

# GENERALIZED BEHAVIOR LEARNING FROM DIVERSE DEMONSTRATIONS

Anonymous authors

Paper under double-blind review

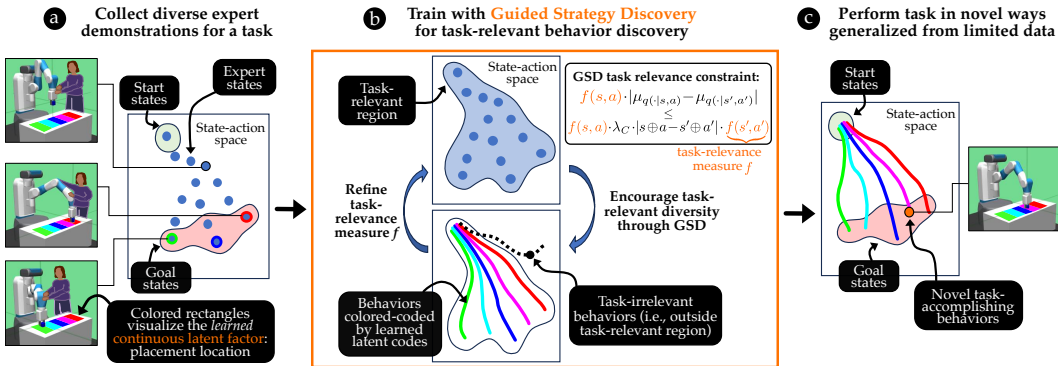


Figure 1: The figure overviews our framework: Guided Strategy Discovery. (a) Given diverse task demonstrations with underlying latent factors, (b) our framework optimizes a task-relevance guided diversity objective, (c) to discover behaviors generalizing to unseen latent factor values.

## ABSTRACT

Diverse behavior policies are valuable in domains requiring quick test-time adaptation or personalized human-robot interaction. Human demonstrations provide rich information regarding task objectives and **factors that govern individual behavior variations**, which can be used to characterize *useful* diversity and learn diverse performant policies. However, we show that prior work that builds naive representations of demonstration heterogeneity fails in generating successful novel behaviors that generalize over **behavior factors**. We propose Guided Strategy Discovery (GSD), which introduces a novel diversity formulation based on a learned task-relevance measure that prioritizes behaviors exploring modeled latent factors. We empirically validate across three continuous control benchmarks for generalizing to in-distribution (interpolation) and out-of-distribution (extrapolation) **factors** that GSD outperforms baselines in novel behavior discovery by  $\sim 21\%$ . Finally, we demonstrate that GSD can generalize striking behaviors for table tennis in a virtual testbed while leveraging human demonstrations collected in the real world.

## 1 INTRODUCTION

Intelligent agents encountered with a novel variation of a learned task should be able to adapt their default decision making to suit the variation at hand. To adapt on-the-fly, agents must learn a concise set of variations to quickly tune their behaviors. Such adaptability is valuable in applications such as few-shot learning (Duan et al., 2017) and personalized human-robot interaction (Wang et al., 2022) where limited examples or interaction must inform compatible approaches for task completion.

Behaviors that adaptable agents learn must be meaningfully diverse such that novel task variations can be addressed sufficiently. In the past, unsupervised reinforcement learning (RL) (Laskin et al., 2021) has been used to learn diverse behaviors or “skills”. However, learned behaviors intended to explore the agent’s environment may not directly be useful towards task completion. Such approaches are further limited by their inability to identify and exhibit meaningful variations that are useful during deployment. While supervised RL can be employed, reward specification to align diverse behaviors with user expectations can be cumbersome (Soares & Fallenstein, 2014).

In contrast to RL, learning from demonstration (LfD) methods enable agents to learn decision-making policies directly from human examples. The distinct **variations** that individuals often **exhibit**, even when pursuing the same task objectives (Sanderson, 1989), **reflect creative ways of task completion**. We assume that these variations<sup>1</sup> are governed by distinct but latent **behavior factors**. These factors impart useful diversity in human behaviors that adaptable agents can exploit.

Prior work in heterogeneous Imitation Learning (IL) (Li et al., 2017; Chi et al., 2023) largely focuses on generating behaviors corresponding to **modes** in training datasets or inferring representations of test behaviors. In this work, we address the challenge of generating behaviors with novel variations. We specifically study the ability to inter- and extrapolate from demonstrations to effectively produce new behaviors, that correspond to **behavior factor values** not seen in the training dataset, while also accomplishing the task. For example, consider a robot quadruped that runs at different speeds. We seek policies that run at 2m/s or 4m/s from demonstrations with speeds of 1m/s and 3m/s. Such a generalization ability can provide task-accomplishing behaviors with desirable characteristics directly through latent prior sampling. **However, generalization with latent behavior factors is challenging**, as we need to accurately identify the latent dimensions along which demonstrations vary and locate individual behaviors in the corresponding space before extending to novel behaviors.

We focus on novel behavior generation in the setting of online IL (Ho & Ermon, 2016) due to its data-sample efficiency. We show that prior approaches that utilize mutual information (MI)-based diversity objectives (Li et al., 2017) fail to produce novel behaviors. We draw inspiration from recent work in unsupervised RL (Park et al., 2022; 2024; 2023) that modify MI-based objectives and structure latent representations in order to induce specific behavioral traits (e.g., high Euclidean or temporal distances between states, controllability, etc). We propose to modify representation learning by restricting the latent space from capturing state-action space regions irrelevant to the task, identified through distillation of **demonstration**-specific occupancy measures. We find that our formulation encourages diversity specifically along traits that vary across demonstrations. We refer to this objective as task-relevant diversity as it produces behaviors that retain task-performance.

We present a novel approach to learn diverse task-accomplishing behaviors from demonstrations that generalize over **latent behavior factors**. Our contributions are four-fold:

- We show the need for a novel formulation of diversity for generalization in IL from diverse demonstrations through experiments in a 2D Point Maze domain (Sec. 4).
- We formulate task-relevant diversity, an objective to encourage diversity along factors of variation among demonstrations by restricting representations from capturing irrelevant regions. We propose Guided Strategy Discovery (GSD), an algorithm that optimizes diversity alongside imitation to discover novel task-accomplishing behaviors (Sec. 5).
- We demonstrate that GSD generalizes to novel behaviors with 21% average error reduction in **behavior factors** (known during evaluation) over four baselines across two splits (interpolation and extrapolation) in three domains spanning robot control, driving, and manipulation (Sec. 6.1).
- We demonstrate that GSD generalizes from physical human demonstrations to capture diverse stroke styles in a simulated Table Tennis domain (Sec. 6.3).

## 2 RELATED WORK

**Generalization in Behavior Learning** Prior works have studied generalization when agents are faced with task specifications from test distributions (Benjamins et al., 2022; Silva et al., 2021; Padalkar et al., 2023; Nair et al., 2022; Shridhar et al., 2023; Xu et al., 2022; Driess et al., 2020), or deployment settings different from training environments (Fu et al., 2017; Kumar et al., 2020; Packer et al., 2018; Kumar et al., 2021; Xie et al., 2023; Cobbe et al., 2019; Osa et al., 2022; Zahavy et al., 2022). In IL, generalization has been considered when demonstrators operate with diverse conditions (Qiu et al., 2023; Tangkaratt et al., 2020; Chen et al., 2021; Paleja et al., 2020; Schrum et al., 2023b; Li et al., 2017; Chen et al., 2020; Wang et al., 2017; Li et al., 2017; Hausman et al., 2017; Peng et al., 2022a). Our work focuses on the latter, where we study heterogeneous demonstrators with latent **behavior factors**. Among these, prior works either attempt to characterize heterogeneity through latent representations (Paleja et al., 2020; Schrum et al., 2023b; Li et al., 2017; Chen et al., 2020), learn performant behaviors from diverse demonstrators (Qiu et al., 2023;

<sup>1</sup>Prior work has attributed **two causes** to demonstration diversity: (i) humans performing tasks in varying sub-optimal ways (Ramachandran & Amir, 2007), and (ii) humans performing tasks optimally in varying styles (Sanderson, 1989). Our work focuses **solely on the latter**.

Tangkaratt et al., 2020; Chen et al., 2021) or simply imitate multiple behaviors (Wang et al., 2017; Li et al., 2017; Hausman et al., 2017; Peng et al., 2022a). To the best of our knowledge, our work is the first to address all three to seek performant behaviors that generalize over **behavior factors**.

**Learning from demonstrations** Prior works in LfD utilize demonstrations to learn rewards (Abbeel & Ng, 2004; Ziebart et al., 2008; Fu et al., 2017; Chen et al., 2020; Ross et al., 2011) or task-performant policies (Ross et al., 2011; Ho & Ermon, 2016; Qiu et al., 2023; Tangkaratt et al., 2020; Chen et al., 2021). **Our work seeks diverse policies, particularly from expert demonstrations with continuous latent factors** (Wang et al., 2017; Li et al., 2017; Hausman et al., 2017; Chen et al., 2020; Peng et al., 2022a). While several works address heterogeneous IL with large datasets Chi et al. (2023), we use environment interaction to tackle covariate shift in low data regimes (Ho & Ermon, 2016; Kostrikov et al., 2018; Reddy et al., 2019; Garg et al., 2021). **Our work belongs to the popular class of adversarial IL (AIL) methods** (Orsini et al., 2021) which model imitation as adversarial game between policies and a discriminator that captures expert occupancy. We specifically study MI-based diversity objectives alongside online adversarial IL (Li et al., 2017; Hausman et al., 2017; Peng et al., 2022a) for learning infinite behaviors with continuous latent spaces. We address the limitations of MI in capturing latent factors specific to demonstrations.

**Diverse behavior learning** Diverse behavior learning has been employed for exploration, pre-training, and generalization to novel environments. **Quality Diversity** (Batra et al., 2024) assumes the availability of task performance metrics and functions for measuring behavior factors. **In contrast, our work is related to unsupervised RL** (Laskin et al., 2021) that learn behaviors without such privileged information. **Among them, we are similar to competence-based methods** (Sharma et al., 2019; Hansen et al., 2019; Park et al., 2022; 2023; 2024) that learn latent spaces to represent heterogeneous behaviors. Works focus on different aspects of diversity, such as state coverage (Eysenbach et al., 2018; Park et al., 2022; Laskin et al., 2022; Mendonca et al., 2021), dynamics (Sharma et al., 2019), controllability (Park et al., 2023; 2024), etc. Our work adopts ideas of regularization (Park et al., 2022; 2023; 2024) for designing diversity objectives to **improve heterogeneous IL**.

**Structured methods for heterogeneous IL** **CASSI** (Li et al., 2023) uses MI-based objectives to learn novel locomotion behaviors from unlabeled data but relies on additional rewards to ensure task completion, unlike our approach which does not depend on rewards. **FLD** (Li et al., 2024) structures latent spaces using differentiable fast Fourier transforms for periodic motions. **In contrast, our approach is domain-agnostic, making it applicable across a broader range of tasks.** **ASE** (Peng et al., 2022b) applies latent sequence modeling combined with hyperspherical priors for smooth motion transitions, and **CALM** (Tessler et al., 2023) builds on ASE with latent-conditioned discriminators. **However, both methods suffer from limitations of naive MI formulations, which our work addresses.**

### 3 PROBLEM STATEMENT AND PRELIMINARIES

We consider an infinite horizon, discounted, and reward-free Markov Decision Process (MDP\(\mathbb{R}\)),  $(S, A, P, \rho_0, \gamma)$ , where  $S$  and  $A$  represent state and action spaces,  $P: S \times A \times S \rightarrow \mathbb{R}$ , the transition probabilities,  $\rho_0: S \rightarrow \mathbb{R}$ , the initial state distribution, and  $\gamma$ , the discount factor. An optimal expert policy  $\pi^\xi$  is governed by a continuous factor,  $\omega \in \Omega$ . The factor space,  $\Omega$ , is split into **disjoint** train and test regions,  $\text{Tr}(\Omega)$  and  $\text{Te}(\Omega)$ , respectively. Given a demonstrations set,  $\mathcal{D}$ , of trajectories  $\tau_i^\xi = \{s_0, a_0, s_1, a_1, \dots\}$ ,  $a_t \sim \pi^\xi(\cdot|s_t, \omega_i)$ ,  $s_{t+1} \sim P(\cdot|s_t, a_t)$ , and  $\omega_i \in \text{Tr}(\Omega)$ , we aim to learn a policy  $\pi$  that captures the expert behavior  $\pi^\xi$  over the entire factor space without access to  $\omega_i$  or  $\Omega$ .

We ground our approach in InfoGAIL (Li et al., 2017), built upon Generative Adversarial Imitation Learning (GAIL) (Ho & Ermon, 2016) to imitate demonstrations:  $J^{\text{GAIL}} := E_{\pi^\xi}[\log D(s, a)] + E_\pi[\log(1 - D(s, a))]$ , where  $D$  is a discriminator that distinguishes between the learned policy,  $\pi$ , and the expert policy,  $\pi^\xi$ . InfoGAIL additionally introduces a latent variable,  $z \in Z$ , to capture **factors** underlying expert demonstrations. InfoGAIL optimizes MI by a variational lower bound (Barber & Agakov, 2004),  $q(z|s, a)$ . We refer to  $q$  as the *decoder* as it infers  $z$  from the state action pair. The InfoGAIL objective is  $J^{\text{InfoGAIL}} := J^{\text{GAIL}} + \lambda_I E_{z, \pi}[\log q(z|s, a)]$ , where  $\lambda_I$  controls the diversity objective weight.

For formulating our diversity objective, we build on ideas from network distillation (Teh et al., 2017; Czarnecki et al., 2019; Chen et al., 2020). MSRD (Chen et al., 2020) learns task rewards from demonstrations,  $\{\zeta_i\}$ , over a finite set of **factors** by employing distillation with AIRL (Fu et al., 2017), a variant of GAIL that recovers a reward function,  $r(s, a)$ . MSRD then distills the reward functions for each **variation**,  $r_{\zeta_i}$ , into a common reward function for the task,  $\tilde{r}_0$ . The

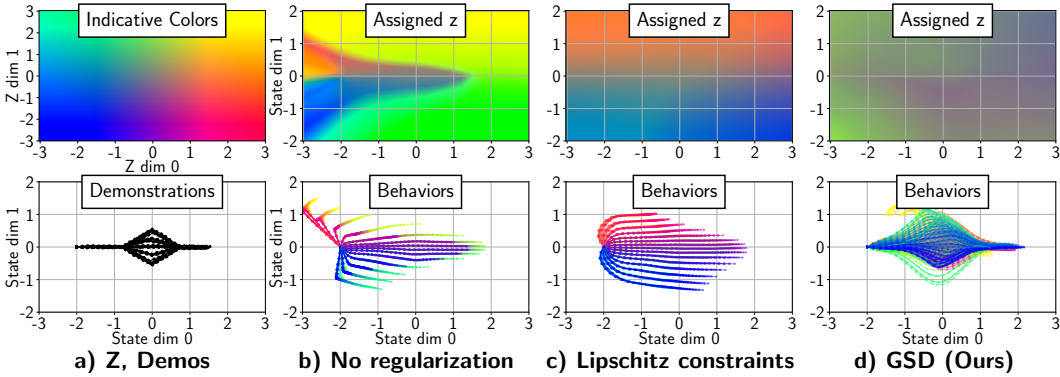


Figure 2: Fig. 2a (Top): Map shows colors assigned to the 2D latent space. Fig. 2a (Bottom): The agent starts at  $(-2, 0)$ , moves to  $(2, 0)$ , passing through  $(-1, 0)$ ,  $(0, \omega)$  and  $(1, 0)$  where  $\omega \in [-1, 1]$ . Figs. 2b, 2c, 2d: Latent vectors assigned to the state space, with a state-only decoder, and policy behaviors, under a high importance weight  $\lambda_I$ , are visualized. Trajectories are colored according to conditioning latent vectors per Fig. 2a (Top). Baselines deviate from demonstrations arbitrarily (Fig. 2b, Bottom) or uniformly (Fig. 2c, Bottom), disregarding the goal. GSD (ours, Fig. 2d) discovers behaviors with novel latent variations (waypoints along  $x = 0$ ) while reaching closer to the goal  $(2, 0)$ .

distillation is done by formulating each reward as,  $r_{\zeta_i}(s, a) := \tilde{r}_0(s, a) + \tilde{r}_{\zeta_i}(s, a)$ . The **factor-specific** residual reward,  $\tilde{r}_{\zeta_i}(s, a)$ , is encouraged to be close to zero with the additional objective,  $J^{\text{MSRD}} := -E_{\zeta_i, \pi}[(\tilde{r}_{\zeta_i}(s, a))^2]$ , to encourage the reward information common across demonstrations pertaining to the task to be represented by the task reward function  $\tilde{r}_0$ .

#### 4 NEED FOR REGULARIZATION

In this section, we show that prior diversity formulations from Li et al. (2017); Park et al. (2022) fail to produce novel, task-accomplishing behaviors, motivating the need for a new formulation.

**No regularization:** InfoGAIL’s diversity objective, MI, promotes diverse behaviors by rewarding the visitation of states associated with distinct latent vectors. In a 2D PointMaze domain with continuous state-action space (see Fig. 2), we show that an increased diversity objective’s weight,  $\lambda_I$ , does not necessarily result in more diverse behaviors that accomplish the task. This result can be attributed to the decoder,  $q$ , a neural network (NN) that assigns latent vectors to state-action pairs,  $\langle s, a \rangle$ . Without regularization, the decoder,  $q$ , produces unconstrained latent assignments, with a high variety in smaller regions (see Fig. 2b, top, several distinct colors close to point  $(-2, 0)$ ). This finding aligns with prior work (Choi et al., 2021; Park et al., 2022). Without regularization, related behaviors with close-by states can be mapped to unrelated far-away regions in the latent space without any meaningful structure. This behavior can cause insufficient (see Fig. 2b, bottom, several trajectories clump together along  $y = 0$ ) or arbitrary (no pattern that governs deviation from demonstrations) behavior diversity.

**Prior regularization methods produce misaligned behaviors:** Prior works in unsupervised RL (Choi et al., 2021; Park et al., 2022) imposed Lipschitz constraints on the decoder, to enforce that for any two state-action pairs,  $\langle s, a \rangle, \langle s', a' \rangle$ , the assigned latent vectors (specifically the mean  $\mu$  of the approximate posterior distribution  $q(\cdot|s, a)$ ), differ by at most the Euclidean distance between the pairs, scaled by  $\lambda_C$ . Formally,  $\|\mu_{q(\cdot|s, a)} - \mu_{q(\cdot|s', a')}\| \leq \lambda_C \cdot \|s \oplus a - s' \oplus a'\|$ , where  $\oplus$  denotes vector concatenation. The Lipschitz constraints ensure smooth latent vector assignments (see Fig. 2c, top), which encourages behaviors to deviate from demonstrations uniformly over the state space. However, resulting behaviors do not necessarily proceed towards the goal (see Fig. 2c, bottom). Other diversity formulations focusing on controllability and temporal reachability (Park et al., 2023; 2024) will face similar issues if the auxiliary objectives are misaligned with behavior heterogeneity. We propose a general formulation that encourages behavior diversity along latent dimensions inferred from the demonstrations, without compromising task performance.

#### 5 OUR METHOD: GUIDED STRATEGY DISCOVERY

We present our approach for achieving generalizable IL from diverse demonstrations.



## 5.1 ENCOURAGING $f$ -RELEVANT DIVERSITY

First, we present a general approach for encouraging diversity selectively within state-action space regions indicated by high energy with respect to a scalar energy function,  $f: S \times A \rightarrow [0, 1]$ .

We design our approach by analysing transitions that occur during learning. Consider a scenario visualized in Fig. 3a, where an exploring agent is at a high  $f$ -energy state-action pair,  $\langle s, a \rangle$ , assigned a latent vector  $z$ , and it enters another pair,  $\langle s', a' \rangle$ . If a different vector,  $z'$  s.t.  $z' \neq z$ , were assigned to  $\langle s', a' \rangle$ , the diversity rewards,  $\log q(z|s, a)$ ,  $\log q(z'|s', a')$ , would encourage behaviors  $\pi(\cdot|\cdot, z)$ ,  $\pi(\cdot|\cdot, z')$ , to visit  $\langle s, a \rangle$  and  $\langle s', a' \rangle$  respectively. The behavior,  $\pi(\cdot|\cdot, z')$ , would be desirable if  $\langle s', a' \rangle$  has high  $f$ -energy, as it would visit a high energy pair different from  $\langle s, a \rangle$ , increasing coverage of high energy regions. On the other hand, if  $\langle s', a' \rangle$  were a low  $f$ -energy pair, the behavior  $\pi(\cdot|\cdot, z')$  visiting a low energy pair would be undesirable. In this case, the assignment for  $\langle s', a' \rangle$  could be constrained close to  $z$ , which would remove the incentive for a behavior distinct from  $\pi(\cdot|\cdot, z)$  to specifically visit  $\langle s', a' \rangle$ .

Selectively allowing distinct latent vector assignments only in high-energy regions encourages behaviors that target these regions, thereby promoting diversity only in high-energy regions. Constraint shown in Eq. 1 formalizes this intuition: For a transition from  $\langle s, a \rangle$  to  $\langle s', a' \rangle$ , the latter’s latent vector can be far from the former’s by at most the Euclidean distance between the two pairs, scaled by  $f$ -energy of the latter and a factor  $\lambda_C$ .

$$\|\mu_{q(\cdot|s,a)} - \mu_{q(\cdot|s',a')}\| \leq \lambda_C \cdot \|s \oplus a - s' \oplus a'\| \cdot f(s', a') \quad (1)$$

The proposed constraint (Eq. 1) disregards the energy of the starting state-action pair,  $f(s, a)$ . The constraint enforces the same latent assignment for pairs in a low  $\rightarrow$  low energy transition and allows different latent assignments for pairs in a low  $\rightarrow$  high energy transition as visualized in Fig. 3b. The enforcement can lead to connected

low energy regions being assigned the same latent vector which is different from that of reachable high energy regions. Distinct latent vectors for low energy regions can result in behaviors visiting those low-energy regions, which is **not** desirable.

We rectify our constraints to prevent latent assignments for low energy regions, by enforcing them only with transitions with high energy starting pairs through scaling the constraint with  $f$ -energy of the starting pair, as shown in Eq. 2. Thus, when the starting pair is of low energy, the constraint implemented with a Lagrange multiplier is less effective due to a smaller violation.

The modified constraints encourage the decoder to effectively use the latent space to solely represent high-energy regions. We refer to the resulting objective as  $f$ -relevant diversity.

$$f(s, a) \cdot [\|\mu_{q(\cdot|s,a)} - \mu_{q(\cdot|s',a')}\|] \leq f(s, a) \cdot [\lambda_C \cdot \|s \oplus a - s' \oplus a'\| \cdot f(s', a')] \quad (2)$$

## 5.2 INFERRING A TASK-RELEVANCE MEASURE FROM DEMONSTRATIONS

Diverse task-accomplishing behaviors can be learned if an appropriate energy function  $f$  can be used to instantiate our  $f$ -relevant diversity objective. **We now present an approach to derive such an  $f$  from demonstrations, utilizing occupancies captured by the discriminator,  $D$ .**

Function  $f$  should indicate regions where agent occupation is favorable for the task while also capturing demonstrators’ heterogeneity. We classify these regions as: (I) regions occupied by all experts, and (II) subspaces where different experts occupy distinct smaller regions. We propose that capturing these two types of regions in the energy function provides us with an objective that encourages task accomplishment and generalization over **behavior factors**.

While the discriminator,  $D$ , could be used to model  $f$ , it would capture type II subspace insufficiently.  $D$  is trained to capture the union of demonstration occupancies, covering all type I regions,

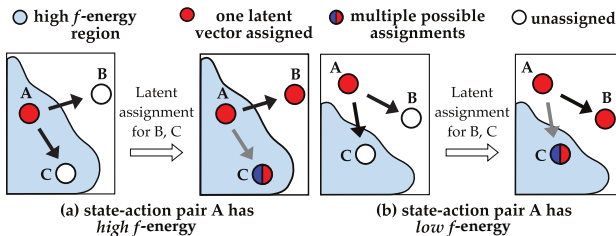


Figure 3: This figure shows a visualization of state transitions that an agent at state-action pair, A, may undergo and the effect of our constraint in Eq. 1. The light blue region indicates high  $f$ -energy and white state-action pairs do not yet have assigned latent vectors. The lightness of the arrows indicates the slackness of our constraint (Eq. 1) from scaling with  $f$ -energy on the right-hand side.

**Algorithm 1** Guided Strategy Discovery**Input:** Dataset of diverse expert demonstrations,  $\mathcal{D} = \{\tau_i^\xi\}$ **Output:** Latent-conditioned policy capturing diverse behaviors,  $\pi$ 

- 1: Initialize policy  $\pi$ , task relevance  $f$ , factor-specific residual  $g$ , bias  $b$ , and decoder  $q$
- 2: **for**  $i \in \{0, 1, 2, \dots\}$  epoch **do**
- 3:   Sample  $z^\pi$  from prior,  $\tau^\pi$  using policy  $\pi(\cdot | \cdot, z^\pi)$ ;  $\tau^\xi$  from  $\mathcal{D}$ , infer  $z^\xi$  using decoder  $q$
- 4:   **Compute discriminator outputs according to the conditioned distillation structure** (Eq. 3):  
 $D(s, a, z) = \sigma(\lambda_S \cdot [f(s, a) + g(s, a, z)] + b)$
- 5:   **Update  $f, g, b$  with gradient ascent on the discriminator objective computed by linearly combining  $J^{\text{GAIL}}$  (Ho & Ermon, 2016) and the distillation objective** (Eq. 4):  
 $\mathbb{E}_{\tau^\xi} [\log D(s, a, z^\xi)] + \mathbb{E}_{\tau^\pi} [\log(1 - D(s, a, z^\pi))] - \mathbb{E}_{\tau^\xi} [(g(s, a, z^\xi))^2] - \mathbb{E}_{\tau^\pi} [(g(s, a, z^\pi))^2]$
- 6:   **Compute the decoder objective using policy behavior samples** (Eysenbach et al., 2018):  
 $\mathbb{E}_{\tau^\pi} [\log \mathcal{N}(z^\pi | \mu_{q(\cdot | s, a)}, \Sigma_{q(\cdot | s, a)})]$
- 7:   **Update decoder with gradient ascent on the objective while enforcing task-relevance constraints** (Eq. 2):  $[\lambda_C \cdot \|s \oplus a - s' \oplus a'\| \cdot f(s', a') - \|\mu_{q(\cdot | s, a)} - \mu_{q(\cdot | s', a')}\|] \cdot f(s, a) \leq 0$
- 8:   **Update policy  $\pi$  with RL using linearly combined behavior imitation and diversity rewards,  $\log(D(s, a, z^\pi))$  and  $\log \mathcal{N}(z^\pi | \mu_{q(\cdot | s, a)}, \Sigma_{q(\cdot | s, a)})$ , respectively.**
- 9: **end for**

but only certain regions in the type II subspace that correspond to training demonstrations.  $f$  modeled in this way would limit behavior discovery beyond the training dataset.

We employ the distillation of demonstration-specific occupancy measures into a common measure, to model  $f$  and capture type II subspaces. **Demonstration-specific occupancies are obtained by conditioning the discriminator,  $D(s, a, z)$ , on the inferred latent vector,  $z$ . It is further split as a linear combination of a latent-independent measure, i.e., our desired energy function,  $f(s, a)$ , and a dependent term,  $g(s, a, z)$ , in the logit space as shown in Eq. 3, where  $g: S \times A \times Z \rightarrow [0, 1]$ ,  $\sigma$  is the logistic function,  $\lambda_S$ , a scaling constant, and  $b$ , a learnable bias.  $\lambda_S$  and  $b$  are introduced to transform the sum of bounded measures and enable  $D$  to capture most of the probability range  $[0, 1]$ . The discriminator is trained with an additional **distillation objective**, as shown in Eq. 4, to minimize the residual,  $g$ , to only capture **demonstration-specific** occupancy.**

$$D(s, a, z) = \sigma(\lambda_S \cdot [f(s, a) + g(s, a, z)] + b) \quad (3)$$

$$J^{\text{R}} := -E_{z, \pi} [(g(s, a, z))^2] \quad (4)$$

The objective,  $J^{\text{R}}$ , encourages  $f$  to fully capture both type I and II regions. Type I regions are captured by  $f$ , as  $g$  is driven to zero where occupancy is common across demonstrations and latent-independent.  $g$  is encouraged to be close to zero even in regions with demonstration-specific occupancy, causing it to capture minimal possible information and distilling the rest into  $f$ . We posit that  $f$  indicates entire subspaces where occupancy is demonstration dependent, i.e., type II subspaces, while  $g$  indicates regions in these subspaces specific to each demonstration. We call our procedure for deriving  $f$  conditioned distillation (ConDist), due to its use of latent conditioning and distillation, similar to prior reward distillation frameworks (Chen et al., 2020).

**Algorithm:** We combine  $f$ -relevant diversity and conditioned distillation to propose Guided Strategy Discovery (GSD). **High-level steps are outlined in Algorithm 1. Detailed steps can be found in Appendix C.** In each epoch, we sample behaviors using the policy conditioned on latent vectors from the prior (Line 3). Latent vectors for demonstrations are inferred using the decoder (Line 3). We use the proposed discriminator structure (Line 4), define the imitation and distillation objectives and update the energy function, residual, and bias, using gradients (Line 5). We optimize the variational lower bound (Line 6) while enforcing with proposed constraints (Line 7) to update the decoder using gradients. Finally, we update the policy with rewards from the discriminator and decoder (Line 8).

We posit that  $f$ -relevant diversity and conditioned distillation are synergistic. An accurate  $f$  function representing demonstrations can guide latent assignments and associated behaviors to generalize beyond demonstrations. A latent space representing diverse demonstrations can help distillation capture regions beyond demonstrations that generalize latent **behavior factors**. Fig. 2d (bottom) shows that with GSD, the learned behaviors in 2D PointMaze capture novel latent variations by passing through waypoints along  $x = 0$  while reaching closer to the goal of (2, 0) better than baselines in Figs. 2b, 2c. In addition, GSD produces a higher fraction of goal-reaching trajectories despite a low weight for the imitation objective and a weaker incentive to match the expert.

## 6 EVALUATION

We present empirical evaluation to answer the following research questions:

1. How do various methods perform in generalization to behaviors with novel latent factors while maintaining task performance? (Sec. 6.1)
2. How do various methods structure behaviors in the latent space? (Sec. 6.2)
3. How do various methods perform in learning diverse task-accomplishing behaviors from real-world human demonstrations? (Sec. 6.3)

**Domains:** Instead of using D3IL benchmarks (Jia et al., 2024), which focus on discrete behavior modes, we use domains with continuous variation and clear task objectives to better evaluate generalization from limited demonstrations. For Sec. 6.1, 6.2, we use HalfCheetah (Wawrzyński, 2009), FetchPickPlace (Plappert et al., 2018) and DriveLaneshift (Leurent, 2018) as they provide well-defined tasks with distinct one-dimensional (1D) factors. We script expert policies based on these factors. In HalfCheetah, the robot runs at various speeds; in FetchPickPlace, the arm places an objects at different locations; and in DriveLaneshift, the ego-car overtakes at varying headway distances. 1D factors help avoid multiple sources of heterogeneity allowing careful examination of learned policies.

**Methods:** We consider InfoGAIL as our base method, representative of approaches that combine online IL with MI-based diversity. Other heterogeneous online IL methods capture finite sets of behaviors and are omitted due to less scope for novel behavior discovery. Comparison with non-heterogeneous online IL methods Garg et al. (2021) is omitted as they do not directly address demonstration diversity. While we incorporate some improvements from online IL Orsini et al. (2021) across all the evaluated methods, a thorough evaluation of their integration with diversity objectives is left for future work. Comparison with offline IL approaches that learn without environment interaction is presented in Appendix D.3. We compare the following variants of InfoGAIL.

- IG: InfoGAIL (Li et al., 2017) with a continuous two dimensional latent variable.
- IG+Lipz: IG with Lipschitz constraints for decoder  $q$  to investigate the uniform diversity.
- IG+Con: IG with a conditioned discriminator  $D(s, a, z)$  structure to investigate the effect of conditioning the discriminator.
- IG+ConDist: IG with our proposed conditioned distillation to investigate the effect of extraction of a task-relevance measure (Eqs. 4, 3).
- IG+ConDist+Lipz: IG+ConDist with Lipschitz to investigate the uniform diversity formulation alongside conditioned distillation.
- **GSD (Ours)**: IG+ConDist with our proposed task-relevant diversity formulation (Eq. 2).

### 6.1 QUANTITATIVE EVALUATION

We investigate whether learned behaviors can represent factor values in the disjoint test region  $\text{Te}(\Omega)$ , after training on demonstrations,  $\mathcal{D}$ , from the train region,  $\text{Tr}(\Omega)$  (see Sec. 3). We consider factors that are measurable from trajectories (only for the sake of evaluation) to assess recovery performance, i.e., how well the learned latent space can represent expert behavior, by comparing desired and measured factor values. When diverse expert behaviors form distinct modes, this framework checks if continuous factors underlying them can be accurately identified and generalized.

**Splits:** We divide the bounded 1D factor range into five consecutive equal-sized intervals:

- **Interpolation:** The first, third, and fifth intervals represent the train region, and the second and fourth are the test region. The split allows us to evaluate the ability to interpolate behaviors to two factor space intervals while providing three non-consecutive intervals to represent the factor.
- **Extrapolation:** The second and fourth intervals represent the train region, while the first and fifth intervals are the test region. We choose two non-consecutive intervals for the train region to have a sparse dataset while providing enough diversity to represent the factor.

These splits evaluate how well the latent space captures factors to interpolate and extrapolate behaviors. We use five demonstrations per interval (details in Appendix B).

**Metrics:** We search for desired behaviors using  $K \in \{10, 20, 30, 40, 50\}$  latent vector samples from the prior  $p_z(\cdot)$ , where  $K$  represents the test time search sample-complexity, varied to investigate how well we generate desired behaviors from limited samples. We roll out policies conditioned on the sampled vectors, measure the factors of the sampled behaviors, and compute the least mean absolute error (MAE) between the desired and the  $K$  measured values, averaging over  $1500/K$  rounds. We

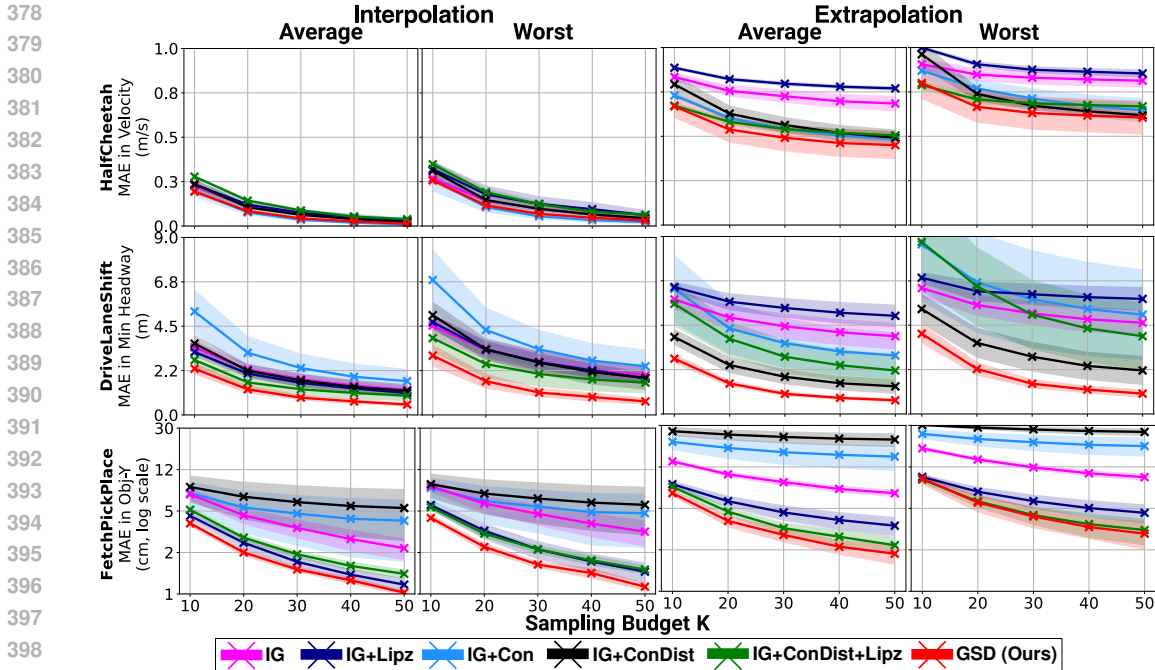


Figure 4: The figure shows average and worst-case recovery errors over the three domains and two factor space splits corresponding to interpolative and extrapolative generalization. Shaded regions are standard errors over five train seeds. GSD outperforms baselines in recovery of unseen latent factors across most domains and splits.

refer to this metric as the recovery error. We consider the midpoints of test intervals as the set of desired values. We report average and worst errors over desired values in the test region, providing estimates of closeness between desired values and closest available behavior’s factor, on average and worst case. We report mean and standard errors over five train seeds in Fig. 4. We evaluate task performance by averaging environment returns over 1500 latent samples. We show recovery and task performance tradeoff in Fig. 5. Exact numbers are provided in Appendix H.

**Lipschitz constraints:** For HalfCheetah (Fig. 4, top row), IG+Lipz and IG+ConDist+Lipz have worse recovery errors compared to IG and IG+ConDist, indicated by the dark blue curve above magenta, and green above black, respectively. For DriveLaneShift (middle row), IG+Lipz and IG+ConDist+Lipz exhibit the same trend against IG and IG+ConDist: for interpolation, errors are improved shown by dark blue falling below magenta and green below black; and for extrapolation, the errors are worsened. For FetchPickPlace (bottom row), IG+Lipz and IG+ConDist+Lipz improve over IG and IG+ConDist indicated by dark blue and green consistently below magenta and black respectively. Lipschitz constraints seem to be benefiting FetchPickPlace alone, which might be due to “uniform diversity” aligning with object-

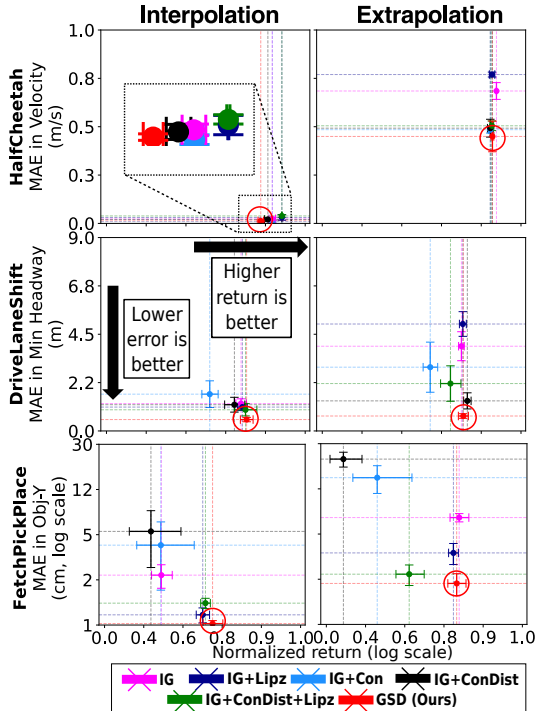


Figure 5: The figure shows the tradeoff between task and recovery performance for three domains and two splits. Error bars show standard errors over five seeds. High returns (x-axis) and low errors (y-axis) are better. GSD (circled in red) improves recovery while retaining or improving task performance across most domains and splits.



432 position factors, that is absent in other domains. This supports our hypothesis that relevant factors  
 433 for diversity must be inferred from demonstrations to benefit all domains.

434 **Conditioning:** For HalfCheetah (top row), IG+Con improves errors compared to IG, indicated by  
 435 the light blue curve below magenta. However, DriveLaneshift and FetchPickPlace, IG+Con seems  
 436 to worsen performance, with light blue largely above magenta in the bottom two rows. Worsened  
 437 errors may be a result of conditioning on a latent variable capturing arbitrary factors.

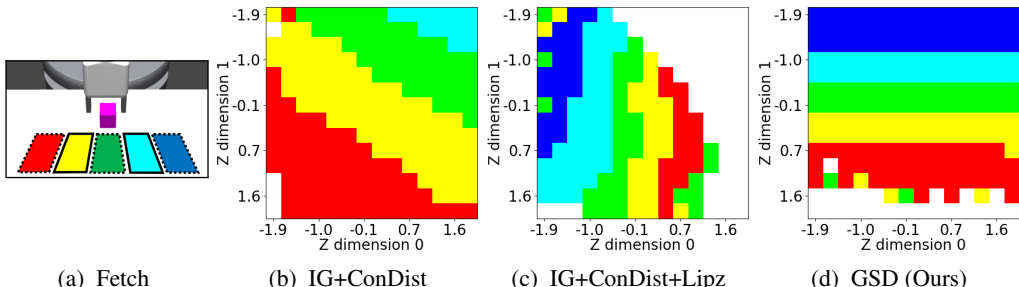
438 **Conditioned Distillation:** For HalfCheetah (top row), IG+ConDist and IG+ConDist+Lipz improve  
 439 over IG and IG+Lipz for extrapolation, indicated by black and green curves below magenta and  
 440 dark blue respectively. They remain on par for interpolation. For DriveLaneshift as well (middle  
 441 row), IG+ConDist and IG+ConDist+Lipz improve over IG and IG+Lipz for extrapolation ( $K \geq 30$ )  
 442 and remain on par for/slightly improve interpolation. For FetchPickPlace (bottom row), the trends  
 443 are interesting. IG+ConDist worsens errors over IG for interpolation and extrapolation, indicated  
 444 by black above magenta. However, with Lipz’s addition, IG+ConDist+Lipz tends close to IG+Lipz  
 445 for interpolation and outperforms it for extrapolation. The patterns firstly suggest that conditioned  
 446 distillation can improve extrapolation performance. In addition, for FetchPickPlace where Lipschitz  
 447 constraints are particularly effective, conditioned distillation can further improve extrapolation.

448 **Task-relevant Diversity:** GSD improves recovery over other approaches across most domains and  
 449 splits, shown by the red curve below others in all plots, except for interpolation with HalfCheetah  
 450 (top row, first two columns). In HalfCheetah (top row), the close performance across methods may  
 451 be attributed to wide differences in gait styles across velocities that are challenging to interpolate or  
 452 extrapolate. In DriveLaneshift (middle row), GSD reduces recovery error considerably over other  
 453 approaches. In FetchPickPlace, GSD is closely matched by IG+Lipz or IG+ConDist+Lipz as Lips-  
 454 chitz constraints already capture relevant factors. Nevertheless, GSD can further improve recovery.

455 **Tradeoff between Task and Recovery Performance:** Across all domains, GSD either matches or  
 456 improves average normalized returns over the latent prior, as indicated in Fig. 5 by the red cross  
 457 generally being aligned with or positioned further right than others in all domain-split combinations  
 458 but one. These results demonstrate the effectiveness of GSD’s task-relevant diversity formulation in  
 459 learning behaviors that reduce recovery error while maintaining task performance.

461 6.2 QUALITATIVE EVALUATION

462 We visualize the nature of the learned latent spaces for extrapolation in FetchPickPlace in Fig. 6.  
 463 IG+ConDist learns a large behavior set for placing the object in the red test region (indicated by red  
 464 cells in Fig. 6b), but ignores the dark-blue test region. While IG+ConDist+Lipz learns behaviors for  
 465 all regions, it learns several that fail to place the object quickly enough (indicated by white cells in  
 466 Fig. 6c). GSD learns behaviors that achieve the task (few white cells in Fig. 6d) while representing  
 467 all place locations equally in proximity to each other (roughly equal number of cells across colors  
 468 nearby each other). GSD exhibits potential for improving the accountability of policy learning by  
 469 enabling well structured latent spaces.



480 Figure 6: Fig.6a shows FetchPickPlace with object placement locations color-coded. Solid and  
 481 dotted boundaries indicate train and test regions respectively. Figs.6b,6c,6d: Policy behaviors are  
 482 shown in the 2D latent space through colors for resulting place-locations shown in 6a. White regions  
 483 indicate failed placements or placements with low task reward. Behaviors with IG+ConDist (6b),  
 484 IG+ConDist+Lipz (6c) either represent the relevant regions disproportionately or fail to accomplish  
 485 the task. **Behaviors with GSD (6d) accomplish the task (low presence of white cells) and represent all regions well (roughly equal number of cells across colors).**

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

### 6.3 EVALUATION WITH REAL-WORLD HUMAN DEMONSTRATIONS

We further evaluate our approach to test scalability to complex tasks with human demonstrations in a Table Tennis (TT) domain. TT represents a dynamic domain that requires precise robot motion and fast reaction times while acting on noisy observations. Our physical setup consists of a Barrett WAM Arm mounted to the ceiling in front of a TT table, and a racquet attached as the arm’s effector. Balls are launched using a Butterfly Amicus launcher at a fixed orientation and velocity with some noise. Balls are detected and tracked using a YOLO object detector and a Kalman Filter. An expert provides kinesthetic demonstrations of push and lob strokes. We recreate the setup in simulation with PyBullet for behavior learning. Ball initialization and observation noise levels in the simulation match real data. Complete details are in Appendix E.

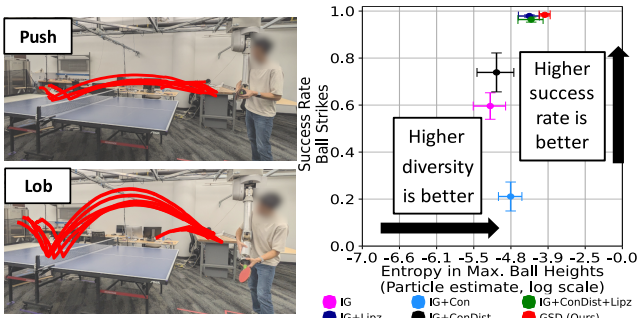


Figure 7: **Left:** The images visualize ball trajectories achieved by an expert kinesthetically demonstrating two types of strokes. **Right:** The figure shows the tradeoff between ball striking success rate and diversity in ball heights achieved. GSD outperforms baselines for both metrics.

While multiple continuous factors may exist underlying TT stroke styles, we evaluate generalization for maximum ball height, which we assume to be one of the underlying continuous factors. We evaluate various methods in simulation for achieving high diversity in ball heights. We compute entropy in ball height values using particle estimates (Singh et al., 2003), after disregarding unsuccessful trials that fail to strike the ball over the table. We report the success rate traded off with diversity in ball heights in Fig. 7 (right). Our method GSD outperforms all baselines in both measures of success rate and entropy.

## 7 CONCLUSION, LIMITATIONS AND FUTURE WORK

We study the problem of generalization from diverse demonstrations over underlying latent factors. We investigate the shortcomings of prior MI-based methods and propose a novel diversity formulation. Our empirical evaluation shows that our approach GSD improves the recovery of factors over the next best baseline (for  $K=50$ ) by 18.3% and 24.6% for interpolation and extrapolation respectively while retaining task performance in three domains with synthetic demonstrations. Our qualitative analysis shows the potential of our approach in making learned policies accountable. Lastly, our experiments with real-world human demonstrations shows that our framework can capture a diverse range of task-accomplishing behaviors in a challenging domain requiring quick response times.

**Limitations:** Our experiments focus on demonstrations with one-dimensional latent factors. Our approach may struggle with higher dimensional or non-Markovian factors, which could require specialized designs for disentangling dimensions or capturing observation history dependence. Scaling to visual domains, where continuous factors must be inferred from sparsely distributed demonstrations, may also be challenging, as simple models for the energy function  $f$  may not generalize well. Furthermore, our assumption that demonstration occupancies correlate with task success may be violated with non-expert demonstrators or partial observability, which may require state estimation models. While current work prioritizes validating our core contributions, we plan to evaluate scalability in future work.

Our evaluation with human demonstrations is further limited to quantitative metrics. We aim to conduct user studies to subjectively evaluate behaviors in human robot interaction settings. Our evaluation is further limited to simulated domains. We aim to explore the efficacy of our diversity formulation for learning novel behaviors in physical robot systems with improved data-sample efficiency. We further aim to explore the theoretical implications of our formulation and its alignment with the imitation objective. Further limitations are discussed in Appendix G.

540 REPRODUCIBILITY STATEMENT

541  
542 Implementation details, hyperparameters and evaluation procedures are detailed in Appen-  
543 dices D, E.2. Data generation and collection is detailed in Appendices A, B, E.1.  
544

545 REFERENCES

546  
547 Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In  
548 *Proceedings of the 21st International Conference on Machine learning*, pp. 1, 2004.  
549

550 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Con-  
551 crete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.

552 Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob  
553 McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience re-  
554 play. *Advances in Neural Information Processing Systems*, 30, 2017.  
555

556 David Barber and Felix Agakov. The IM algorithm: a variational approach to Information Maxi-  
557 mization. *Advances in Neural Information Processing Systems*, 16(320):201, 2004.  
558

559 Sumeet Batra, Bryon Tjanaka, Stefanos Nikolaidis, and Gaurav Sukhatme. Quality diversity for  
560 robot learning: Limitations and future directions. In *Proceedings of the Genetic and Evolutionary  
561 Computation Conference Companion*, pp. 587–590, 2024.

562 Carolin Benjamins, Theresa Eimer, Frederik Schubert, Aditya Mohan, André Biedenkapp, Bodo  
563 Rosenhahn, Frank Hutter, and Marius Lindauer. Contextualize me—the case for context in rein-  
564 forcement learning. *arXiv preprint arXiv:2202.04500*, 2022.  
565

566 Marcel Binz and Dominik M. Endres. Where do human heuristics come from? *ArXiv*,  
567 abs/1902.07580, 2019. URL [https://api.semanticscholar.org/CorpusID:  
568 67769766](https://api.semanticscholar.org/CorpusID:67769766).

569 Gisela Böhm and Hans-Rüdiger Pfister. How people explain their own and others’ behavior: a  
570 theory of lay causal explanations. *Frontiers in psychology*, 6:109763, 2015.  
571

572 Stephen P Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge university press, 2004.

573 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and  
574 Wojciech Zaremba. Openai gym. arxiv. *arXiv preprint arXiv:1606.01540*, 10, 2016.  
575

576 Zhangjie Cao, Yilun Hao, Mengxi Li, and Dorsa Sadigh. Learning feasibility to imitate demonstra-  
577 tors with different dynamics. *arXiv preprint arXiv:2110.15142*, 2021.  
578

579 Letian Chen, Rohan Paleja, Muyleng Ghuy, and Matthew Gombolay. Joint goal and strategy infer-  
580 ence across heterogeneous demonstrators via reward network distillation. In *Proceedings of the  
581 2020 ACM/IEEE International Conference on human-robot interaction*, pp. 659–668, 2020.

582 Letian Chen, Rohan Paleja, and Matthew Gombolay. Learning from suboptimal demonstration via  
583 self-supervised reward regression. In *Conference on robot learning*, pp. 1262–1277. PMLR, 2021.  
584

585 Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shu-  
586 ran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint  
587 arXiv:2303.04137*, 2023.

588 Jongwook Choi, Archit Sharma, Honglak Lee, Sergey Levine, and Shixiang Shane Gu. Variational  
589 empowerment as representation learning for goal-based reinforcement learning. *arXiv preprint  
590 arXiv:2106.01404*, 2021.  
591

592 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
593 reinforcement learning from human preferences. *Advances in Neural Information Processing  
Systems*, 30, 2017.

- 594 Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generaliza-  
595 tion in reinforcement learning. In *International Conference on Machine Learning*, pp. 1282–1289.  
596 PMLR, 2019.
- 597 Antonio Coronato, Muddasar Naeem, Giuseppe De Pietro, and Giovanni Paragliola. Reinforcement  
598 learning for intelligent healthcare applications: A survey. *Artificial intelligence in medicine*, 109:  
599 101964, 2020.
- 600  
601 Wojciech M Czarnecki, Razvan Pascanu, Simon Osindero, Siddhant Jayakumar, Grzegorz Swirszcz,  
602 and Max Jaderberg. Distilling policy distillation. In *The 22nd International Conference on artifi-*  
603 *cial intelligence and statistics*, pp. 1331–1340. PMLR, 2019.
- 604  
605 Allan Dafoe, Edward Hughes, Yoram Bachrach, Tantum Collins, Kevin R McKee, Joel Z  
606 Leibo, Kate Larson, and Thore Graepel. Open problems in cooperative ai. *arXiv preprint*  
607 *arXiv:2012.08630*, 2020.
- 608  
609 Anca Dragan and Siddhartha Srinivasa. Generating legible motion. *Frontiers in psychology*, 2013.
- 610  
611 Danny Driess, Jung-Su Ha, and Marc Toussaint. Deep visual reasoning: Learning to predict  
612 action sequences for task and motion planning from an initial scene image. *arXiv preprint*  
613 *arXiv:2006.05398*, 2020.
- 614  
615 Yan Duan, Marcin Andrychowicz, Bradley Stadie, OpenAI Jonathan Ho, Jonas Schneider, Ilya  
616 Sutskever, Pieter Abbeel, and Wojciech Zaremba. One-shot imitation learning. *Advances in*  
617 *Neural Information Processing Systems*, 30, 2017.
- 618  
619 Marco Ewerton, Guilherme Maeda, Gerrit Kollegger, Josef Wiemeyer, and Jan Peters. Incremental  
620 imitation learning of context-dependent motor skills. In *IEEE-RAS 16th International Conference*  
621 *on Humanoid Robots (Humanoids)*, pp. 351–358. IEEE, 2016.
- 622  
623 Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need:  
624 Learning skills without a reward function. In *International Conference on Learning Representa-*  
625 *tions*, 2018.
- 626  
627 Chelsea Finn, Paul Christiano, Pieter Abbeel, and Sergey Levine. A connection between generative  
628 adversarial networks, inverse reinforcement learning, and energy-based models. *arXiv preprint*  
629 *arXiv:1611.03852*, 2016.
- 630  
631 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation  
632 of deep networks. In *International Conference on machine learning*, pp. 1126–1135. PMLR,  
633 2017a.
- 634  
635 Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imita-  
636 tion learning via meta-learning. In *Conference on Robot Learning*, pp. 357–368. PMLR, 2017b.
- 637  
638 Pete Florence, Corey Lynch, Andy Zeng, Oscar A Ramirez, Ayzaan Wahid, Laura Downs, Adrian  
639 Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In  
640 *Conference on Robot Learning*, pp. 158–168. PMLR, 2022.
- 641  
642 Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse rein-  
643 forcement learning. *arXiv preprint arXiv:1710.11248*, 2017.
- 644  
645 Kanishk Gandhi, Siddharth Karamcheti, Madeline Liao, and Dorsa Sadigh. Eliciting compatible  
646 demonstrations for multi-human imitation learning. In *Conference on Robot Learning*, pp. 1981–  
647 1991. PMLR, 2023.
- 648  
649 Yapeng Gao, Jonas Tebbe, and Andreas Zell. Optimal stroke learning with policy gradient approach  
650 for robotic table tennis. *Applied Intelligence*, 53(11):13309–13322, 2023.
- 651  
652 Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, and Stefano Ermon. Iq-learn:  
653 Inverse soft-q learning for imitation. *Advances in Neural Information Processing Systems*, 34:  
654 4028–4039, 2021.



- 648 Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence minimization  
649 perspective on imitation learning methods. In *Conference on Robot Learning*, pp. 1259–1277.  
650 PMLR, 2020.
- 651 Diego Gomez, Michael Bowling, and Marlos C Machado. Proper laplacian representation learning.  
652 *arXiv preprint arXiv:2310.10833*, 2023.
- 653 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
654 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in Neural Informa-  
655 tion Processing Systems*, 27, 2014.
- 656 Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv  
657 preprint arXiv:1611.07507*, 2016.
- 658 Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Im-  
659 proved training of wasserstein gans. *Advances in Neural Information Processing Systems*, 30,  
660 2017.
- 661 Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with  
662 deep energy-based policies. In *International Conference on machine learning*, pp. 1352–1361.  
663 PMLR, 2017.
- 664 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
665 maximum entropy deep reinforcement learning with a stochastic actor. In *International Confer-  
666 ence on machine learning*, pp. 1861–1870. PMLR, 2018.
- 667 Steven Hansen, Will Dabney, Andre Barreto, Tom Van de Wiele, David Warde-Farley, and  
668 Volodymyr Mnih. Fast task inference with variational intrinsic successor features. *arXiv preprint  
669 arXiv:1906.05030*, 2019.
- 670 Mahta HassanPour Zonoozi and Vahid Seydi. A survey on adversarial domain adaptation. *Neural  
671 Processing Letters*, 55(3):2429–2469, 2023.
- 672 Karol Hausman, Yevgen Chebotar, Stefan Schaal, Gaurav Sukhatme, and Joseph J Lim. Multi-modal  
673 imitation learning from unstructured demonstrations using generative adversarial nets. *Advances  
674 in Neural Information Processing Systems*, 30, 2017.
- 675 Alexander Herzog, Kanishka Rao, Karol Hausman, Yao Lu, Paul Wohlhart, Mengyuan Yan, Jessica  
676 Lin, Montserrat Gonzalez Arenas, Ted Xiao, Daniel Kappler, et al. Deep rl at scale: Sorting waste  
677 in office buildings with a fleet of mobile manipulators. *arXiv preprint arXiv:2305.03270*, 2023.
- 678 Jonathan Ho and Stefano Ermon. Generative Adversarial Imitation Learning. *Advances in Neural  
679 Information Processing Systems*, 29, 2016.
- 680 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are uni-  
681 versal approximators. *Neural networks*, 2(5):359–366, 1989.
- 682 Maxence Hussonnois, Thommen George Karimpanal, and Santu Rana. Controlled diversity with  
683 preference: Towards learning a diverse set of desired skills. *arXiv preprint arXiv:2303.04592*,  
684 2023.
- 685 Sagar Imambi, Kolla Bhanu Prakash, and GR Kanagachidambaresan. Pytorch. *Programming with  
686 TensorFlow: Solution for Edge Computing Applications*, pp. 87–104, 2021.
- 687 Rohit Jena, Changliu Liu, and Katia Sycara. Augmenting gail with bc for sample efficient imitation  
688 learning. In *Conference on Robot Learning*, pp. 80–90. PMLR, 2021.
- 689 Xiaogang Jia, Denis Blessing, Xinkai Jiang, Moritz Reuss, Atalay Donat, Rudolf Lioutikov, and  
690 Gerhard Neumann. Towards diverse behaviors: A benchmark for imitation learning with human  
691 demonstrations. *arXiv preprint arXiv:2402.14606*, 2024.
- 692 Liyiming Ke, Sanjiban Choudhury, Matt Barnes, Wen Sun, Gilwoo Lee, and Siddhartha Srinivasa.  
693 Imitation learning as f-divergence minimization. In *Algorithmic Foundations of Robotics XIV:  
694 Proceedings of the 14th Workshop on the Algorithmic Foundations of Robotics 14*, pp. 313–329.  
695 Springer, 2021.

- 702 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
703 *arXiv:1312.6114*, 2013.  
704
- 705 Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of zero-shot gener-  
706 alisation in deep reinforcement learning. *Journal of Artificial Intelligence Research*, 76:201–264,  
707 2023.
- 708 W Bradley Knox and Peter Stone. Tamer: Training an agent manually via evaluative reinforcement.  
709 In *7th IEEE International Conference on development and learning*, pp. 292–297. IEEE, 2008.  
710
- 711 W Bradley Knox and Peter Stone. Interactively shaping agents via human reinforcement: The tamer  
712 framework. In *Proceedings of the 5th International Conference on Knowledge Capture*, pp. 9–16,  
713 2009.
- 714 Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan Tompson.  
715 Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial  
716 imitation learning. *arXiv preprint arXiv:1809.02925*, 2018.
- 717 Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. RMA: Rapid motor adaptation for  
718 legged robots. *arXiv preprint arXiv:2107.04034*, 2021.  
719
- 720 Saurabh Kumar, Aviral Kumar, Sergey Levine, and Chelsea Finn. One solution is not all you need:  
721 Few-shot extrapolation via structured maxent rl. *Advances in Neural Information Processing*  
722 *Systems*, 33:8198–8210, 2020.
- 723 Michael Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang, Lerrel  
724 Pinto, and Pieter Abbeel. Urlb: Unsupervised reinforcement learning benchmark. *arXiv preprint*  
725 *arXiv:2110.15191*, 2021.  
726
- 727 Michael Laskin, Hao Liu, Xue Bin Peng, Denis Yarats, Aravind Rajeswaran, and Pieter Abbeel. Cic:  
728 Contrastive intrinsic control for unsupervised skill discovery. *arXiv preprint arXiv:2202.00161*,  
729 2022.
- 730 Edouard Leurent. An environment for autonomous driving decision-making. <https://github.com/eleurent/highway-env>, 2018.  
731
- 732 Chenhao Li, Sebastian Blaes, Pavel Kolev, Marin Vlastelica, Jonas Frey, and Georg Martius. Ver-  
733 satile skill control via self-supervised adversarial imitation of unlabeled mixed motions. In *2023*  
734 *IEEE international conference on robotics and automation (ICRA)*, pp. 2944–2950. IEEE, 2023.  
735
- 736 Chenhao Li, Elijah Stanger-Jones, Steve Heim, and Sangbae Kim. Fld: Fourier latent dynamics for  
737 structured motion representation and learning. *arXiv preprint arXiv:2402.13820*, 2024.
- 738 Yunzhu Li, Jiaming Song, and Stefano Ermon. Infogail: Interpretable imitation learning from visual  
739 demonstrations. *Advances in Neural Information Processing Systems*, 30, 2017.
- 740 Minghuan Liu, Tairan He, Minkai Xu, and Weinan Zhang. Energy-based imitation learning. *arXiv*  
741 *preprint arXiv:2004.09395*, 2020.  
742
- 743 Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and  
744 Pierre Sermanet. Learning latent plans from play. In *Conference on robot learning*, pp. 1113–  
745 1132. PMLR, 2020.
- 746 Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-  
747 Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline  
748 human demonstrations for robot manipulation. In *Conference on Robot Learning*, pp. 1678–1690.  
749 PMLR, 2022.
- 750 Russell Mendonca, Oleh Rybkin, Kostas Daniilidis, Danijar Hafner, and Deepak Pathak. Discover-  
751 ing and achieving goals via world models. *Advances in Neural Information Processing Systems*,  
752 34:24379–24391, 2021.  
753
- 754 Josh Merel, Leonard Hasenclever, Alexandre Galashov, Arun Ahuja, Vu Pham, Greg Wayne,  
755 Yee Whye Teh, and Nicolas Heess. Neural probabilistic motor primitives for humanoid control.  
*arXiv preprint arXiv:1811.11711*, 2018.

- 756 Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization  
757 for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.  
758
- 759 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle-  
760 mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level  
761 control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.  
762
- 763 Eduardo F Morales and Claude Sammut. Learning to fly by combining reinforcement learning with  
764 behavioural cloning. In *Proceedings of the 21st International Conference on Machine learning*,  
765 pp. 76, 2004.
- 766 Suraj Nair, Eric Mitchell, Kevin Chen, Silvio Savarese, Chelsea Finn, et al. Learning language-  
767 conditioned robot behavior from offline data and crowd-sourced annotation. In *Conference on*  
768 *Robot Learning*, pp. 1303–1315. PMLR, 2022.
- 769 Tianwei Ni, Harshit Sikchi, Yufei Wang, Tejus Gupta, Lisa Lee, and Ben Eysenbach. f-irl: Inverse  
770 reinforcement learning via state marginal matching. In *Conference on Robot Learning*, pp. 529–  
771 551. PMLR, 2021.  
772
- 773 OpenAI. o1. <https://openai.com/index/learning-to-reason-with-llms/>,  
774 2024. Learning to reason with Large Language Models.  
775
- 776 Manu Orsini, Anton Raichuk, Léonard Hussenot, Damien Vincent, Robert Dadashi, Sertan Girgin,  
777 Matthieu Geist, Olivier Bachem, Olivier Pietquin, and Marcin Andrychowicz. What matters for  
778 adversarial imitation learning? *Advances in Neural Information Processing Systems*, 34:14656–  
779 14668, 2021.
- 780 Takayuki Osa, Joni Pajarinen, Gerhard Neumann, J Andrew Bagnell, Pieter Abbeel, Jan Peters, et al.  
781 An algorithmic perspective on imitation learning. *Foundations and Trends® in Robotics*, 7(1-2):  
782 1–179, 2018.  
783
- 784 Takayuki Osa, Voot Tangkaratt, and Masashi Sugiyama. Discovering diverse solutions in deep  
785 reinforcement learning by maximizing state–action-based mutual information. *Neural Networks*,  
786 152:90–104, 2022.
- 787 Charles Packer, Katelyn Gao, Jernej Kos, Philipp Krähenbühl, Vladlen Koltun, and Dawn Song.  
788 Assessing generalization in deep reinforcement learning. *arXiv preprint arXiv:1810.12282*, 2018.  
789
- 790 Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander  
791 Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic  
792 learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.  
793
- 794 Rohan Paleja, Andrew Silva, Letian Chen, and Matthew Gombolay. Interpretable and personalized  
795 apprenticeship scheduling: Learning interpretable scheduling policies from heterogeneous user  
796 demonstrations. *Advances in Neural Information Processing Systems*, 33:6417–6428, 2020.
- 797 Seohong Park, Jongwook Choi, Jaekyeom Kim, Honglak Lee, and Gunhee Kim. Lipschitz-  
798 constrained unsupervised skill discovery. In *International Conference on Learning Representations*,  
799 2022.  
800
- 801 Seohong Park, Kimin Lee, Youngwoon Lee, and Pieter Abbeel. Controllability-aware unsupervised  
802 skill discovery. In *International Conference on Machine Learning*, pp. 27225–27245. PMLR,  
803 2023.
- 804 Seohong Park, Oleh Rybkin, and Sergey Levine. METRA: Scalable unsupervised RL with metric-  
805 aware abstraction. In *The Twelfth International Conference on Learning Representations*, 2024.  
806
- 807 Tim Pearce, Tabish Rashid, Anssi Kanervisto, Dave Bignell, Mingfei Sun, Raluca Georgescu, Ser-  
808 gio Valcarcel Macua, Shan Zheng Tan, Ida Momennejad, Katja Hofmann, et al. Imitating human  
809 behaviour with diffusion models. In *The Eleventh International Conference on Learning Representations (ICLR 2023)*, 2023.

- 810 Jian-Wei Peng, Min-Chun Hu, and Wei-Ta Chu. An imitation learning framework for generat-  
811 ing multi-modal trajectories from unstructured demonstrations. *Neurocomputing*, 500:712–723,  
812 2022a.
- 813  
814 Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-scale  
815 reusable adversarial skill embeddings for physically simulated characters. *ACM Transactions On*  
816 *Graphics (TOG)*, 41(4):1–17, 2022b.
- 817 Matthias Plappert, Marcin Andrychowicz, Alex Ray, Bob McGrew, Bowen Baker, Glenn Pow-  
818 ell, Jonas Schneider, Josh Tobin, Maciek Chociej, Peter Welinder, et al. Multi-goal reinforce-  
819 ment learning: Challenging robotics environments and request for research. *arXiv preprint*  
820 *arXiv:1802.09464*, 2018.
- 821 Mariya Popova, Olexandr Isayev, and Alexander Tropsha. Deep reinforcement learning for de novo  
822 drug design. *Science advances*, 4(7):eaap7885, 2018.
- 823  
824 Yiwen Qiu, Jialong Wu, Zhangjie Cao, and Mingsheng Long. Out-of-dynamics imitation learning  
825 from multimodal demonstrations. In *Conference on Robot Learning*, pp. 1071–1080. PMLR,  
826 2023.
- 827 Deepak Ramachandran and Eyal Amir. Bayesian inverse reinforcement learning. In *IJCAI*, vol-  
828 ume 7, pp. 2586–2591, 2007.
- 829  
830 Siddharth Reddy, Anca D Dragan, and Sergey Levine. Sqil: Imitation learning via reinforcement  
831 learning with sparse rewards. *arXiv preprint arXiv:1905.11108*, 2019.
- 832  
833 Allen Z Ren, Justin Lidard, Lars L Ankile, Anthony Simeonov, Pulkit Agrawal, Anirudha Majum-  
834 dar, Benjamin Burchfiel, Hongkai Dai, and Max Simchowitz. Diffusion policy policy optimiza-  
835 tion. *arXiv preprint arXiv:2409.00588*, 2024.
- 836  
837 Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and struc-  
838 tured prediction to no-regret online learning. In *Proceedings of the fourteenth International Con-*  
839 *ference on artificial intelligence and statistics*, pp. 627–635. JMLR Workshop and Conference  
Proceedings, 2011.
- 840 Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint*  
841 *arXiv:1609.04747*, 2016.
- 842  
843 Stuart Russell. *Human compatible: Artificial intelligence and the problem of control*. Penguin,  
844 2019.
- 845 Penelope M Sanderson. The human planning and scheduling role in advanced manufacturing sys-  
846 tems: An emerging human factors domain. *Human Factors*, 31(6):635–666, 1989.
- 847  
848 Mariah L Schrum, Erin Hedlund-Botti, and Matthew Gombolay. Reciprocal MIND MELD: Improv-  
849 ing learning from demonstration via personalized, reciprocal teaching. In *Conference on Robot*  
850 *Learning*, pp. 956–966. PMLR, 2023a.
- 851  
852 Mariah L Schrum, Emily Sumner, Matthew C Gombolay, and Andrew Best. Maveric: A data-driven  
853 approach to personalized autonomous driving. *arXiv preprint arXiv:2301.08595*, 2023b.
- 854  
855 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-  
856 dimensional continuous control using generalized advantage estimation. *arXiv preprint*  
*arXiv:1506.02438*, 2015.
- 857  
858 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
859 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 860  
861 Nur Muhammad Shafiullah, Zichen Cui, Ariuntuya Arty Altanzaya, and Lerrel Pinto. Behavior  
862 transformers: Cloning  $k$  modes with one stone. *Advances in Neural Information Processing*  
863 *Systems*, 35:22955–22968, 2022.
- 864  
865 Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware  
866 unsupervised discovery of skills. *arXiv preprint arXiv:1907.01657*, 2019.



- 864 Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for  
865 robotic manipulation. In *Conference on Robot Learning*, pp. 785–799. PMLR, 2023.
- 866
- 867 Andrew Silva, Nina Moorman, William Silva, Zulfiqar Zaidi, Nakul Gopalan, and Matthew Gom-  
868 bolay. Lancon-learn: Learning with language to enable generalization in multi-task manipulation.  
869 *IEEE Robotics and Automation Letters*, 7(2):1635–1642, 2021.
- 870 David Silver, Satinder Singh, Doina Precup, and Richard S Sutton. Reward is enough. *Artificial*  
871 *Intelligence*, 299:103535, 2021.
- 872
- 873 Özgür Şimşek and Andrew G Barto. Using relative novelty to identify useful temporal abstractions in  
874 reinforcement learning. In *Proceedings of the 21st International Conference on Machine learning*,  
875 pp. 95, 2004.
- 876 Harshinder Singh, Neeraj Misra, Vladimir Hnizdo, Adam Fedorowicz, and Eugene Demchuk. Near-  
877 est neighbor estimates of entropy. *American journal of mathematical and management sciences*,  
878 23(3-4):301–321, 2003.
- 879
- 880 Nate Soares and Benja Fallenstein. Aligning superintelligence with human interests: A technical  
881 research agenda. *Machine Intelligence Research Institute (MIRI) technical report*, 8, 2014.
- 882 Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep  
883 conditional generative models. *Advances in neural information processing systems*, 28, 2015.
- 884
- 885 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 886
- 887 Andrew Szot, Amy Zhang, Dhruv Batra, Zsolt Kira, and Franziska Meier. Bc-irl: Learning general-  
888 izable reward functions from demonstrations. *arXiv preprint arXiv:2303.16194*, 2023.
- 889 Adrien Ali Taiga, Rishabh Agarwal, Jesse Farebrother, Aaron Courville, and Marc G Bellemare.  
890 Investigating multi-task pretraining and generalization in reinforcement learning. In *The 11th*  
891 *International Conference on Learning Representations*, 2022.
- 892
- 893 Chen Tang, Ben Abbatematteo, Jiaheng Hu, Rohan Chandra, Roberto Martín-Martín, and Peter  
894 Stone. Deep reinforcement learning for robotics: A survey of real-world successes. *arXiv preprint*  
895 *arXiv:2408.03539*, 2024.
- 896 Voot Tangkaratt, Bo Han, Mohammad Emtiyaz Khan, and Masashi Sugiyama. Variational imitation  
897 learning with diverse-quality demonstrations. In *Proceedings of the 37th International Conference*  
898 *on Machine Learning*, pp. 9407–9417, 2020.
- 899
- 900 Yee Teh, Victor Bapst, Wojciech M Czarnecki, John Quan, James Kirkpatrick, Raia Hadsell, Nicolas  
901 Heess, and Razvan Pascanu. Distral: Robust multitask reinforcement learning. *Advances in*  
902 *Neural Information Processing Systems*, 30, 2017.
- 903 Chen Tessler, Yoni Kasten, Yunrong Guo, Shie Mannor, Gal Chechik, and Xue Bin Peng. Calm:  
904 Conditional adversarial latent models for directable virtual characters. In *ACM SIGGRAPH 2023*  
905 *Conference Proceedings*, pp. 1–9, 2023.
- 906
- 907 Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in*  
908 *Neural Information Processing Systems*, 30, 2017.
- 909
- 910 Chen Wang, Claudia Pérez-D’Arpino, Danfei Xu, Li Fei-Fei, Karen Liu, and Silvio Savarese. Co-  
911 gail: Learning diverse strategies for human-robot collaboration. In *Conference on Robot Learn-*  
912 *ing*, pp. 1279–1290. PMLR, 2022.
- 913
- 914 Yawei Wang and Xiu Li. Reward function shape exploration in adversarial imitation learning: an  
915 empirical study. In *2021 IEEE International Conference on Artificial Intelligence and Computer*  
916 *Applications (ICAICA)*, pp. 52–57. IEEE, 2021.
- 917
- 918 Ziyu Wang, Josh S Merel, Scott E Reed, Nando de Freitas, Gregory Wayne, and Nicolas Heess.  
919 Robust imitation of diverse behaviors. *Advances in Neural Information Processing Systems*, 30,  
920 2017.

- 918 Paweł Wawrzyński. A cat-like robot real-time learning to run. In *Adaptive and Natural Computing*  
919 *Algorithms: 9th International Conference, Kuopio, Finland, Revised Selected Papers 9*, pp. 380–  
920 390. Springer, 2009.
- 921  
922 Annie Xie, Dylan Losey, Ryan Tolsma, Chelsea Finn, and Dorsa Sadigh. Learning latent representa-  
923 tions to influence multi-agent interaction. In *Conference on robot learning*, pp. 575–588. PMLR,  
924 2021.
- 925 Annie Xie, Lisa Lee, Ted Xiao, and Chelsea Finn. Decomposing the generalization gap in imitation  
926 learning for visual robotic manipulation. *arXiv preprint arXiv:2307.03659*, 2023.
- 927  
928 Mengdi Xu, Yikang Shen, Shun Zhang, Yuchen Lu, Ding Zhao, Joshua Tenenbaum, and Chuang  
929 Gan. Prompting decision transformer for few-shot policy generalization. In *International Con-*  
930 *ference on machine learning*, pp. 24631–24645. PMLR, 2022.
- 931 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey  
932 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.  
933 In *Conference on Robot Learning*, pp. 1094–1100. PMLR, 2020.
- 934  
935 Luyao Yuan, Xiaofeng Gao, Zilong Zheng, Mark Edmonds, Ying Nian Wu, Federico Rossano,  
936 Hongjing Lu, Yixin Zhu, and Song-Chun Zhu. In situ bidirectional human-robot value alignment.  
937 *Science Robotics*, 7(68):eabm4183, 2022.
- 938 Tom Zahavy, Yannick Schroecker, Feryal Behbahani, Kate Baumli, Sebastian Flennerhag, Shaobo  
939 Hou, and Satinder Singh. Discovering policies with domino: Diversity optimization maintaining  
940 near optimality. *arXiv preprint arXiv:2205.13521*, 2022.
- 941  
942 Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual  
943 manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
- 944  
945 Zeyu Zhu and Huijing Zhao. A survey of deep rl and il for autonomous driving policy learning.  
*IEEE Transactions on Intelligent Transportation Systems*, 23(9):14043–14065, 2021.
- 946  
947 Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse  
948 reinforcement learning. In *AAAI*, volume 8, pp. 1433–1438. Chicago, IL, USA, 2008.
- 949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

## A POINT MAZE

The PointMaze domain considered in Sec. 4 is presented in Fig. 8. PointMaze is a two-dimensional navigation environment with continuous state and action spaces. The state vector represents the agent’s current location’s x- and y- coordinates in  $[-\infty, \infty]^2$ . The action space is a velocity command, a two-dimensional vector in  $[-1, 1]^2$ . The episode length is fixed to 25 steps. Expert demonstrations are collected from a PD controller parameterized with a one-dimensional (1D) factor,  $\omega$  that determines the waypoint through which the agent passes  $(0, \omega)$  on its way to the goal  $(2, 0)$ .

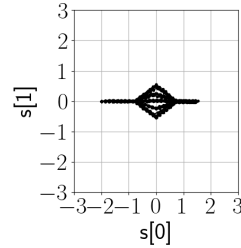


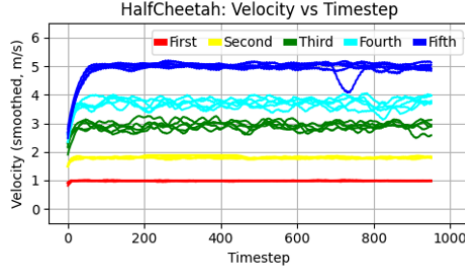
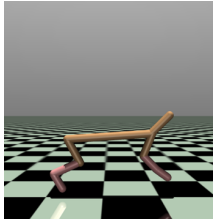
Figure 8: The figure visualizes expert demonstrations in PointMaze with varying waypoints along  $x = 0$ .

## B DOMAINS AND DEMONSTRATIONS

The bounded factor range is divided into 5 intervals for each domain, as explained in Sec. 6. For each interval, we add Gaussian noise to the mean value of the interval to generate five samples. We condition the expert policy on the five samples to obtain five demonstrations for each interval.

### B.1 HALFCHEETAH

The HalfCheetah environment considered in Sec. 6 is from the gym library (Brockman et al., 2016). The observation vector consists of the positions and velocities of the robot joints, along with height and velocities in the vertical and horizontal directions. The reward is modified as shown in Eq. 5, where  $r_t$  is the reward



at the time,  $t$ , and  $x_t$  is the position of the center of mass of the robot along the x-axis at time  $t$ , and  $I$ , the indicator function that outputs 1 if and only if (iff) its argument evaluates to true. The undiscounted episode return counts the number of steps in which the cheetah moves forward by a non-zero amount. The return is normalized using the range,  $[0, 1050]$ . The episode length is fixed at 1000 steps. The environment is stochastic with the robot initialized at random configurations.

$$r_{step} = I(x_{t+1} - x_t > 0) \tag{5}$$

The factor is the mean velocity measured as the net change in the x-coordinate over the elapsed time. Due to environment stochasticity, we use five sampled trajectories per conditioning latent vector during evaluation and consider the mean value. Demonstrations consist of the robot running at different mean velocities  $[1, 2, 3, 4, 5]$  m/s, collected using RL policies trained using SAC (Haarnoja et al., 2018) and auxiliary rewards for target velocities.

### B.2 DRIVELANESHIFT

The DriveLaneshift environment is built from the highway-env library (Leurent, 2018). The highway consists of two lanes. The scenario includes the ego-car in the right lane controlled by the agent, and another car in front, in the same lane, that maintains a constant speed of 25 m/s. The task of the ego-car is to shift to the left lane, overtake the other car, and reach the target speed of 30 m/s. The reward at each step is as shown in Eq. 6, where  $r_t$  is the reward at the time,  $t$ ,  $b_{onroad}$ , evaluates to true iff the car is within the road bounds at time  $t$ ,  $b_{safe}$  evaluates to true iff the ego-car has not crashed until time  $t$ ,  $b_{leftlane}$  evaluates to true iff the ego-car is in the left lane at time  $t$ ,  $v_t$  is the speed at time  $t$ , and  $clip(x, a, b)$  clips the value  $x$  to lie between  $a$  and  $b$ . The return is normalized

using the range  $[0, 175]$ . The state vector includes positions (absolute for the ego-car, relative for the other), velocities, heading angles, and longitudinal, latitudinal, and angular offsets to the closest lane for both cars. The episode length is fixed at 50 steps. The environment is deterministic.

$$r_{step} = I(b_{onroad}) + I(b_{safe}) + I(b_{leftlane}) + clip\left(\frac{|v_t - 30|}{5}, 0, 1\right) \quad (6)$$

The factor is the min headway distance, i.e., the distance between the ego-car and the other, at which the ego-car shifts to the left lane before overtaking. Demonstrations consist of the ego-car performing overtaking maneuvers at varying min headway distances  $[10.92, 18.28, 25.62, 32.91, 40.27]$  m, collected using a scripted PD controller.

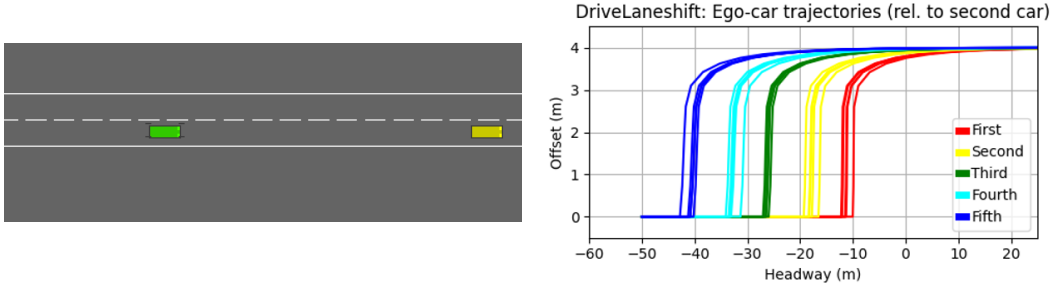


Figure 10: **Top:** The images visualize the highway overtaking scenario with the ego-car (green) starting behind the other car (yellow) in the right lane. **Bottom:** The figure visualizes the position of the ego-car relative to the other car as recorded in the demonstrations. The trajectories are colored to indicate the factor interval in which they belong.

### B.3 FETCHPICKPLACE

The FetchPickPlace environment considered is from the gym library (Brockman et al., 2016). The task is to move the object from its initial location on the table at  $[1.20, 0.75]$  to along the line  $x = 1.40$ , with reward at each step measured as shown in Eq. 7, where  $r_t$  is the reward at the time,  $t$ , and  $x_t$  is the position of the center of mass of the object along the x-axis at time  $t$ . The return is normalized using the range  $[-20, -5]$ . The state vector includes the end effector position and velocity, object position and velocity, finger gripper position and velocity, and object position relative to the gripper. The episode length is fixed at 100 steps. The environment is deterministic.

$$r_{step} = -|x_{t+1} - x_t| \quad (7)$$

The factor is the y-coordinate of the final object position. Demonstrations consist of the robot arm picking the object up from the initial location and placing it at the target x-coordinate and varying y-coordinates,  $0.75 + [-0.32, -0.16, 0, 0.16, 0.32]$  m, collected using a scripted state-based PD controller.

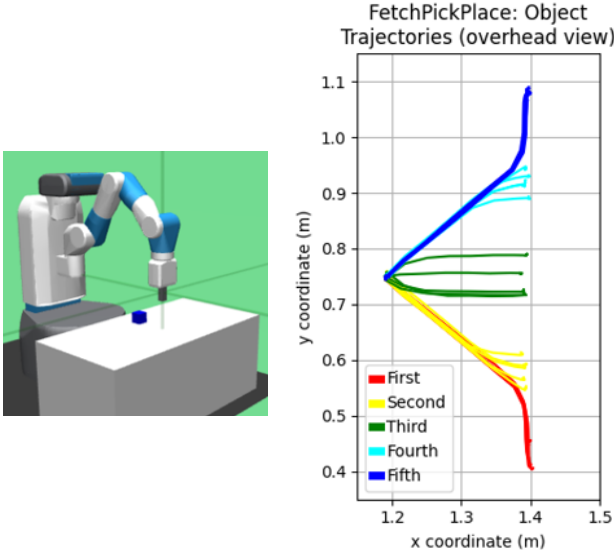


Figure 11: **Left:** The images visualize the Fetch robot with an object on the table. **Right:** The figure shows the object trajectories (from the top down) recorded in the demonstrations. The trajectories are colored to indicate the factor interval in which they belong.

## C ALGORITHM - DETAILED VERSION

A detailed version of Algorithm 1 can be found in Algorithm 2, with objectives and gradient steps for all components explicitly written down.



**Algorithm 2** Guided Strategy Discovery**Input:**  $\mathcal{D} = \{\tau_i^\xi\}$ **Output:**  $\pi$ 

- 1: Initialize policy  $\pi$ , task relevance  $f$ , factor-specific residual  $g$ , decoder  $q$ , with parameters  $\theta_\pi, \theta_f, \theta_g, \theta_q$ , Lagrange multiplier  $\lambda$ , bias  $b$ , and learning rates  $\eta_\pi, \eta_f, \eta_g, \eta_q, \eta_\lambda, \eta_b$
- 2: **for**  $i \in \{0, 1, 2, \dots\}$  epoch **do**
- 3:   Sample  $z^\pi$  from prior,  $\tau^\pi$  using policy  $\pi(\cdot|\cdot, z^\pi)$ ;  $\tau^\xi$  from  $\mathcal{D}$ , infer  $z^\xi$  using decoder  $q$
- 4:   Define objective for functions  $f, g$  and bias  $b$ :  

$$D(s, a, z) = \sigma(\lambda_S \cdot [f(s, a) + g(s, a, z)] + b)$$

$$J^I \leftarrow \mathbb{E}_{\tau^\xi} [\log D(s, a, z^\xi)] + \mathbb{E}_{\tau^\pi} [\log(1 - D(s, a, z^\pi))] - \mathbb{E}_{\tau^\xi} [(g(s, a, z^\xi))^2] - \mathbb{E}_{\tau^\pi} [(g(s, a, z^\pi))^2]$$
- 5:   Update  $f, g, b$  using gradients:  $[\theta_f, \theta_g, b] := [\theta_f, \theta_g, b] + [\eta_f \nabla_{\theta_f} J^I, \eta_g \nabla_{\theta_g} J^I, \eta_b \nabla_b J^I]$
- 6:   Define objective for decoder  $q$ :  

$$\delta(s, a, s', a') \leftarrow [\lambda_C \cdot \|s \oplus a - s' \oplus a'\| \cdot f(s', a') - \|\mu_{q(\cdot|s, a)} - \mu_{q(\cdot|s', a')}\|] \cdot f(s, a)$$

$$q_L(z, s, a) \leftarrow \mathcal{N}(z | \mu_{q(\cdot|s, a)}, \Sigma_{q(\cdot|s, a)})$$

$$J^E \leftarrow \mathbb{E}_{\tau^\pi} [\log q_L(z, s, a) + \lambda \cdot \min(\delta(s, a, s', a'), \epsilon)]$$
- 7:   Update decoder  $q$  and  $\lambda$  using gradients:  $[\theta_q, \lambda] := [\theta_q, \lambda] + [\eta_q \nabla_{\theta_q} J^E, -\eta_\lambda \nabla_\lambda J^E]$
- 8:   Update policy  $\pi$  with RL using rewards:  $r(s, a, z) = \log(D(s, a, z)) + \lambda_I \cdot \log q_L(z, s, a)$
- 9: **end for**

**D EVALUATION****D.1 IMPLEMENTATION**

We implement our approach on top of the public code-base for VILD (Tangkaratt et al., 2020) that implements adversarial IL algorithms using PyTorch (Imambi et al., 2021): [github.com/voot-t/vild\\_code](https://github.com/voot-t/vild_code). We use implementation tricks from their codebase to ensure convergence across methods, such as gradient penalty (Gulrajani et al., 2017) with a weight of 0.1 for the discriminator/task relevance, and the positive logarithmic function (Wang & Li, 2021) for discriminator rewards, i.e.,  $r(s, a) = -\log(1 - D(s, a))$  instead of  $r(s, a) = \log(D(s, a))$ .

We use a normal prior for the latent space across all approaches. The decoder  $q$  outputs the mean and diagonal elements of the covariance matrix of the approximate posterior distribution, which is assumed to be Gaussian.

**Conditioned Discriminator** To infer latent code  $\tau^\xi$  for a demonstration trajectory  $\tau^\xi$ , we make a simplifying assumption that the posterior distributions across transitions are independent. Thus, the product of the individual distributions gives us the posterior distribution for the demonstration trajectory.

We add expert transitions with mismatched demonstration latent vectors as fake samples to the discriminator dataset to ensure that the conditioned discriminator,  $D(s, a, z)$ , does not ignore the input latent vector. We upsample “real” data points in a batch to avoid imbalanced classes for discriminator gradient updates.

**Decoder Regularization** We perform spectral normalization (Miyato et al., 2018) using the PyTorch function `nn.utils.parametrizations.spectral_norm`. We scale the inputs to the decoder to implement Lipschitz constraint scaling with  $\lambda_S$ .

**D.2 DOMAIN-SPECIFIC VARIATIONS AND TUNING**

The hyperparameters used in our optimization are listed in Tables 1, 2. Each method is independently tuned for  $\lambda_I$  (and  $\lambda_C$  for Lipz, GSD) over the specified ranges, to maximize MAE over the test split for  $K=10$  over averaged over four rounds of evaluation and five train seeds. All hyperparameters omitted from the tables are set to default values from our base implementation.

**D.3 COMPARISON AGAINST OFFLINE IL APPROACHES**

We compare our approach against offline IL approaches that learn solely from data without any environment interaction to provide a comprehensive evaluation. We rely on the implementations

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153

Hyperparameter	Value
NN update minibatch size	256
Policy learning rate	3e-4
Entropy bonus	0.0001
Gamma	0.99
GAE coefficient (Schulman et al., 2015)	0.97
NN architectures	FCN
Policy activation	Tanh
BC warmstart epochs	10000
Disc. activation	Tanh
Disc. learning rate	1e-3
Disc. gradient steps	5
Dec. hidden dimensions	[100, 256]
Dec. activation	ReLU
Dec. gradient steps	5
$\lambda_S$	10
$b$ initial value	-5
Constraint slack ( $\epsilon$ )	1e-6
Lambda learning rate	1e-3
Optimizers	Adam

1154  
1155  
1156  
1157

Table 1: The table contains the list of hyperparameters, common across domains and generalization settings. GAE: Generalized Advantage Estimation, NN: Neural Network, FCN: Fully connected network, BC: Behavior cloning, Disc.: Discriminator, Dec.: Decoder

1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176

Hyperparam. \ Domain	HalfCheetah	DriveLaneshift	FetchPickPlace
Env. steps	1.5e7	0.5e7	1e7
RL algorithm	PPO	PPOBC*	PPOBC*
BC halflife, weight (PPOBC)	-	0.1, 0.2	0.1, 0.1
NN update interval (steps)	10000	1000	10000
NN hidden dimensions	[100, 100]	[100, 100]	[32, 32]
Observation norm. (w/ demos.)	False	False	True
Policy weight decay	(0, 1e-3)	1e-4	(1e-4, 5e-5)
Dec. learning rate	1e-3	1e-4	1e-3
Lambda initial value	(100, 500)	500	5000
Dec. gradient norm clip	$\infty$	25	$\infty$
Dec. rewards clip	$[-\infty, \infty]$	$[-20, 5]$	$[-\infty, \infty]$
Distillation objective weight	(0.02, 0.001)	(0.001, 0.001)	(0.0001, 0.0005)
$\lambda_S$ sweep list	[0.02, 0.05, 0.1, 0.2]	[0.1, 0.5, 1.0, 5.0]	[0.1, 0.5, 1.0, 5.0]
$\lambda_I$ sweep list	[(0.9, 0.8), [0.8, 0.7]]	[(0.99, 0.97), [0.99, 0.97]]	[(0.9, 0.8), [0.8, 0.7]]

1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

Table 2: The table contains the list of hyperparameters, specific to each domain (indicated by the column) and generalization setting (indicated by a 2-tuple (left, right), where left and right correspond to interpolation and extrapolation respectively). \*PPOBC (Jena et al., 2021) augments the policy objective with a behavior cloning (BC) loss term which improves learning stability without directly affecting the discriminator or decoder. We highlight that using the BC loss term comes at no additional human cost, as demonstrations are already available in the IL setting. Furthermore, we make no assumptions about the demonstrations’ behavior factors either and use the decoder network to infer the latent factor.

open-sourced by D3IL (Jia et al., 2024). To maintain focus on multimodal action distribution modeling, we exclude architectures that incorporate state histories or predict action sequences, such as action chunking (Zhao et al., 2023). We use fully connected neural network backbones with two

hidden layers each containing 100 units. Unless otherwise specified, we use default hyperparameters that are most common across tasks in D3IL. We use demonstrations corresponding to held-out factor intervals as the validation dataset for early stopping.

- Behavior cloning (BC): The NN takes state as input and outputs actions. The NN is trained with mean squared error (MSE) loss.
- BC with VAE (BC-VAE): We utilize a state-conditioned encoder-decoder setup to model the action distribution (Sohn et al., 2015). We use a latent space dimension of size 2, similar to our approach. We use a weight of 5.0 for scaling KL-divergence loss.
- Implicit BC (IBC): Florence et al. (2022) propose energy models that implicitly capture the action distribution at each state. The action is inferred by optimizing the energy function using Markov chain Monte-carlo sampling at each inference step.
- K-Means Discretization (BeT): Shafiqullah et al. (2022) propose an approach to capture multimodal action distributions using a learned discretization with  $K$  predicted action means and offsets. We use  $K = 64$ .
- Diffusion Policy (Diffusion): Chi et al. (2023); Pearce et al. (2023) propose modeling action distributions with a diffusion model conditioned on the state. We use timestep embeddings of size 16 and 24 denoising steps.

We show the recovery and performance trade-off of offline IL and our approach in Fig. 12. With regard to task performance (x-axis), for HalfCheetah (top row) and FetchPickPlace (bottom), GSD outperforms offline approaches (except for interpolation in HalfCheetah) indicated by the red cross being positioned further right than others. For DriveLaneShift (middle row), offline approaches other than BC and BC-VAE are competitive with GSD. The result suggests that in domains with complex dynamics like HalfCheetah and FetchPickPlace, environment interaction is necessary for task completion when learning from few demonstrations.

With regard to recovery performance (y-axis), for HalfCheetah (top row) and FetchPickPlace (bottom), GSD outperforms all offline approaches, indicated by the red cross being positioned below others. For interpolation in DriveLaneShift (middle row, left), approaches IBC, BeT and Diffusion are comparable to GSD. However, for extrapolation (middle row, right), GSD outperforms all methods. Poor performance of offline approaches in domains with complex dynamics like HalfCheetah and FetchPickPlace may be attributed to the absence of environment interaction. In simpler domains, IBC, BeT or Diffusion may be able to interpolate diverse behaviors. However, for extrapolation, environment interaction is necessary, even for simpler domains.

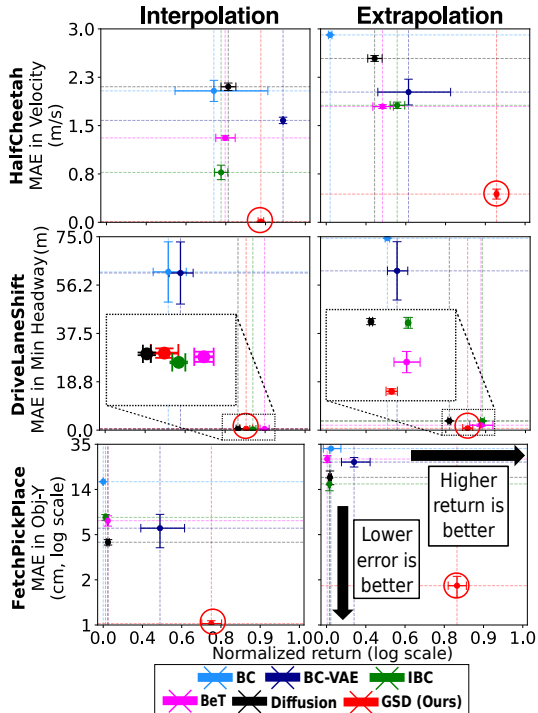


Figure 12: The figure shows the tradeoff between task and recovery performance for three domains and two splits for offline IL approaches and GSD (indicated with red circles). Error bars show standard errors over five seeds.

## E EVALUATION WITH HUMAN DEMONSTRATIONS IN TABLE TENNIS

### E.1 SETUP

**Demonstrations:** The WAM arm is enabled with gravity compensation for collecting kinesthetic demonstrations. Messages published to a ROS interface are collected for two seconds (starting after the ball is detected to move over the table) which is enough time to capture the return trajectory. Joint states ( $\mathbb{R}^7$ ) and ball positions ( $\mathbb{R}^3$ ) are matched over time, and concatenated to construct state

vectors ( $\mathbb{R}^{10}$ ). Action vectors are constructed by calculating the displacements at corresponding timesteps. We collect five demonstrations for each of the two-stroke types considered.

**Simulation:** We use position control for the WAM arm at 100hz with the control gains tuned to match real robot demonstrations visually when replaying action commands open loop. We tune ball flight parameters such that the ball flight paths (before being struck) visually match those from the real demonstrations when launched from a similar position, velocity, and noise as the ball launcher. We add Gaussian noise to the ball positions in the observation vector to mimic real recorded ball positions. We use an episode length of 200 steps that corresponds to a real-life execution period of two seconds.

## E.2 EVALUATION

We detect if the ball has been returned by checking if the velocity component along the long edge of the TT table has reversed. Once a returning ball is detected, we check if the ball remains above the table plane and within 10 cm beyond the sides of the table for the following 0.5 seconds. If the trajectory satisfies both criteria, we deem it a success. We calculate the factor value for each successful return which is the maximum height the ball reaches in the return trajectory.

We sample five trajectories per sampled latent vector due to the stochastic nature of the ball observations. However, not every sampled trajectory for a particular latent vector is guaranteed to succeed due to the stochasticity in the domain and optimization. Thus, we consider a latent vector successful if the ball is returned in at least three out of five trials. We consider the factor value for the successful latent vector to be the mean of the values of the successful trajectories.

For each method and train seed, we sample 200 latent vectors from the prior, and sample five trajectories per latent vector. We report the fraction of successful latent vectors to evaluate if behaviors can accomplish the underlying task. Among the set of latent vectors, we subsample 100 successful latent vectors (after ensuring each method has at least 50% success rate) and report the entropy (based on particle estimates (Singh et al., 2003)) among the calculated factors using the equation shown in Eq. 8 where  $V = \{v_i\}_{i=0}^M$  is the set of factor values  $v_i$ ,  $M = 100$ ,  $K = 50$  and  $\text{NNe}_{K,V}(v_i)$  returns the  $K$ th nearest neighbor to  $v_i$  from the set of values  $V$ . The entropy measure is up to a proportionality constant, as we use it to compare diversity achieved in return ball trajectory heights across methods.

$$H_K(V) = \frac{1}{M} \sum_{i=0}^M \log \|v_i - \text{NNe}_{K,V}(v_i)\| \quad (8)$$

## F GENERALITY OF $f$ -RELEVANT DIVERSITY

Our  $f$ -relevant diversity formulation discussed in Sec. 5.1 is designed to encourage behavior diversity with respect to any defined energy measure  $f$ . We briefly demonstrate the generality of our formulation in the simple 2D PointMaze domain with a user-defined energy function as shown in Fig. 13. Our formulation has a potential application in diverse solution discovery (Kumar et al., 2020; Osa et al., 2022), where a bounded form of the estimated  $Q$  function can be used as  $f$  to encourage diversity in high value regions.

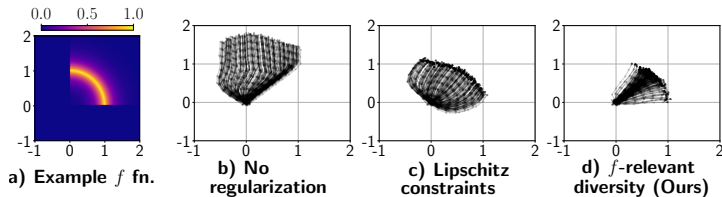


Figure 13: The figure visualizes behaviors in 2D PointMaze learned with a predefined energy function  $f$  shown in 13a.  $f$  is used to specify rewards in 13b-d and additionally formulate our diversity objective in 13d. With no regularization (13b) or Lipschitz constraints (13c), trajectories visit low  $f$ -energy regions of the state-space. However, with  $f$ -relevant diversity (ours, 13d), a larger portion of trajectories cover diverse high energy states.

## G FURTHER LIMITATIONS AND FUTURE WORK

Our work pertains to the realm of adversarial IL frameworks that employ diversity objectives in the form of MI. The generalization capabilities of other multimodal IL frameworks based on non-adversarial IL should be explored.

The scope of generalization considered in this work pertains to variations in the demonstrated expert **behavior factors**. Generalization to altered environment dynamics, adversarial perturbations, etc., could also be considered in the context of IL.

Our approach requires a task-relevance measure,  $f$ , which we derived from demonstration occupancies. IL approaches that do not explicitly model expert occupancy (Reddy et al., 2019; Garg et al., 2021) may not be readily compatible for integration with our regularization. However,  $Q$  functions learned during policy optimization may act as suitable alternatives for  $f$ . Our regularization is implemented through approximately enforced constraints using Lagrange multipliers. Approximate enforcement may permit spurious behaviors that drastically vary from demonstrations. Parametric approaches akin to spectral normalization for Lipschitz continuity (Miyato et al., 2018) are desirable.

We employ distillation to extract a task-relevance measure. Other approaches that learn generalizable reward functions (Szot et al., 2023) could also be explored. Our work pertains to online adversarial IL frameworks that employ diversity objectives in the form of MI. The efficacy of our constraints could be explored with offline and non-MI-based diversity frameworks.

## H EVALUATION RESULTS

We provide the numerical figures for recovery errors used to plot the graphs in Fig. 4 below. We further abbreviate Con, ConDist and Lipz as CO, CD and LZ respectively due to width constraints.

### HalfCheetah: Interpolation, Average:

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$0.258 \pm 0.007$	$0.147 \pm 0.006$	$0.106 \pm 0.004$	$0.080 \pm 0.004$	$0.063 \pm 0.004$
IG+LZ	$0.265 \pm 0.009$	$0.156 \pm 0.010$	$0.111 \pm 0.008$	$0.089 \pm 0.008$	$0.069 \pm 0.006$
IG+CO	$0.229 \pm 0.015$	$0.114 \pm 0.007$	$0.074 \pm 0.004$	$0.058 \pm 0.003$	$0.049 \pm 0.003$
IG+CD	$0.266 \pm 0.014$	$0.142 \pm 0.007$	$0.101 \pm 0.006$	$0.077 \pm 0.005$	$0.060 \pm 0.005$
IG+CD+LZ	$0.307 \pm 0.004$	$0.177 \pm 0.005$	$0.124 \pm 0.004$	$0.091 \pm 0.003$	$0.077 \pm 0.003$
GSD (Ours)	$0.226 \pm 0.007$	$0.120 \pm 0.003$	$0.081 \pm 0.003$	$0.065 \pm 0.001$	$0.052 \pm 0.002$

### HalfCheetah: Interpolation, Worst:

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$0.298 \pm 0.004$	$0.177 \pm 0.002$	$0.134 \pm 0.003$	$0.103 \pm 0.003$	$0.076 \pm 0.003$
IG+LZ	$0.348 \pm 0.020$	$0.212 \pm 0.021$	$0.159 \pm 0.019$	$0.131 \pm 0.017$	$0.099 \pm 0.013$
IG+CO	$0.291 \pm 0.029$	$0.142 \pm 0.014$	$0.092 \pm 0.007$	$0.072 \pm 0.007$	$0.061 \pm 0.005$
IG+CD	$0.339 \pm 0.024$	$0.181 \pm 0.013$	$0.133 \pm 0.011$	$0.102 \pm 0.009$	$0.082 \pm 0.008$
IG+CD+LZ	$0.373 \pm 0.007$	$0.223 \pm 0.010$	$0.158 \pm 0.008$	$0.117 \pm 0.006$	$0.100 \pm 0.006$
GSD (Ours)	$0.287 \pm 0.007$	$0.151 \pm 0.005$	$0.106 \pm 0.004$	$0.084 \pm 0.002$	$0.069 \pm 0.002$

### HalfCheetah: Extrapolation, Average:

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$0.841 \pm 0.010$	$0.766 \pm 0.014$	$0.737 \pm 0.016$	$0.710 \pm 0.017$	$0.697 \pm 0.019$
IG+LZ	$0.891 \pm 0.004$	$0.829 \pm 0.005$	$0.804 \pm 0.005$	$0.788 \pm 0.005$	$0.779 \pm 0.005$
IG+CO	$0.741 \pm 0.007$	$0.618 \pm 0.002$	$0.560 \pm 0.004$	$0.525 \pm 0.006$	$0.505 \pm 0.007$
IG+CD	$0.801 \pm 0.024$	$0.642 \pm 0.020$	$0.581 \pm 0.019$	$0.538 \pm 0.019$	$0.513 \pm 0.020$
IG+CD+LZ	$0.686 \pm 0.008$	$0.597 \pm 0.008$	$0.559 \pm 0.010$	$0.539 \pm 0.010$	$0.525 \pm 0.011$
GSD (Ours)	$0.682 \pm 0.028$	$0.556 \pm 0.032$	$0.512 \pm 0.033$	$0.484 \pm 0.034$	$0.472 \pm 0.034$

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

**HalfCheetah: Extrapolation, Worst:**

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$0.910 \pm 0.013$	$0.854 \pm 0.016$	$0.837 \pm 0.017$	$0.828 \pm 0.017$	$0.821 \pm 0.017$
IG+LZ	$0.999 \pm 0.004$	$0.909 \pm 0.007$	$0.880 \pm 0.009$	$0.869 \pm 0.009$	$0.860 \pm 0.010$
IG+CO	$0.875 \pm 0.019$	$0.779 \pm 0.021$	$0.725 \pm 0.024$	$0.685 \pm 0.027$	$0.661 \pm 0.030$
IG+CD	$0.963 \pm 0.031$	$0.750 \pm 0.017$	$0.686 \pm 0.012$	$0.656 \pm 0.008$	$0.634 \pm 0.007$
IG+CD+LZ	$0.797 \pm 0.012$	$0.719 \pm 0.011$	$0.699 \pm 0.011$	$0.689 \pm 0.011$	$0.682 \pm 0.011$
GSD (Ours)	$0.808 \pm 0.039$	$0.679 \pm 0.038$	$0.646 \pm 0.040$	$0.632 \pm 0.041$	$0.621 \pm 0.041$

**DriveLaneshift: Interpolation, Average:**

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$3.435 \pm 0.148$	$2.305 \pm 0.132$	$1.808 \pm 0.122$	$1.489 \pm 0.105$	$1.277 \pm 0.098$
IG+LZ	$3.193 \pm 0.096$	$2.101 \pm 0.097$	$1.631 \pm 0.101$	$1.345 \pm 0.078$	$1.111 \pm 0.073$
IG+CO	$5.246 \pm 0.500$	$3.153 \pm 0.375$	$2.377 \pm 0.333$	$1.936 \pm 0.302$	$1.712 \pm 0.278$
IG+CD	$3.619 \pm 0.200$	$2.246 \pm 0.151$	$1.749 \pm 0.154$	$1.419 \pm 0.163$	$1.221 \pm 0.154$
IG+CD+LZ	$2.785 \pm 0.196$	$1.652 \pm 0.172$	$1.302 \pm 0.156$	$1.114 \pm 0.142$	$0.985 \pm 0.128$
GSD (Ours)	$2.343 \pm 0.134$	$1.299 \pm 0.094$	$0.877 \pm 0.061$	$0.685 \pm 0.058$	$0.530 \pm 0.045$

**DriveLaneshift: Interpolation, Worst:**

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$4.528 \pm 0.262$	$3.294 \pm 0.226$	$2.693 \pm 0.228$	$2.280 \pm 0.193$	$2.008 \pm 0.190$
IG+LZ	$4.702 \pm 0.179$	$3.315 \pm 0.187$	$2.675 \pm 0.192$	$2.253 \pm 0.152$	$1.859 \pm 0.140$
IG+CO	$6.843 \pm 0.690$	$4.304 \pm 0.507$	$3.329 \pm 0.441$	$2.754 \pm 0.397$	$2.469 \pm 0.380$
IG+CD	$5.057 \pm 0.307$	$3.329 \pm 0.249$	$2.657 \pm 0.269$	$2.171 \pm 0.282$	$1.878 \pm 0.271$
IG+CD+LZ	$3.895 \pm 0.314$	$2.581 \pm 0.310$	$2.079 \pm 0.291$	$1.797 \pm 0.269$	$1.637 \pm 0.241$
GSD (Ours)	$3.009 \pm 0.246$	$1.709 \pm 0.171$	$1.138 \pm 0.104$	$0.911 \pm 0.106$	$0.687 \pm 0.079$

**DriveLaneshift: Extrapolation, Average:**

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$5.807 \pm 0.231$	$4.889 \pm 0.269$	$4.446 \pm 0.288$	$4.144 \pm 0.297$	$3.943 \pm 0.302$
IG+LZ	$6.435 \pm 0.124$	$5.692 \pm 0.193$	$5.375 \pm 0.224$	$5.134 \pm 0.244$	$4.977 \pm 0.257$
IG+CO	$6.347 \pm 0.769$	$4.335 \pm 0.584$	$3.588 \pm 0.567$	$3.166 \pm 0.542$	$2.967 \pm 0.521$
IG+CD	$3.902 \pm 0.185$	$2.480 \pm 0.185$	$1.886 \pm 0.172$	$1.559 \pm 0.166$	$1.399 \pm 0.168$
IG+CD+LZ	$5.588 \pm 0.511$	$3.809 \pm 0.400$	$2.902 \pm 0.369$	$2.472 \pm 0.362$	$2.206 \pm 0.369$
GSD (Ours)	$2.803 \pm 0.108$	$1.545 \pm 0.074$	$1.019 \pm 0.053$	$0.815 \pm 0.046$	$0.695 \pm 0.048$

**DriveLaneshift: Extrapolation, Worst:**

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$6.386 \pm 0.219$	$5.527 \pm 0.246$	$5.102 \pm 0.268$	$4.805 \pm 0.289$	$4.624 \pm 0.291$
IG+LZ	$6.908 \pm 0.125$	$6.222 \pm 0.224$	$6.071 \pm 0.244$	$5.923 \pm 0.259$	$5.836 \pm 0.269$
IG+CO	$8.599 \pm 1.302$	$6.673 \pm 1.140$	$5.827 \pm 1.107$	$5.324 \pm 1.069$	$5.033 \pm 1.026$
IG+CD	$5.326 \pm 0.338$	$3.606 \pm 0.341$	$2.907 \pm 0.339$	$2.438 \pm 0.322$	$2.206 \pm 0.328$
IG+CD+LZ	$8.723 \pm 1.105$	$6.473 \pm 0.890$	$5.038 \pm 0.790$	$4.340 \pm 0.764$	$3.945 \pm 0.766$
GSD (Ours)	$4.067 \pm 0.225$	$2.283 \pm 0.156$	$1.534 \pm 0.105$	$1.238 \pm 0.091$	$1.034 \pm 0.094$

**FetchPickPlace: Interpolation, Average:**

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$0.071 \pm 0.004$	$0.045 \pm 0.004$	$0.034 \pm 0.003$	$0.027 \pm 0.003$	$0.022 \pm 0.002$
IG+LZ	$0.044 \pm 0.002$	$0.024 \pm 0.001$	$0.016 \pm 0.001$	$0.012 \pm 0.001$	$0.010 \pm 0.001$
IG+CO	$0.072 \pm 0.009$	$0.053 \pm 0.010$	$0.046 \pm 0.010$	$0.041 \pm 0.011$	$0.040 \pm 0.011$
IG+CD	$0.083 \pm 0.011$	$0.067 \pm 0.012$	$0.060 \pm 0.012$	$0.055 \pm 0.012$	$0.052 \pm 0.012$
IG+CD+LZ	$0.050 \pm 0.001$	$0.027 \pm 0.001$	$0.019 \pm 0.001$	$0.015 \pm 0.001$	$0.012 \pm 0.001$
GSD (Ours)	$0.037 \pm 0.001$	$0.020 \pm 0.001$	$0.014 \pm 0.000$	$0.011 \pm 0.000$	$0.008 \pm 0.000$



1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

**FetchPickPlace: Interpolation, Worst:**

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$0.084 \pm 0.005$	$0.058 \pm 0.005$	$0.046 \pm 0.005$	$0.037 \pm 0.004$	$0.031 \pm 0.004$
IG+LZ	$0.056 \pm 0.003$	$0.032 \pm 0.002$	$0.021 \pm 0.002$	$0.016 \pm 0.001$	$0.013 \pm 0.001$
IG+CO	$0.081 \pm 0.009$	$0.061 \pm 0.011$	$0.054 \pm 0.011$	$0.048 \pm 0.011$	$0.047 \pm 0.011$
IG+CD	$0.088 \pm 0.010$	$0.072 \pm 0.012$	$0.064 \pm 0.012$	$0.059 \pm 0.012$	$0.056 \pm 0.012$
IG+CD+LZ	$0.053 \pm 0.002$	$0.030 \pm 0.001$	$0.021 \pm 0.001$	$0.017 \pm 0.001$	$0.014 \pm 0.001$
GSD (Ours)	$0.042 \pm 0.001$	$0.022 \pm 0.001$	$0.015 \pm 0.001$	$0.013 \pm 0.001$	$0.009 \pm 0.000$

**FetchPickPlace: Extrapolation, Average:**

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$0.137 \pm 0.003$	$0.103 \pm 0.003$	$0.087 \pm 0.003$	$0.075 \pm 0.003$	$0.068 \pm 0.003$
IG+LZ	$0.083 \pm 0.003$	$0.057 \pm 0.003$	$0.045 \pm 0.003$	$0.038 \pm 0.003$	$0.034 \pm 0.003$
IG+CO	$0.210 \pm 0.016$	$0.183 \pm 0.018$	$0.167 \pm 0.018$	$0.158 \pm 0.018$	$0.152 \pm 0.018$
IG+CD	$0.264 \pm 0.012$	$0.246 \pm 0.013$	$0.233 \pm 0.014$	$0.225 \pm 0.014$	$0.220 \pm 0.015$
IG+CD+LZ	$0.078 \pm 0.003$	$0.046 \pm 0.002$	$0.032 \pm 0.002$	$0.026 \pm 0.002$	$0.022 \pm 0.002$
GSD (Ours)	$0.068 \pm 0.002$	$0.037 \pm 0.002$	$0.027 \pm 0.002$	$0.021 \pm 0.002$	$0.018 \pm 0.002$

**FetchPickPlace: Extrapolation, Worst:**

Model	$K = 10$	$K = 20$	$K = 30$	$K = 40$	$K = 50$
IG	$0.182 \pm 0.004$	$0.143 \pm 0.004$	$0.120 \pm 0.003$	$0.105 \pm 0.003$	$0.097 \pm 0.003$
IG+LZ	$0.097 \pm 0.003$	$0.071 \pm 0.004$	$0.057 \pm 0.004$	$0.049 \pm 0.004$	$0.044 \pm 0.004$
IG+CO	$0.249 \pm 0.012$	$0.223 \pm 0.015$	$0.208 \pm 0.017$	$0.196 \pm 0.017$	$0.191 \pm 0.018$
IG+CD	$0.303 \pm 0.003$	$0.286 \pm 0.006$	$0.275 \pm 0.008$	$0.266 \pm 0.009$	$0.260 \pm 0.010$
IG+CD+LZ	$0.092 \pm 0.003$	$0.057 \pm 0.003$	$0.042 \pm 0.004$	$0.035 \pm 0.004$	$0.031 \pm 0.004$
GSD (Ours)	$0.094 \pm 0.005$	$0.056 \pm 0.005$	$0.041 \pm 0.004$	$0.033 \pm 0.003$	$0.028 \pm 0.004$