preprint arXiv:1806.01946*, 2018.
- [4] Andrew G Barto and Sridhar Mahadevan. Recent advances in hierarchical reinforcement learning. *Discrete event dynamic systems*, 13(1):41–77, 2003.
- [5] Andres Campero, Roberta Raileanu, Heinrich Küttler, Joshua B. Tenenbaum, Tim Rocktäschel,
   and Edward Grefenstette. Learning with amigo: Adversarially motivated intrinsic goals. In
   9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria,
   May 3-7, 2021. OpenReview.net, 2021.
- [6] Elliot Chane-Sane, Cordelia Schmid, and Ivan Laptev. Goal-conditioned reinforcement learn ing with imagined subgoals. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*,
   volume 139 of *Proceedings of Machine Learning Research*, pages 1430–1440. PMLR, 2021.
- [7] Crystal Chao, Maya Cakmak, and Andrea L Thomaz. Towards grounding concepts for transfer
   in goal learning from demonstration. In *2011 IEEE International Conference on Development and Learning (ICDL)*, volume 2, pages 1–6. IEEE, 2011.
- [8] Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Sa haria, Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample effi ciency of grounded language learning. *arXiv preprint arXiv:1810.08272*, 2018.
- [9] Maxime Chevalier-Boisvert, Lucas Willems, and Suman Pal. Minimalistic gridworld environment for openai gym. https://github.com/maximecb/gym-minigrid, 2018.
- [10] Ugo Dal Lago, Marco Pistore, and Paolo Traverso. Planning with a language for extended
   goals. In *AAAI/IAAI*, pages 447–454, 2002.
- [11] Peter Dayan and Geoffrey E Hinton. Feudal reinforcement learning. Advances in neural
   *information processing systems*, 5, 1992.
- [12] Peter Dayan and Geoffrey E. Hinton. Feudal reinforcement learning. In Stephen Jose Hanson,
  Jack D. Cowan, and C. Lee Giles, editors, *Advances in Neural Information Processing Systems*5, *[NIPS Conference, Denver, Colorado, USA, November 30 December 3, 1992]*, pages 271–
  278. Morgan Kaufmann, 1992.
- [13] Thomas G. Dietterich. Hierarchical reinforcement learning with the MAXQ value function
   decomposition. J. Artif. Intell. Res., 13:227–303, 2000.
- [14] Zach Dwiel, Madhavun Candadai, Mariano J. Phielipp, and Arjun K. Bansal. Hierarchical
   policy learning is sensitive to goal space design. *CoRR*, abs/1905.01537, 2019.

- [15] Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. C-learning: Learning to
   achieve goals via recursive classification. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021.* OpenReview.net, 2021.
- Inverse reinforcement learning for vision-based instruction following. *arXiv preprint arXiv:1902.07742*, 2019.
- [17] Anirudh Goyal, Riashat Islam, Daniel Strouse, Zafarali Ahmed, Hugo Larochelle, Matthew M.
   Botvinick, Yoshua Bengio, and Sergey Levine. Infobot: Transfer and exploration via the
   information bottleneck. In *7th International Conference on Learning Representations, ICLR* 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.
- [18] Herbert P Grice. Logic and conversation. In Speech acts, pages 41–58. Brill, 1975.
- [19] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR, 2018.
- [20] Robert I Jennrich. Asymptotic properties of non-linear least squares estimators. *The Annals of Mathematical Statistics*, 40(2):633–643, 1969.
- Yiding Jiang, Shixiang Shane Gu, Kevin P Murphy, and Chelsea Finn. Language as an abstraction for hierarchical deep reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- [22] Leslie Pack Kaelbling. Learning to achieve goals. In *IJCAI*, volume 2, pages 1094–8. Citeseer,
   1993.
- [23] Leslie Pack Kaelbling. Learning to achieve goals. In Ruzena Bajcsy, editor, *Proceedings of the 13th International Joint Conference on Artificial Intelligence. Chambéry, France, August 28 September 3, 1993*, pages 1094–1099. Morgan Kaufmann, 1993.
- [24] Leslie Pack Kaelbling et al. An architecture for intelligent reactive systems. *Reasoning about actions and plans*, pages 395–410, 1987.
- [25] Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A
   survey. *Journal of artificial intelligence research*, 4:237–285, 1996.
- [26] Khimya Khetarpal, Zafarali Ahmed, Gheorghe Comanici, David Abel, and Doina Precup.
   What can i do here? a theory of affordances in reinforcement learning. In *International Con- ference on Machine Learning*, pages 5243–5253. PMLR, 2020.
- [27] Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical
   deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *Advances in neural information processing systems*, 29, 2016.
- [28] 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.
- [29] Andrew Levy, George Dimitri Konidaris, Robert Platt Jr., and Kate Saenko. Learning multi level hierarchies with hindsight. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019.* OpenReview.net, 2019.
- [30] Dianbo Liu, Alex M Lamb, Kenji Kawaguchi, Anirudh Goyal, Chen Sun, Michael C Mozer,
   and Yoshua Bengio. Discrete-valued neural communication. *Advances in Neural Information Processing Systems*, 34, 2021.

- [31] Bogdan Mazoure, Remi Tachet des Combes, Thang Long Doan, Philip Bachman, and R Devon
   Hjelm. Deep reinforcement and infomax learning. *Advances in Neural Information Processing Systems*, 33:3686–3698, 2020.
- [32] Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state
   abstraction and provably efficient rich-observation reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 6961–6971. PMLR,
   2020.
- [33] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
   Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [34] Andrew W Moore and Christopher G Atkeson. Prioritized sweeping: Reinforcement learning
   with less data and less time. *Machine learning*, 13(1):103–130, 1993.
- [35] Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman,
  Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Process- ing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS*2018, December 3-8, 2018, Montréal, Canada, pages 3307–3317, 2018.
- [36] Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Near-optimal representation
   learning for hierarchical reinforcement learning. *arXiv preprint arXiv:1810.01257*, 2018.
- [37] Ofir Nachum, Haoran Tang, Xingyu Lu, Shixiang Gu, Honglak Lee, and Sergey Levine. Why
  does hierarchy (sometimes) work so well in reinforcement learning? *CoRR*, abs/1909.10618,
  2019.
- [38] Suraj Nair and Chelsea Finn. Hierarchical foresight: Self-supervised learning of long-horizon
   tasks via visual subgoal generation. In *8th International Conference on Learning Representa- tions, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020.* OpenReview.net, 2020.
- [39] Soroush Nasiriany, Vitchyr Pong, Steven Lin, and Sergey Levine. Planning with goal conditioned policies. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence
   d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 14814–14825, 2019.
- [40] 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.
- [41] Vitchyr Pong, Shixiang Gu, Murtaza Dalal, and Sergey Levine. Temporal difference models:
   Model-free deep rl for model-based control. *arXiv preprint arXiv:1802.09081*, 2018.
- [42] Vitchyr Pong, Shixiang Gu, Murtaza Dalal, and Sergey Levine. Temporal difference models:
   Model-free deep RL for model-based control. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- [43] Roberta Raileanu and Tim Rocktäschel. RIDE: rewarding impact-driven exploration for
   procedurally-generated environments. In 8th International Conference on Learning Repre sentations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020.
- [44] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images
   with vq-vae-2. Advances in neural information processing systems, 32, 2019.

- [45] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approx imators. In *International conference on machine learning*, pages 1312–1320. PMLR, 2015.
- [46] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In Francis R. Bach and David M. Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 1312–1320. JMLR.org, 2015.
- [47] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust
   region policy optimization. In *International conference on machine learning*, pages 1889–
   1897. PMLR, 2015.
- [48] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
   policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [49] Wolfram Schultz. Predictive reward signal of dopamine neurons. *Journal of neurophysiology*, 80(1):1–27, 1998.
- [50] Max Schwarzer, Ankesh Anand, Rishab Goel, R Devon Hjelm, Aaron Courville, and Philip
   Bachman. Data-efficient reinforcement learning with self-predictive representations. *arXiv preprint arXiv:2007.05929*, 2020.
- [51] Max Schwarzer, Ankesh Anand, Rishab Goel, R. Devon Hjelm, Aaron C. Courville, and Philip
   Bachman. Data-efficient reinforcement learning with self-predictive representations. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021.* OpenReview.net, 2021.
- [52] Max Schwarzer, Nitarshan Rajkumar, Michael Noukhovitch, Ankesh Anand, Laurent Charlin,
  R. Devon Hjelm, Philip Bachman, and Aaron C. Courville. Pretraining representations for
  data-efficient reinforcement learning. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N.
  Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Informa- tion Processing Systems 34: Annual Conference on Neural Information Processing Systems*2021, NeurIPS 2021, December 6-14, 2021, virtual, pages 12686–12699, 2021.
- [53] David Silver, Satinder Singh, Doina Precup, and Richard S Sutton. Reward is enough. *Artificial Intelligence*, 299:103535, 2021.
- [54] Satinder P Singh and Richard S Sutton. Reinforcement learning with replacing eligibility
   traces. *Machine learning*, 22(1):123–158, 1996.
- [55] 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.
- [56] Sainbayar Sukhbaatar, Emily Denton, Arthur Szlam, and Rob Fergus. Learning goal embed dings via self-play for hierarchical reinforcement learning. *CoRR*, abs/1811.09083, 2018.
- 474 [57] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press,
   475 2018.
- [58] Richard S Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M Pilarski, Adam
   White, and Doina Precup. Horde: A scalable real-time architecture for learning knowledge
   from unsupervised sensorimotor interaction. In *The 10th International Conference on Au- tonomous Agents and Multiagent Systems-Volume 2*, pages 761–768, 2011.
- [59] Richard Stuart Sutton. *Temporal credit assignment in reinforcement learning*. PhD thesis,
   University of Massachusetts Amherst, 1984.
- [60] Gerald Tesauro. Practical issues in temporal difference learning. *Advances in neural informa- tion processing systems*, 4, 1991.

- [61] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- [62] Aad W. van der Vaart and Jon A. Wellner. Weak Convergence and Empirical Processes.
   Springer New York, 1996.
- [63] Vivek Veeriah, Junhyuk Oh, and Satinder Singh. Many-goals reinforcement learning. *CoRR*,
   abs/1806.09605, 2018.
- 490 [64] Marco A. Wiering and Jürgen Schmidhuber. Hq-learning. Adapt. Behav., 6(2):219–246, 1997.
- [65] Mengjiao Yang and Ofir Nachum. Representation matters: Offline pretraining for sequential
   decision making. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th Interna- tional Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume
   139 of *Proceedings of Machine Learning Research*, pages 11784–11794. PMLR, 2021.
- [66] Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Reinforcement learning with
   prototypical representations. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*,
   volume 139 of *Proceedings of Machine Learning Research*, pages 11920–11931. PMLR, 2021.
- [67] Lunjun Zhang, Ge Yang, and Bradly C. Stadie. World model as a graph: Learning latent
   landmarks for planning. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*,
   volume 139 of *Proceedings of Machine Learning Research*, pages 12611–12620. PMLR, 2021.
- [68] Tianren Zhang, Shangqi Guo, Tian Tan, Xiaolin Hu, and Feng Chen. Generating
   adjacency-constrained subgoals in hierarchical reinforcement learning. In Hugo Larochelle,
   Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Informa- tion Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.*

# 508 Checklist

- The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or [N/A]. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:
- Did you include the license to the code and datasets? [Yes] See Section ??.
- Did you include the license to the code and datasets? [No] The code and the data are proprietary.
- Did you include the license to the code and datasets? [N/A]
- Please do not modify the questions and only use the provided macros for your answers. Note that the
  Checklist section does not count towards the page limit. In your paper, please delete this instructions
  block and only keep the Checklist section heading above along with the questions/answers below.
- 520 1. For all authors...

521	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's
522	contributions and scope? [Yes]. Main claims are justified with a set of illustrative
523	examples and experiments; both standard and toy examples included for experiments
524	to justify the key claims.

(b) Did you describe the limitations of your work? [Yes]

526 527 528 529 530 531 532		(c) (d)	Did you discuss any potential negative societal impacts of your work? [N/A] Besides the general impact of machine learning methods we do not foresee and identify any particular and immediate negative societal impact of our work. We do not conduct any experiments that could be in itself harmful for society but understand that the applicability of this method is very broad. Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
533	2.	If yo	u are including theoretical results
534		(a)	Did you state the full set of assumptions of all theoretical results? [Yes]
535		(b)	Did you include complete proofs of all theoretical results? [Yes]
536	3.	If yo	u ran experiments
537 538 539		(a)	Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [Yes]. Code will be provided along with the supplementary materials for reproducibility.
540 541		(b)	Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] . More details to be included in supplemetary section
542 543		(c)	Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]. All experiments are run across 3 to 5 random seeds.
544 545 546		(d)	Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]. This will be mentioned in supplementary and acknowledgements
547	4.	If yo	u are using existing assets (e.g., code, data, models) or curating/releasing new assets
548		(a)	If your work uses existing assets, did you cite the creators? [N/A]
549		(b)	Did you mention the license of the assets? [N/A]
550 551		(c)	Did you include any new assets either in the supplemental material or as a URL? [N/A]
552 553		(d)	Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
554 555		(e)	Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? $[N/A]$
556	5.	If yo	u used crowdsourcing or conducted research with human subjects
557 558		(a)	Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
559 560		(b)	Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
561 562		(c)	Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]