

# Zero-Shot Policy Transfer with Disentangled Task Representation of Meta-Reinforcement Learning

Zheng Wu<sup>1,\*</sup>, Yichen Xie<sup>1,\*</sup>, Wenzhao Lian<sup>2</sup>, Changhao Wang<sup>1</sup>, Yanjiang Guo<sup>3</sup>,  
Jianyu Chen<sup>3</sup>, Stefan Schaal<sup>2</sup> and Masayoshi Tomizuka<sup>1</sup>

**Abstract**—Humans are capable of abstracting various tasks as different combinations of multiple attributes. This perspective of compositionality is vital for human rapid learning and adaption since previous experiences from related tasks can be combined to generalize across novel compositional settings. In this work, we aim to achieve zero-shot policy generalization of Reinforcement Learning (RL) agents by leveraging the task compositionality. Our proposed method is a meta-RL algorithm with disentangled task representation, explicitly encoding different aspects of the tasks. Policy generalization is then performed by inferring unseen compositional task representations via the obtained disentanglement without extra exploration. The evaluation is conducted on three simulated tasks and a challenging real-world robotic insertion task. Experimental results demonstrate that our proposed method achieves policy generalization to unseen compositional tasks in a zero-shot manner.

## I. INTRODUCTION

Policy transfer has been a long-standing challenge in the robotics community. Reasoning over existing task experiences and transferring knowledge to novel combinatorial tasks with familiar elements is vital for human rapid learning and adaption. For instance, the images in Figure 1(a) can be abstracted as combinations of (*digit, color*), and humans can effortlessly imagine other unseen combinatorial images. In this work, we aim to enable similar ability in the context of reinforcement learning (RL). Specifically, we consider the generalization scenarios where tasks-of-interest can be described by a set of degrees of variation (DoVs) and our goal is to achieve policy transfer across unseen compositional tasks, namely, tasks of unseen combinations of existing DoVs. Considering the example depicted in Figure 1(b), the robot policy is influenced not only by the *goal* that can be directly measured, but also by other *environment*-related properties of the system, such as robot model, friction, control gain, etc. When facing an unseen compositional task, e.g., (*real, front*) in Figure 1(b), we aim to combine the knowledge obtained from other related tasks, e.g., (*real, left*), (*sim1, front*), to directly get the policy for the new task.

The idea of policy transfer across novel task combinations of known DoVs is also explored in [1], where the authors proposed a modular policy architecture to decompose the policy into goal-specific and robot-specific modules and demonstrated zero-shot policy transfer to unseen robot-goal combinations using existing modules. Specifically, the

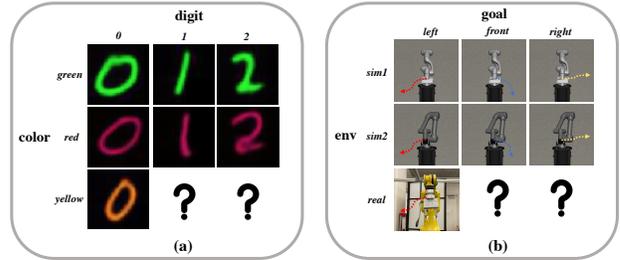


Fig. 1: Illustration of unseen compositional tasks with existing degrees of variation. (a) While (*yellow, 1*) image is unseen, both *yellow* and *1* are present in other images. (b) While (*real, front*) is a novel task, both *real* and *front* are present in other tasks.

authors decompose the observations into goal-related and robot-related observations, which are fed as input to the corresponding neural network modules to output desired actions. However, this method requires training an individual neural network based module for each possible value of DoV, making it computationally expensive to scale up when certain DoV have many possible values. For example, given 20 different possible goals and 10 different possible robots, this method needs to train 30 different modules in total. Additionally, one implicit assumption of [1] is that observation alone can capture the differences between tasks with different DoV combinations. However, there exist certain scenarios where the available observations fail to reflect the differences between tasks. For example, when learning an RL agent to push objects of different weights using image as observation, the dynamics variations between tasks are not captured by observation  $o_t$ , but instead by observed transitions  $(o_t, a_t, o_{t+1}, r_t)$ . Solely performing task decomposition *w.r.t.* the observation space limits the potential to be applied in such scenarios.

In contrast, the key insight of our work is to achieve task decomposition from the perspective of task representation instead of policy architecture as in [1]. Specifically, we build upon the framework of context-based meta-RL framework [2], [3], [4], [5], [6] which jointly learns a latent task embedding from observed transitions and a policy conditioned on the inferred task representation. The task decomposition is achieved by disentangling latent task embedding as the concatenation of embeddings for each task DoV. For this purpose, we develop a task disentanglement regularization objective for meta-training (Section III-C). The

\* The authors are equally contributed.

<sup>1</sup>University of California, Berkeley, Berkeley, CA, USA. <sup>2</sup>Intrinsic Innovation LLC, Mountain View, CA, USA. <sup>3</sup>Tsinghua University, Beijing, China.

achieved disentanglement in the latent task space allows us to infer novel task representations without extra experiences by combining the existing representations of task DoVs in the unseen tasks, thus achieving zero-shot policy generalization (Section III-D).

To evaluate our approach, we extend three traditional meta-RL tasks to compositional tasks composed of certain DoVs and test the generalization performance of our method on the unseen compositional tasks. Experimental results indicate that the generalized policies obtained from our method in a zero-shot manner outperform other baselines. We also demonstrate that our proposed approach can be seamlessly applied as a zero-shot sim-to-real generalization method. We evaluate the sim-to-real algorithm on challenging real-world tilted peg-in-hole tasks. Results demonstrate the effectiveness of our approach.

## II. RELATED WORK

### A. RL Policy Generalization across Similar Tasks

Policy generalization has been a long-standing challenge in robotics community. While RL [7] achieves promising results on robotic manipulation tasks [8], [9], the impedance of deploying RL algorithms lies at the numerous data required during training. To this end, many efforts have been devoted to generalizing the existing policies to similar tasks. A common application of policy generalization is sim-to-real [10], where a policy is learned in simulation and then transferred to the real world. This is usually achieved by domain randomization, which randomizes the tasks with a wide task distribution in simulation and assumes the real-world task is captured by this distribution [11], [12], [13]. Meta-RL approaches [3], [4], [6], [14], [15], [16], [17] offer another perspective of policy generalization by learning to learn from a distribution of tasks, enabling the agent to quickly adapt to novel tasks with limited explorations. Context-based meta-RL methods [3], [4], which view meta-RL as task inference and learn a hidden task variable from collected experience, have demonstrated the ability to adapt meta-policy to real-world manipulation tasks [4], [18]. However, those methods still require collecting online rollout data for task inference and are therefore not zero-shot. Zero-shot policy transfer is particularly useful in circumstances when online interactions are expensive; for example, the task objects are fragile such as the thru-hole components in PCB assembly. In this work, we improve context-based meta-RL methods for policy transfer in a zero-shot manner by considering the task decomposition in the latent task representation.

### B. Task Decomposition in Robotics

Task decomposition is widely investigated in the robotics field to ease the learning process and reuse the knowledge from similar tasks. Most of the research efforts focus on decomposing long-horizon manipulation tasks temporally into several sub-tasks and learning sub-policy for each sub-task [19], [20], [21], [22], [23], [24]. The learned sub-policies that can be re-arranged in a novel way to obtain the policy

for unseen long-horizon tasks. There are also some works that consider task decomposition from other perspectives. For example, [25] performs task decomposition by leveraging the correspondence in the provided analogy and decompose a task as a (action, object) tuple. [26] decomposes the task in the feature space as principle and non-principle features. [27] reasons task decomposition from the object-centric way. The task decomposition considered in our work is most similar to [1], in which tasks are decomposed by a set of degrees of variation. By leveraging knowledge from different tasks with common degrees of variation, we aim to directly obtain policies for novel combinatorial tasks with existing degrees of variation. However, unlike [1] that performs task decomposition through a modular policy architecture, we achieve task decomposition in the context of meta-RL by disentangling the latent task representation.

## III. OUR PROPOSED APPROACH

In this section, we first give a brief review of prior works on context-based meta-RL algorithms in Section III-A. This provides the background to formalize our problem setup in Section III-B. Afterwards, we describe our method in two parts: how to enable task decomposition in the latent task representation during training in Section III-C, and zero-shot policy generalization leveraging the achieved decomposition in Section III-D.

### A. Preliminary: Context-based Meta-RL

Meta-RL typically assumes a distribution of tasks  $p(\mathcal{T})$ . Each task  $\mathcal{T}$  is modeled as an independent Markov decision process (MDP),  $\mathcal{M} \equiv (\mathcal{S}, \mathcal{A}, \mathcal{P}, R)$ , where  $\mathcal{S}$  corresponds to the state space,  $\mathcal{A}$  corresponds to the action space,  $\mathcal{P}$  denotes the transition probability, and  $R$  represents the real-valued reward. Meta-RL algorithms learn the task-conditioning policy from training tasks sampled from  $p(\mathcal{T})$ , and adapt the learned policy to new tasks. Specifically, we focus on one perspective of meta-RL, *i.e.* context-based meta-RL, which views meta-learning as task inference [3], [4]. These methods, with an encoder  $q_\phi$ , map each task  $\mathcal{T}_i$  into a task embedding  $\mathbf{z}_i$  based on the past experience on this task, *e.g.* context  $\mathbf{c}_i = \{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}'_t, r_t)\}_{t=1\dots N}$  [3]. The policy  $\pi_\theta$  is conditioned on the task embedding so that it can be adapted to different tasks, *i.e.*,

$$\mathbf{a} \sim \pi_\theta(\cdot | \mathbf{s}, \mathbf{z}_i), \quad \text{where } \mathbf{z}_i \sim q_\phi(\cdot | \mathbf{c}_i). \quad (1)$$

While our method is agnostic to the specific context-based meta-RL algorithms, in our implementation, we build on the framework of probabilistic embeddings for actor-critic RL (PEARL) [3], which optimizes the objective,

$$\mathbb{E}_{\mathcal{T}_i} \left[ \mathbb{E}_{\mathbf{z}_i \sim q_\phi(\mathbf{z}_i | \mathbf{c}_i)} \left[ R(\mathcal{T}_i, \mathbf{z}_i) + \beta D_{\text{KL}}(q_\phi(\mathbf{z}_i | \mathbf{c}_i) \| p(\mathbf{z}_i)) \right] \right], \quad (2)$$

where  $R(\mathcal{T}_i, \mathbf{z}_i)$  denotes the addition of actor and critic objectives as defined in [3]. We refer the readers [3] for more technical details.

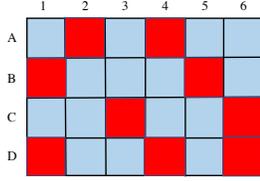


Fig. 2: Illustration of the targeted policy generalization setting in our work. We assume the tasks-of-interest can be described as combinations of degrees of variation (DoV). Given experience from training task set  $\mathbf{T}_{train}$  (blue), we aim to generalize across tasks (in red) whose DoV combinations are unseen but each DoV label is present in  $\mathbf{T}_{train}$ .

### B. Problem Formulation

Our goal is to achieve zero-shot policy transfer across unseen compositional tasks. We define compositional task as the task equipped with a predefined set of degrees of variation (DoVs) [1]. Loosely speaking, a task DoV describes one varied aspect of the task. For example, in Fig. 1(b), there are two DoVs describing the tasks, *i.e.* *environment* and *goal*, and the DoV *goal* has three unique values, *i.e.* *left* (red trajectory), *front* (blue trajectory), and *right* (yellow trajectory). Formally, considering the task distribution  $p(\mathcal{T})$ , each  $\mathcal{T}_i$  is equipped with a unique combination  $\mathcal{C}_i$  of  $M$  task DoVs as  $\mathcal{C}_i = (y_1^i, \dots, y_M^i)$ , where  $y_j^i$  is the label of the  $j$ -th DoV of task  $\mathcal{T}_i$ . The  $M$  DoVs are problem-dependent and identical for all the tasks  $\mathcal{T}_i$  in the task set  $\mathbf{T}$ , while each DoV can have different labels and the task DoV combination  $\mathcal{C}_i$  for each task  $\mathcal{T}_i$  is unique. We assume the DoV labels,  $\mathcal{C}_i$ , of each task  $\mathcal{T}_i$  in  $\mathbf{T}$  are given. Our problem setting aims at generalization across tasks whose DoV combinations are unseen but the DoV labels are present in training tasks  $\mathbf{T}_{train}$ , as illustrated in Figure 2.

### C. Task Decomposition via Latent Task Embedding Disentanglement

In this section, we describe how we achieve task decomposition *w.r.t.* task DoV combinations in the context of meta-reinforcement learning (Meta-RL). Figure 3 illustrates our framework. The key insight is to achieve task decomposition by enforcing the disentanglement of the latent task representation in [3] during meta-training. While traditional meta-RL algorithms learn the task embedding  $\mathbf{z}_i$  with a single encoder  $q_\phi$  as Eq. 1, we replace it with  $M$  separate DoV encoders  $q_{\phi_1}, q_{\phi_2}, \dots, q_{\phi_M}$ , where  $q_{\phi_j}(\mathbf{z}_i^{y_j} | \mathbf{c}_i)$  learns the DoV embedding for the  $j$ -th DoV of task  $\mathcal{T}_i$ . The task embedding is obtained by concatenating all  $M$  DoV embeddings as  $\mathbf{z}_i = [\mathbf{z}_i^{y_1}, \dots, \mathbf{z}_i^{y_M}]$ , and the policy  $\pi_\theta$  is conditioned on the disentangled task embedding:

$$\mathbf{a} \sim \pi_\theta(\cdot | \mathbf{s}, \mathbf{z}_i) \quad \text{where } \mathbf{z}_i = [\mathbf{z}_i^{y_1}, \dots, \mathbf{z}_i^{y_M}] \quad (3)$$

$$\mathbf{z}_i^{y_j} \sim q_{\phi_j}(\cdot | \mathbf{c}_i), j = 1, 2, \dots, M.$$

For different tasks  $\mathcal{T}_u, \mathcal{T}_v$  sharing the same label in the  $j$ -th DoV, *i.e.*,  $y_j^u = y_j^v = y_j$ , their corresponding DoV embedding should be regularized to follow an identical probability distribution. One naive way to implement this

principle is to bring close the corresponding DoV embedding of different tasks with shared DoV labels within the sampled data batch. However, we found empirically that this strategy is easily prone to the affection of noise, making the training process unstable. To this end, we introduce a target distribution  $p(\mathbf{z}_{target}^{y_j})$  for the  $j$ -th task DoV with  $y_j$  label, and minimize the Kullback–Leibler (KL) divergence between  $q_{\phi_j}(\mathbf{z}_i^{y_j} | \mathbf{c}_i)$  and  $p(\mathbf{z}_{target}^{y_j})$ . Inspired by [28], [29], Exponential Moving Average (EMA) is incorporated into our training scheme to update this target distribution during training, achieving robustness against outliers during training. For each DoV label  $y_j$  of the  $j$ -th DoV, we maintain an individual target distribution  $p(\mathbf{z}_{target}^{y_j})$ . Following PEARL [3], our DoV embedding is modeled as Gaussian distribution. So we maintain the EMA for both the mean  $\bar{\boldsymbol{\mu}}_j^{y_j}$  and variance  $\bar{\boldsymbol{\Sigma}}_j^{y_j}$  of  $p(\mathbf{z}_{target}^{y_j})$ . This task disentanglement regularization is written as follows for each task  $\mathcal{T}_i$ , where  $\lambda$  is a hyperparameter to balance different loss items.

$$\mathcal{L}_{dis}^i = \lambda \sum_{j=1}^M D_{KL}(q_{\phi_j}(\mathbf{z}_i^{y_j} | \mathbf{c}_i) || p(\mathbf{z}_{target}^{y_j})) \quad (4)$$

$$p(\mathbf{z}_{target}^{y_j}) = \mathcal{N}(\bar{\boldsymbol{\mu}}_j^{y_j}, \bar{\boldsymbol{\Sigma}}_j^{y_j})$$

where  $q_{\phi_j}(\mathbf{z}_i^{y_j} | \mathbf{c}_i) = \mathcal{N}(\boldsymbol{\mu}_i^{y_j}, \boldsymbol{\Sigma}_i^{y_j})$  is the current prediction for task  $\mathcal{T}_i$ .  $\bar{\boldsymbol{\mu}}_j^{y_j}$  and  $\bar{\boldsymbol{\Sigma}}_j^{y_j}$  are updated in each training step  $t$  with EMA:

$$\begin{aligned} \bar{\boldsymbol{\mu}}_j^{y_j, t} &\leftarrow \tau \cdot \bar{\boldsymbol{\mu}}_j^{y_j, t-1} + (1 - \tau) \cdot \boldsymbol{\mu}_i^{y_j} \\ \bar{\boldsymbol{\Sigma}}_j^{y_j, t} &\leftarrow \tau \cdot \bar{\boldsymbol{\Sigma}}_j^{y_j, t-1} + (1 - \tau) \cdot \boldsymbol{\Sigma}_i^{y_j}, \end{aligned} \quad (5)$$

where  $\tau = 0.99$  in our implementation. The task disentanglement regularization is applied during meta-training, as summarized in Algorithm 1, where red lines highlight the differences compared with PEARL.

### D. Zero-Shot Policy Generalization through Task Decomposition

When generalizing policies to novel tasks, traditional meta-RL algorithms usually requires adequate explorations on the new tasks. However, exploration can be costly in some scenarios, *e.g.*, the task objects are fragile such as the thru-hole components in PCB assembly. In contrast, we can achieve zero-shot policy generalization leveraging the task decomposition obtained from Section III-C. We identify two scenarios of zero-shot policy generalization as follows.

a) S1: Test tasks are defined by combinations of DoV labels already seen in training tasks  $\mathbf{T}_{train}$ . In this case, we directly concatenate the EMA  $\{\mathbf{z}_{target}^{y_j}\}_{j=1}^M$  of the corresponding DoV embeddings obtained from training as the test task representation. The meta-policy conditioned on the concatenated task representation is used as the generalized policy for the new task.

b) S2: Test tasks have DoV label that does not exist in the training tasks  $\mathbf{T}_{train}$ , denoted as,  $y_j = y^*$ . Similar to PEARL meta-test, our method first collects context in one single task  $\mathcal{T}_{test}$  with the unseen DoV label and acquires the latent task representation  $\mathbf{z}_{test}$ . Then the embedding of the DoV label,  $\mathbf{z}_{test}^{y_j^*}$ , is inferred as the corresponding dimensions

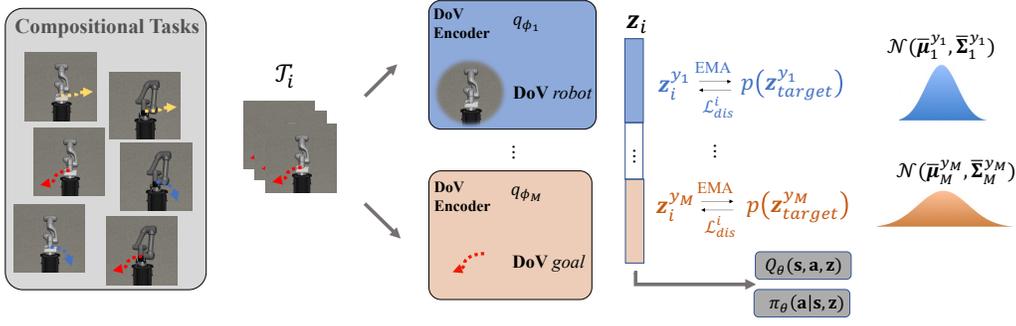


Fig. 3: We achieve task decomposition via disentanglement in the latent task space. The disentanglement is achieved by encoding different task DoVs with individual DoV encoders and enforcing identical task DoVs across different tasks have similar embeddings.

**Algorithm 1** Our meta-training procedure. Differences with PEARL are highlighted in **red**.

**Require:** Batches of compositional training tasks  $\{\mathcal{T}_i\}_{i=1\dots T}$  with corresponding task DoV combinations  $\{\mathcal{C}_i\}_{i=1\dots T}$ , learning rate  $\alpha_1, \alpha_2, \alpha_3$

- 1: Init. replay buffer  $\mathcal{B}_i$  for each training task
- 2: **while** not done **do**
- 3:   **for** each task  $\mathcal{T}_i$  **do**
- 4:     Initialize context  $\mathbf{c}_i = \{\}$
- 5:     **for**  $k = 1, \dots, K$  **do**
- 6:       Sample  $\mathbf{z}_i^{y_j} \sim q_{\phi_j}(\cdot|\mathbf{c}_i), j = 1, \dots, M$
- 7:       Concatenate task embedding  $\mathbf{z}_i = [\mathbf{z}_i^{y_1}, \dots, \mathbf{z}_i^{y_M}]$
- 8:       Gather data from  $\pi_{\theta}(\mathbf{a}|\mathbf{s}, \mathbf{z}_i)$  and add to  $\mathcal{B}_i$  and update  $\mathbf{c}_i = \{(s_t, \mathbf{a}_t, s'_t, \mathbf{r}_t)\}_{t:1\dots N} \sim \mathcal{B}_i$
- 9:     **end for**
- 10:   **end for**
- 11:   **for** step in training steps **do**
- 12:     **for** each  $\mathcal{T}_i$  **do**
- 13:       Sample context batch  $\mathbf{c}_i \sim \mathcal{S}_c(\mathcal{B}_i)$  and RL batch  $b_i \sim \mathcal{B}_i$
- 14:       Sample  $\mathbf{z}_i^{y_j} \sim q_{\phi_j}(\cdot|\mathbf{c}_i), j = 1, \dots, M$
- 15:       Concatenate task embedding  $\mathbf{z}_i = [\mathbf{z}_i^{y_1}, \dots, \mathbf{z}_i^{y_M}]$
- 16:       Compute  $\mathcal{L}_{actor}^i, \mathcal{L}_{critic}^i, \mathcal{L}_{KL}^i$  same as PEARL
- 17:        $\mathcal{L}_{dis}^i = \lambda \sum_{j=1}^M D_{KL}(q_{\phi_j}(\mathbf{z}_i^{y_j}|\mathbf{c}_i)||p(\mathbf{z}_{target}^{y_j}))$ , where  $p(\mathbf{z}_{target}^{y_j}) = \mathcal{N}(\bar{\mu}_j^{y_j}, \bar{\Sigma}_j^{y_j})$
- 18:        $\bar{\mu}_j^{y_j, t} \leftarrow \tau \cdot \bar{\mu}_j^{y_j, t-1} + (1 - \tau) \cdot \mu_j^{y_j}$
- 19:        $\bar{\Sigma}_j^{y_j, t} \leftarrow \tau \cdot \bar{\Sigma}_j^{y_j, t-1} + (1 - \tau) \cdot \Sigma_j^{y_j}$
- 20:     **end for**
- 21:      $\phi_j \leftarrow \phi_j - \alpha_1 \nabla_{\phi_j} \sum_i (\mathcal{L}_{critic}^i + \mathcal{L}_{KL}^i + \mathcal{L}_{dis}^i), j = 1, \dots, M$
- 22:     Update  $\theta_{\pi}, \theta_Q$  same as PEARL
- 23:   **end for**
- 24: **end while**

extracted from  $\mathbf{z}_{test}$ . After obtaining the embedding of this unseen DoV label, we repeat the steps in **S1** to achieve generalization across other test tasks with  $y_j = y^*$ . While

our method does require collecting experiences from one single task in test task set, the generalization across other test tasks is still **zero-shot**. As demonstrated in Section IV-B, this can be readily used as a sim-to-real generalization method. For instance, we meta-train a policy to complete 10 similar manipulation tasks IN different simulation environments. When generalizing to the real world, our method only requires interacting with the real environment on one single task, and is able to generalize to the other 9 tasks in a zero-shot fashion.

## IV. EXPERIMENTAL RESULTS

In the experiments, we aim to investigate the following questions. 1) Does our method described in Section III-C effectively disentangle the learned task representation? 2) How does the disentanglement affect the model training for compositional tasks? 3) Can our method achieve zero-shot policy generalization across novel tasks through task decomposition?

### A. Simulation Experiments

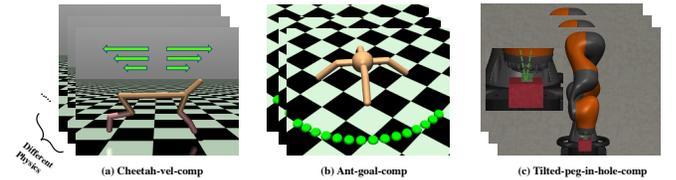


Fig. 4: Three set of compositional tasks for algorithm evaluation. (a) Cheetah-vel-comp with varying *goal velocities* and *physics*. (b) Ant-goal-comp with varying *goal positions* and *physics*. (c) Tilted-peg-in-hole-comp with varying *hole tilting directions* and *physics*.

a) *Compositional Task Setup.*: To study the targeted policy transfer problem in this work, we build three set of simulated compositional tasks shown in Figure 4. Each set of compositional tasks has two task degrees of variation (DoVs), *i.e.*, *goal* and *physics*. For (a) Cheetah-vel-comp and (b) Ant-goal-comp, we extend the traditional Cheetah-vel and Ant-goal from [15] by randomizing some physical parameters, *i.e.*, *body mass*, *body inertia*, *damping* and *friction*.

Specifically, we randomly sample 20 different sets of physical parameters and 20 different goal velocities for Cheetah-vel-comp, leading to 400 compositional tasks in total. Thus each task has a unique (*physics*, *goal*) DoV combination. Similarly for Ant-goal-comp, 15 different groups of physical parameters and 15 different goal 2D locations are randomly sampled. In (c) Tilted-peg-in-hole-comp, we create challenging tilted peg-in-hole tasks based on robosuite [30]. The tasks are performed using a 7-Degree-of-Freedom (DoF) KUKA LBR IIWA robot with operation space control [31]. We adopt a square peg-hole task with  $0.5\text{mm}$  clearance and tilted the hole surface plane  $5^\circ$  towards different directions. Tilted peg-in-hole tasks are prone to jamming issues during execution, thus posing more challenges to robust policy learning. The observation space is 6-dimensional end-effector pose and the actions are the desired end-effector poses. The compositional tasks are composed of 4 different hole tilting directions towards  $\pm x$  and  $\pm y$  axis separately as well as 20 different sets of physical parameters (*i.e.* friction between peg and hole, damping and controller step size).

b) *Training Tasks Selection.*: The training task set  $\mathbf{T}_{train}$  is a subset of the whole compositional task set  $\mathbf{T}$ . When selecting  $\mathbf{T}_{train}$ , we randomly select tasks in  $\mathbf{T}$  with probability  $\alpha \in (0, 1)$ , while ensuring each task DoV label in  $\mathbf{T}$  appear at least once in  $\mathbf{T}_{train}$ . For example, in Figure 2, the 15 tasks in blue are selected for training ( $\alpha = 0.625$ ). In the implementation, we use  $\alpha = 0.5$ , which is discussed in Section IV-C. Other tasks in  $\mathbf{T}$  are set aside as novel test tasks, *i.e.*  $\mathbf{T}_{test} = \mathbf{T} \setminus \mathbf{T}_{train}$ .

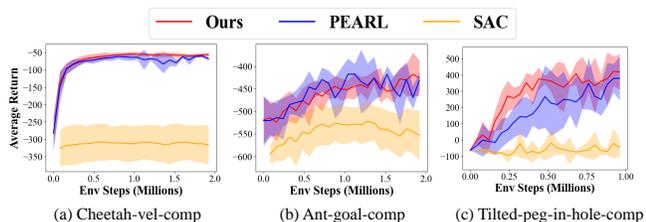


Fig. 5: Average return on test tasks during the training period. Our method and PEARL apply the same meta-test strategy, where the first two trajectories are aggregated into the context for task inference.

c) *Baselines.*: We compare our method against two baselines, PEARL [3] and SAC [32]. PEARL directly learns the task embedding from training tasks without exploiting the inherent task combinatorial structure. During generalization, *i.e.*, meta-test, PEARL requires extra explorations on each test task, so it is not a zero-shot approach. The SAC agent is trained with all the training tasks together where a training task is randomly sampled in each episode. This can be viewed as a variant of domain randomization. The learned SAC policy is directly applied to test tasks. It is worth noting that [1] can not be directly applied as a baseline because solely decomposing tasks *w.r.t.* observations in the three tasks fail to capture the dynamics differences, as discussed in Section I.

d) *Effect on Training.*: Figure 5 shows the performance of the learned policies on  $\mathbf{T}_{test}$  as training proceeds, across 5 random seeds. Our method achieves comparable meta-test results compared to PEARL, indicating that the enforced regularization in latent task representation (Section III-C) does not deteriorate the learning process of meta-RL. It is worth noting that the aim of our proposed method is to improve the performance of zero-shot policy transfer across test tasks, instead of meta-test where the experience of the test tasks is given.

e) *Zero-Shot Evaluation.*: The evaluation is conducted for the two scenarios **S1** and **S2** (Section III-D) respectively. In **S1**,  $\mathbf{T}_{test} = \mathbf{T} \setminus \mathbf{T}_{train}$  is used as the test tasks. While in **S2**, we first sample task attributes that do not exist in  $\mathbf{T}$  and the test tasks are set as the combinations of the unseen attributes and other seen attributes, *e.g.*, (*unseen physics*, *seen goal*), (*unseen goal*, *seen physics*). The experimental results across 5 random seeds are shown in Figure 6. We compare our zero-shot performance (*Zero-shot (ours)*) with multiple baselines: 1) *SAC*: A single model is trained on all training tasks with SAC [20]. 2) *Prior (PEARL)*: The same  $\mathcal{N}(0, 1)$  prior distribution is applied to all the task embeddings  $\mathbf{z}_i$  of PEARL model [3]. 3) *Meta-test (PEARL)*: We apply the same meta-test for PEARL as in [3]. To ensure enough contexts on the test task, the first 8 trajectories are aggregated into the context. 4) *Meta-test (ours)*: The same meta-test with PEARL is applied to the model trained with our method. It is worth mentioning that only 1) and 2) can achieve zero-shot generalization on novel tasks, while 3) and 4) require extra exploration. As shown in Figure 6, our zero-shot generalization significantly outperforms the two zero-shot baselines, *SAC* and *Prior (PEARL)*. Interestingly, it is also superior to the two meta-test baselines, although our generalization does not explore on the test task. We hypothesize that the reason is our method benefits from the EMA embedding  $\mathbf{z}_{target}^{y_j}$  used to compose the embedding  $\mathbf{z}_u$  of the novel task, serving as a temporal ensemble of information accumulated during training. In contrast, purely collecting context from a single test task may suffer from its inherent bias and easily get trapped in the local optimum.

f) *Visualization of Latent Embedding.*: To investigate if the proposed method in Section III-C effectively disentangle the encodings in latent space, Figure 8 visualizes sampled test task representations output by our method and PEARL on Cheetah-vel-comp. The embeddings of test tasks with same *goal* or *physics* are in the same color. All the embeddings are reduced to 2-d with PCA [33] for visualization. It is observed that our method can distinguish different labels of each attribute better than PEARL. Different goals are distinctly separable by our method. For different physics, despite their effect being indirect, similar physics parameters still tend to induce closer embeddings.

## B. Real Experiments

As described in Section III-D, our method can be applied as a sim-to-real algorithm, by treating real-world tasks as tasks with unseen *physics* attribute. We evaluate the method

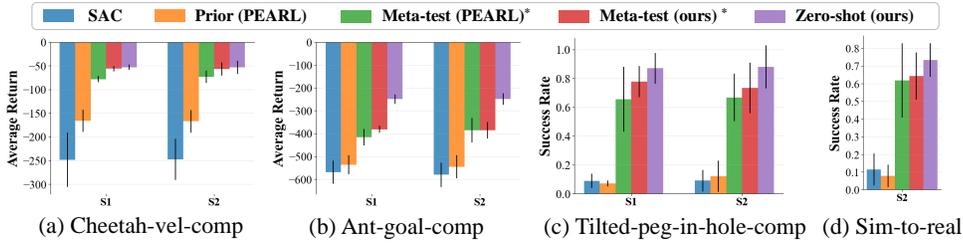


Fig. 6: Evaluation results of policy generalization on test tasks (\* not a zero-shot method).

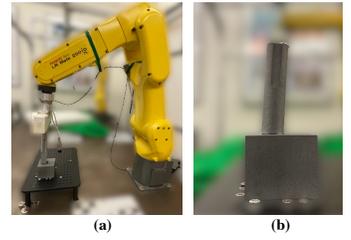


Fig. 7: (a) Real-world setup. (b) 5° tilted hole for experiment.

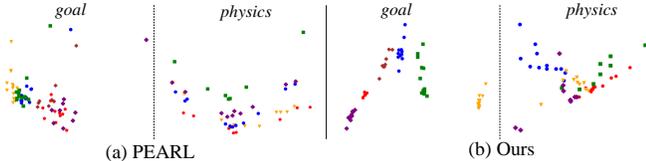


Fig. 8: Visualization of latent task embedding of our method compared to PEARL. The embeddings of tasks with identical *goal* or *physics* labels are in the same color.

on a real-world 6-DoF FANUC Mate 200iD robot with the similar setup as Tilted-peg-in-hole-comp in Section IV-A, as shown in Figure 7. The evaluation of real world is identical to the evaluation of **S2** in simulation (Section IV-A), where, after meta-training in simulation, the meta-policy collects context for one held-out task to infer the real-world physics embedding. The obtained physics embedding is combined with other goal embeddings learned from simulation to generalize across the other three tasks in the real world. Each task is used as the held-out task once for thorough evaluation. When evaluating the generalized policy, the evaluated task is executed 30 trials with 3 different random seeds (10 trials for each random seed). The results are reported in Figure 6(d). Our generalized policy achieves a success rate of 73.3% over different unseen tasks.

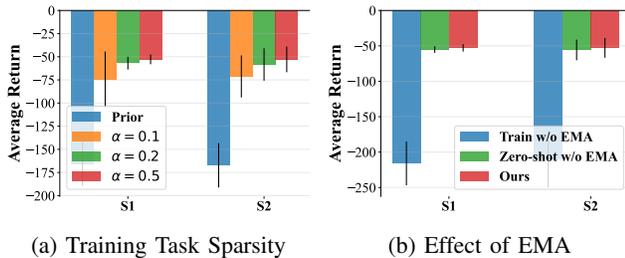


Fig. 9: Results of ablation studies.

### C. Ablation Studies

*a) Sparsity of Training Tasks.*: In our previous experiments, we sample 50% of compositional tasks (*i.e.*  $\alpha = 0.5$ , see Section IV-A (b) for details) for training, so that each task DoV label would appear multiple times during training. To examine the limit of our method, we try to reduce the number of training tasks, *i.e.*  $\alpha = 0.2$  or  $0.1$ , on Cheetah-vel-comp. In this case, each task DoV label would appear

only 4 or 2 times on average among training tasks. Our zero-shot generalization performance is visualized in Figure 9a. Despite the performance decrease with fewer training tasks, our approach still notably outperforms the *Prior* even with the fewest training tasks ( $\alpha = 0.1$ ). This indicates our method learns meaningful disentangled DoV embeddings even with sparse training tasks.

*b) Effect of EMA.*: To justify the role of Exponential Moving Average (EMA) in our training and zero-shot generalization, we consider the following alternative strategies: 1) *training w/o EMA*: Equivalently, we set  $\tau = 0$  in Eq. 5. In this case, the task DoV embedding is regularized by the training task with the same DoV label sampled in the last training batch. 2) *zero-shot generalization w/o EMA*: We try another strategy to get the corresponding DoV embeddings for an unseen task  $\mathcal{T}_u$  instead of using EMA in **S1** of Section III-D. The average DoV embeddings from training tasks with the same task DoV label with  $\mathcal{T}_u$  is adopted instead, *i.e.*  $\mathbf{z}_v^{y_j} = \text{mean}(\{\mathbf{z}_v^{y_j} | \mathcal{T}_v \in \mathbf{T}_{train}, y_j^u = y_j^v\})$ . The results are shown in Figure 9b. Training without EMA fails because of the large noise in the regularization of the DoV embeddings. EMA “smoothes out” the noise, serving as a temporal ensemble of context accumulated during training. In addition, it provides unbiased DoV embeddings for zero-shot generalization, as indicated by the generalization performance gap with and without EMA.

## V. CONCLUSION AND FUTURE WORK

In this paper, we propose a meta-RL algorithm to achieve zero-shot policy generalization across unseen compositional tasks. Our method learns disentangled task representation by encoding each task DoV separately. The DoV embeddings obtained from training tasks enable zero-shot policy generalization on test tasks by composing the task embeddings of novel tasks. Our approach achieves significantly superior performance on simulation tasks and can be seamlessly applied to sim-to-real generalization.

The experiments conducted in this work are mainly the combination of two task DoVs, *i.e.* *physics* and *goal*. In future work, we hope to extend to more diverse attributes, allowing for wider applications of zero-shot policy generalization. Besides, one implicit assumption of modeling task decomposition via disentanglement is the specified task DoVs are independent of each other, limiting the application scenarios. Discovering other ways to achieve task decompo-

sition without this assumption is another interesting direction to explore.

## REFERENCES

- [1] C. Devin, A. Gupta, T. Darrell, P. Abbeel, and S. Levine, "Learning modular neural network policies for multi-task and multi-robot transfer," pp. 2169–2176, 2017.
- [2] C. Finn, K. Xu, and S. Levine, "Probabilistic model-agnostic meta-learning," *Advances in neural information processing systems*, vol. 31, 2018.
- [3] K. Rakelly, A. Zhou, C. Finn, S. Levine, and D. Quillen, "Efficient off-policy meta-reinforcement learning via probabilistic context variables," pp. 5331–5340, 2019.
- [4] T. Z. Zhao, A. Nagabandi, K. Rakelly, C. Finn, and S. Levine, "Meld: Meta-reinforcement learning from images via latent state models," *arXiv preprint arXiv:2010.13957*, 2020.
- [5] J. Humpalik, A. Galashov, L. Hasenclever, P. A. Ortega, Y. W. Teh, and N. Heess, "Meta reinforcement learning as task inference," *arXiv preprint arXiv:1905.06424*, 2019.
- [6] L. Zintgraf, K. Shiarlis, M. Igl, S. Schulze, Y. Gal, K. Hofmann, and S. Whiteson, "Varibad: A very good method for bayes-adaptive deep rl via meta-learning," *arXiv preprint arXiv:1910.08348*, 2019.
- [7] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [8] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke *et al.*, "Scalable deep reinforcement learning for vision-based robotic manipulation," pp. 651–673, 2018.
- [9] M. A. Lee, Y. Zhu, K. Srinivasan, P. Shah, S. Savarese, L. Fei-Fei, A. Garg, and J. Bohg, "Making sense of vision and touch: Self-supervised learning of multimodal representations for contact-rich tasks," pp. 8943–8950, 2019.
- [10] W. Zhao, J. P. Queralta, and T. Westerlund, "Sim-to-real transfer in deep reinforcement learning for robotics: a survey," pp. 737–744, 2020.
- [11] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," pp. 23–30, 2017.
- [12] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," pp. 3803–3810, 2018.
- [13] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel, "Asymmetric actor critic for image-based robot learning," *arXiv preprint arXiv:1710.06542*, 2017.
- [14] Y. Duan, J. Schulman, X. Chen, P. L. Bartlett, I. Sutskever, and P. Abbeel, "RL2: Fast reinforcement learning via slow reinforcement learning," *arXiv preprint arXiv:1611.02779*, 2016.
- [15] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," pp. 1126–1135, 2017.
- [16] A. Gupta, R. Mendonca, Y. Liu, P. Abbeel, and S. Levine, "Meta-reinforcement learning of structured exploration strategies," *Advances in neural information processing systems*, vol. 31, 2018.
- [17] L. Zintgraf, K. Shiarlis, V. Kurin, K. Hofmann, and S. Whiteson, "Fast context adaptation via meta-learning," pp. 7693–7702, 2019.
- [18] G. Schoettler, A. Nair, J. A. Ojea, S. Levine, and E. Solowjow, "Meta-reinforcement learning for robotic industrial insertion tasks," pp. 9728–9735, 2020.
- [19] J. Andreas, D. Klein, and S. Levine, "Modular multitask reinforcement learning with policy sketches," pp. 166–175, 2017.
- [20] T. Haarnoja, V. Pong, A. Zhou, M. Dalal, P. Abbeel, and S. Levine, "Composable deep reinforcement learning for robotic manipulation," pp. 6244–6251, 2018.
- [21] D.-A. Huang, S. Nair, D. Xu, Y. Zhu, A. Garg, L. Fei-Fei, S. Savarese, and J. C. Niebles, "Neural task graphs: Generalizing to unseen tasks from a single video demonstration," pp. 8565–8574, 2019.
- [22] Y. Jiang, S. S. Gu, K. P. Murphy, and C. Finn, "Language as an abstraction for hierarchical deep reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [23] T. Kipf, Y. Li, H. Dai, V. Zambaldi, A. Sanchez-Gonzalez, E. Grefenstette, P. Kohli, and P. Battaglia, "Compile: Compositional imitation learning and execution," pp. 3418–3428, 2019.
- [24] A. Zhou, V. Kumar, C. Finn, and A. Rajeswaran, "Policy architectures for compositional generalization in control," *arXiv preprint arXiv:2203.05960*, 2022.
- [25] J. Oh, S. Singh, H. Lee, and P. Kohli, "Zero-shot task generalization with multi-task deep reinforcement learning," pp. 2661–2670, 2017.
- [26] C. Zhang, H. Zhang, and L. E. Parker, "Feature space decomposition for effective robot adaptation," pp. 441–448, 2015.
- [27] C. Devin, D. Geng, P. Abbeel, T. Darrell, and S. Levine, "Compositional plan vectors," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [28] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," *Advances in neural information processing systems*, vol. 30, 2017.
- [29] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar *et al.*, "Bootstrap your own latent: A new approach to self-supervised learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 271–21 284, 2020.
- [30] Y. Zhu, J. Wong, A. Mandlekar, and R. Martín-Martín, "robosuite: A modular simulation framework and benchmark for robot learning," 2020.
- [31] O. Khatib, "A unified approach for motion and force control of robot manipulators: The operational space formulation," *IEEE Journal on Robotics and Automation*, vol. 3, no. 1, pp. 43–53, 1987.
- [32] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," pp. 1861–1870, 2018.
- [33] H. Abdi and L. J. Williams, "Principal component analysis," *Wiley interdisciplinary reviews: computational statistics*, vol. 2, no. 4, pp. 433–459, 2010.