

WATCH LESS, DO MORE: IMPLICIT SKILL DISCOVERY FOR VIDEO-CONDITIONED POLICY

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we study the problem of video-conditioned policy learning. While previous works mostly focus on learning policies that perform a single skill specified by the given video, we take a step further and aim to learn a policy that can perform multiple skills according to the given video, and generalize to unseen videos by recombining these skills. To solve this problem, we propose our algorithm, **Watch-Less-Do-More**, an information bottleneck-based imitation learning framework for implicit skill discovery and video-conditioned policy learning. In our method, an information bottleneck objective is employed to control the information contained in the video representation, ensuring that it only encodes information relevant to the current skill (**Watch-Less**). By discovering potential skills from training videos, the learned policy is able to recombine them and generalize to unseen videos to achieve compositional generalization (**Do-More**). To evaluate our method, we perform extensive experiments in various environments and show that our algorithm substantially outperforms baselines (up to **2x**) in terms of compositional generalization ability.

1 INTRODUCTION

As large language models (LLMs) have demonstrated remarkable zero-shot and few-shot generalization abilities (Brown et al., 2020; Ouyang et al., 2022), the research focus of decision-making policies has also shifted from addressing a specific task, such as mastering an environment via reinforcement learning (Sutton, 2018) or replicating a dataset via imitation learning (Hussein et al., 2017), to completing diverse tasks based on given instructions. These instructions can be treated as goals for the decision-making models, and encompass modalities such as text (Nair et al., 2022; Carta et al., 2023), goal image (Yadav et al., 2023b;a), or future state (Cui et al., 2022; Lee et al., 2024). To achieve such goal-conditioned policies, a variety of methods have been proposed and achieved great success across multiple domains (Liu et al., 2022). However, the aforementioned goal specifications often overlook dynamic information, such as the ordering of task completion or the method of task completion (if there are many). In contrast, video offers a natural way to represent these details, thereby leading to a line of research exploring video-conditioned policy learning (Eze & Crick, 2024b).

Existing methods for video-conditioned policy learning have been applied to various scenarios, including robotic manipulation (Chane-Sane et al., 2023; Shin et al., 2023; Jiang et al., 2023), navigation (Zhou et al., 2024), and autonomous driving (Shin et al., 2024). Taking different videos as input, the learned policy can be deployed to perform different skills to solve corresponding tasks. However, these methods often consider only the video demonstration of a single task (Chane-Sane et al., 2023) and only the object-level generalization (Jiang et al., 2023). In real-world applications, we often want the learned policy to perform a set of different skills to achieve a combination of multiple tasks.

When the video demonstration of a task combination is given, an ideal policy should directly perform skills as demonstrated in the given video. To train such a policy, researchers have explored skill-based imitation learning methods (Xu et al., 2023; Wang et al., 2023; Shin et al., 2023; 2024). However, these methods often require explicit video segmentation annotations, or videos of another embodiment to train a skill-based policy, which greatly increases the difficulty of the data collection process. Therefore, in this paper, we consider whether it is possible to learn a video-conditioned

054 policy that can perform multiple skills without these requirements. Moreover, as we consider videos
 055 of different task combinations, we also expect the learned policy to achieve compositional general-
 056 ization (Lin et al., 2023), that is, it can still perform well on task combinations that have not been
 057 during training.

058 To fulfill such an expectation, we propose our algorithm, **Watch-Less-Do-More (WL-DM)**, an in-
 059 formation bottleneck-based imitation learning framework for implicit skill discovery and video-
 060 conditioned policy learning. For a video-conditioned policy, the given video can be considered as a
 061 sequence of tasks. As the policy can only work on one task at a time, it should be able to perform
 062 well by focusing only on the current task, instead of the entire task sequence. Based on this intuition,
 063 WL-DM employs the information bottleneck method (Tishby et al., 2000) to control the information
 064 contained in the video representation. This is accomplished by 1) minimizing the mutual informa-
 065 tion (Cover, 1999) between video and its representation to reduce the information contained in the
 066 video representation and 2) maximizing the mutual information between video representation and
 067 the current skill to preserve enough information related to the current task. To better understand
 068 the effect of this information bottleneck method, we further build a theoretical connection between
 069 the proposed method and the intuition behind it. Using this method, WL-DM makes the learned
 070 policy only to consider the current task, which achieves the *implicit* video segmentation without
 071 requiring explicit video segmentation annotations. The advantage of considering only the current
 072 task can be related to the compositional generalization ability. When an unseen video is given, the
 073 video-conditioned policy learned by WL-DM can implicitly decompose the unseen video into seen
 074 tasks, and perform corresponding skills, thus facilitating the compositional generalization ability
 075 of the learned video-conditioned policy. To further validate our algorithm, we propose a practical
 076 implementation of our method and conduct various empirical evaluations across diverse environ-
 077 ments. The experimental results indicate that WL-DM achieves substantially better compositional
 generalization ability than baselines, demonstrating the effectiveness of our method.

078 Our contributions can be summarized as follows:

- 079 • We propose our method, Watch-Less-Do-More (WL-DM), an information bottleneck-
 080 based imitation learning framework for implicit skill discovery and video-conditioned pol-
 081 icy learning, where two different mutual information terms work together to ensure the
 082 video representation contains only information related to the current task.
- 083 • The intuition behind WL-DM is that the optimal policy should behave similarly when con-
 084 ditioned on all tasks and when conditioned on only the current task. To better explain our
 085 method, we further build a theoretical connection between WL-DM and this intuition.
- 086 • We propose a practical implementation of our algorithm and perform empirical evaluations
 087 in Frank Kitchen (Gupta et al., 2020) and Meta world (Yu et al., 2020) to demonstrate the
 088 effectiveness of WL-DM. The experimental results indicate that WL-DM achieves (up to
 089 **2x**) better compositional generalization ability compared to baselines.

092 2 RELATED WORK

094 2.1 LEARNING FROM VIDEOS

095 Using massive Internet data to train language models has been proven to be successful and has re-
 096 sulted in a trend of research on large language models (Brown et al., 2020; Touvron et al., 2023).
 097 Inspired by this success, researchers have begun to pay attention to another type of data wildly avail-
 098 able on the Internet, video data, and produced a series of studies on learning from videos (McCarthy
 099 et al., 2024; Eze & Crick, 2024a). For decision-making models, video data can be used in various
 100 ways, such as reward function learning (Escontrela et al., 2023; Sermanet et al., 2018; Chen et al.,
 101 2021a), dynamic model learning (Baker et al., 2022), representation learning (Nair et al., 2023),
 102 and policy learning (Jang et al., 2022; Jiang et al., 2023; Chane-Sane et al., 2023; Shin et al., 2023;
 103 2024). Our paper belongs to the last category, that is, using video demonstrations as instructions to
 104 learn a video-conditioned policy. It is worth noting that previous work in this category often focuses
 105 only on demonstration videos containing a single task (Chane-Sane et al., 2023), or requires aligned
 106 data of other modalities (Jang et al., 2022; Shin et al., 2023; 2024). This can be attributed to the
 107 lack of clear goal labels in demonstration videos (McCarthy et al., 2024). Therefore, when dealing
 with videos containing multiple tasks, we often need to introduce information in other modalities

108 to provide segmentation annotations for the video, to distinguish the tasks to be completed at each
109 stage (Shin et al., 2023; 2024). Unlike previous work, in this paper, we attempt to directly learn a
110 video-conditioned policy capable of handling videos containing multiple tasks, without introducing
111 additional segmentation annotations.

112 2.2 ONE-SHOT IMITATION LEARNING

113 One-shot imitation learning was originally introduced in Duan et al. (2017), where the goal of this
114 problem is to learn a policy that can quickly adapt to a new task given a single demonstration. For
115 one-shot imitation learning, we can achieve it through different learning methods such as meta-
116 learning (Duan et al., 2017; Finn et al., 2017), semi-supervised learning (Wu et al., 2024), and
117 imitation learning (Jang et al., 2022; Cui et al., 2022; Jiang et al., 2023). Specifically, our method
118 falls into the last category: we assume the existence of an imitation learning dataset paired with
119 video demonstrations, such that we can use this dataset to train a video-conditioned policy.

120 One-shot demonstrations can be presented in various formats, such as trajectories (Cui et al., 2022;
121 Lee et al., 2024), videos (Dasari & Gupta, 2021; Jain et al., 2024; Wang et al., 2023; Xu et al., 2023;
122 Chane-Sane et al., 2023), multimodal information (Jiang et al., 2023; Shin et al., 2023; 2024), etc. In
123 this paper, we consider adapting to new tasks through video demonstrations, that is, one-shot video
124 imitation learning. In previous work, video demonstrations often only include a single task, and the
125 adaptation to new tasks mainly focuses on differences at the embodiment and object level (locations,
126 textures, etc.) (Dasari & Gupta, 2021; Mandi et al., 2022; Chane-Sane et al., 2023). Unlike these
127 studies, we consider video demonstrations containing multiple tasks and focus on adaptation at the
128 level of task combination. For this setting, previous work generally assumes the existence of data
129 corresponding to another embodiment (Wang et al., 2023; Xu et al., 2023) or assumes information
130 in other modalities to provide video segmentation annotations (Shin et al., 2023; 2024). Unlike
131 these works, we do not assume additional data and learn a video-conditioned policy that can finish
132 multiple tasks solely through the information contained in the videos.

133 2.3 COMPOSITIONAL GENERALIZATION

134 Due to the compositional nature of natural language, most previous work considers the composi-
135 tional generalization problem over language instructions. For example, Oh et al. (2017) proposed a
136 method based on hierarchical reinforcement learning that enables the policy to generalize to unseen
137 command combinations and longer command sequences at test time. Stengel-Eskin et al. (2022)
138 combined the transformer model and the masking mechanism to obtain generalization over object
139 combinations. The attention mechanism for compositional generalization was further investigated
140 by Spilsbury & Ilin (2022), and a method utilizing sparse factored attention for goal identifica-
141 tion was proposed. Modular architecture is another way to induce compositional generalization.
142 Carvalho et al. (2023) proposed modular successor features to enhance the compositional general-
143 ization ability, and Logeswaran et al. (2023) directly considered an additive decomposition of the
144 state-action value function to obtain the generalization ability over language instructions.

145 Unlike these studies, we consider the generalization across different task combinations based on
146 video demonstrations. During training, we only have access to a subset of task combinations and
147 their corresponding video demonstrations. Our goal is to enable the policy to decompose different
148 tasks from the videos and acquire skills to solve these tasks. At test time, the policy is expected to
149 reproduce an unseen video demonstration by combining a set of skills learned in the training set.
150 This setting has been studied by Wang et al. (2023); Xu et al. (2023); Shin et al. (2023; 2024). How-
151 ever, Wang et al. (2023); Xu et al. (2023) focused on the cross-embodiment scenario, thus requiring
152 video data from another embodiment, and Shin et al. (2023; 2024) required language information to
153 provide segmentation annotations for videos. Unlike them, our method incorporates an information
154 bottleneck-based objective to achieve implicit video segmentation and skill discovery, without the
155 need for other sources of information.

156 3 PROBLEM FORMULATION

157 In this paper, we consider the video-conditioned policy learning problem. This problem can be
158 formulated as a special case of the goal-conditioned Markov Decision Process (MDP) (Nasiriany

et al., 2019) and defined by a tuple $\langle \mathcal{S}, \mathcal{G}, \mathcal{A}, P, R, \rho_0, \gamma \rangle$. Similar to the general MDP, \mathcal{S} is the set of states, \mathcal{A} is the set of actions, $P(s_{t+1}|s_t, a_t)$ is the transition probability, ρ_0 is the initial state distribution and γ is the discount factor. Additionally, we have \mathcal{G} as the set of goals, which will also affect the reward function $R(s_t, a_t, g)$. For a goal-conditioned policy $\pi(a_t|s_t, g)$ with a given goal g , we want it to maximize the following objective:

$$\mathcal{J}(\pi) = \mathbb{E}_{s_0 \sim \rho_0, a \sim \pi, s' \sim P} \left[\sum_t \gamma^t r(s_t, a_t, g) \right].$$

As we focus on the video-conditioned policy learning, we assume our goals to be videos $\mathcal{G} = \mathcal{V}$, such that a goal-conditioned policy $\pi(a_t|s_t, g)$ becomes a video-conditioned policy $\pi(a_t|s_t, v)$. Moreover, we consider the case where each video $v = (k_0, \dots, k_N)$ contains multiple tasks $k \in \mathcal{T}$, where N is the number of tasks and \mathcal{T} is the set of all possible tasks. To evaluate the compositional generalization ability of the video-conditioned policy, we assume two video sets $\mathcal{V}_{\text{train}}$ and $\mathcal{V}_{\text{test}}$, such that there is no overlapping between the train video set and test video set $\mathcal{V}_{\text{train}} \cap \mathcal{V}_{\text{test}} = \emptyset$ and both video sets contain all possible tasks $\bigcup_{v \in \mathcal{V}_{\text{train}}, k \in v} k = \bigcup_{v \in \mathcal{V}_{\text{test}}, k \in v} k = \mathcal{T}$. The video-conditioned policy will be trained in $\mathcal{V}_{\text{train}}$ to maximize $\mathbb{E}_{v \sim \mathcal{P}_{\text{train}}} \mathcal{J}(\pi)$ and will be tested in $\mathcal{V}_{\text{test}}$ in terms of $\mathbb{E}_{v \sim \mathcal{P}_{\text{test}}} \mathcal{J}(\pi)$, where $\mathcal{P}_{\text{train}}$ and $\mathcal{P}_{\text{test}}$ are uniform distributions across $\mathcal{V}_{\text{train}}$ and $\mathcal{V}_{\text{test}}$ respectively.

4 METHOD

In this section, we introduce our method, Watch-Less-Do-More (WL-DM). The intuition behind our method is that we want the video-conditioned policy to make decisions relying not on the entire video, but only on information related to the current task, thereby achieving implicit video segmentation and skill discovery. To achieve this, we propose an information bottleneck-based objective and theoretically establish the connection between this objective and our intuition. By decomposing training videos into a combination of different skills, the video-conditioned policy can handle unseen videos by recombining these skills to complete the required task combinations demonstrated in the unseen video.

4.1 INTUITION: FOCUSING ON THE CURRENT TASK

As formulated in Section 3, we assume that each video v contains N tasks $[k_0, \dots, k_N]$ that need to be completed and the completion of these tasks is independent. In this case, we further assume a training set $\mathcal{D} = \{\tau_i, v_i\}$, where $\tau_i = (s_0, a_0^*, \dots, s_T, a_T^*)_i$ is the expert trajectory corresponding to the video $v_i = (f_0, \dots, f_T)$ and f_i is the video frame at each timestep. Given such a dataset, we can easily learn a video-conditioned policy $\pi(a_t|s_t, v)$ through imitation learning (Hussein et al., 2017) that can complete different task combinations given different videos, at least within the coverage of the training set. For example, the policy can be trained via the following behavior-cloning loss:

$$\begin{aligned} \mathcal{L}_{\text{BC}}(\theta, \phi) &= -\mathbb{E}_{s_t, a_t^*, v \sim \mathcal{D}} \left[\log \pi(a_t^*|s_t, v) \right] \\ &= -\mathbb{E}_{s_t, a_t^*, v \sim \mathcal{D}} \left[\log f_\theta(a_t^*|s_t, g_\phi(v)) \right], \end{aligned} \quad (1)$$

where g_ϕ is the video encoder and f_θ is the action decoder.

A potential problem with this training method is that, when the size of our training set is limited, the learned policy can easily overfit (Ying, 2019) to videos in the training set. This problem causes the learned policy to focus too much on the details of these videos to distinguish them completely, while ignoring the fact that these videos are composed through elements of the same task set. In such a case, when an unseen video is given, i.e., an unseen combination of tasks, the performance of the learned policy may decrease dramatically due to the overfitting issue. To address this problem, we need to focus on the fact that all videos are composed through elements of the same task set \mathcal{T} . Even for those unseen videos, although the corresponding task combinations are not included in the training set, each task that constitutes them has already been covered in the training set. Therefore, if we can decompose videos into individual tasks and train the policy based on the decomposed tasks, such that $\pi(a_t|s_t, v) = \pi(a_t|s_t, v_{\text{cur}})$, where v_{cur} is the video segment corresponding to k_{cur} and $v = (v_{\text{cur}}, v_{\text{other}})$, the policy can then handle unseen videos as all the skills demonstrated in the

216 videos have been covered and trained in the training set, which is commonly known as compositional
 217 generalization (Lin et al., 2023). However, such a task-level video segmentation annotation could be
 218 inaccessible in many cases. In this paper, we do not assume this kind of annotations as in previous
 219 work (Shin et al., 2023; 2024). To achieve a similar effect, we propose an information bottleneck-
 220 based method, allowing the policy to implicitly decompose demonstration videos, enabling it to rely
 221 only on the information related to the current skill when making decisions, thereby achieving the
 222 compositional generalization ability.

223 4.2 INFORMATION BOTTLENECK FOR VIDEO-CONDITIONED POLICY LEARNING

224 As described in Section 4.1, when all tasks can be completed independently, a video-conditioned
 225 policy $\pi(a_t|s_t, v) = \pi(a_t|s_t, v_{\text{cur}})$ can achieve the compositional generalization ability. However,
 226 a video contains information about not only the current task k_{cur} but also other tasks k_{other} . There-
 227 fore, we need an additional objective to train the video encoder g_ϕ , such that it produces a similar
 228 representation for v and v_{cur} , and we can then ensure $\pi(a_t|s_t, v) = \pi(a_t|s_t, v_{\text{cur}})$. To achieve this,
 229 we need to reduce the mutual information between video representations h_v and the video segments
 230 of other tasks v_{other} , and the reason can be seen from the following theorem:

231 **Theorem 1.** *If we have $\text{MI}(h_v; v_{\text{other}} | s, v_{\text{cur}}) = 0$, then $D_{\text{KL}}(\pi(a|s, v) || \pi(a|s, v_{\text{cur}})) = 0$ for all*
 232 *state-video pairs $(s, v) \in \mathcal{S} \times \mathcal{V}$ with non-zero probability $P(s, v) > 0$.*

233 *Proof.* See Appendix B. □

234 This theorem suggests the necessity of reducing the information of other tasks contained in the
 235 video representation. However, since we do not assume any video segmentation annotation, we
 236 cannot directly obtain segments corresponding to the current task and other tasks from the video, and
 237 therefore cannot directly manipulate the mutual information. Hence, we use a constructive method
 238 to manipulate the information in the video representation indirectly. Specifically, we first minimize
 239 the mutual information between the video representation h_v and the entire video v to minimize the
 240 amount of information contained in the video representation. At the same time, we maximize the
 241 mutual information between the video representation and some approximation of the current skill
 242 (which will be discussed later). Since different skills are required for performing different tasks,
 243 we can in this way indirectly ensure that the video representation still retains a certain amount of
 244 information about the current task. Putting these two terms together, we can construct the following
 245 objective, which is often referred to as the information bottleneck (Tishby et al., 2000):

$$246 \mathcal{L}_{IB} = \text{MI}(h_v; v | s) - \alpha \text{MI}(h_v; z | s), \quad (2)$$

247 where z is an approximated representation of the current skill, and α is the coefficient for the trade-
 248 off between two mutual information terms. As discussed in Tishby & Zaslavsky (2015), the infor-
 249 mation bottleneck is often used to learn a compact representation, which in our case is to dismiss
 250 the irrelevant part k_{other} and retain the relevant part k_{cur} . In the following two sections, we discuss
 251 how to compute this objective in practice.

252 4.3 MINIMIZING MUTUAL INFORMATION WITH VIDEO

253 The first term in Equation (2) is to minimize the mutual information between video representation
 254 h_v and the entire video v . By expanding this term, we have:

$$255 \text{MI}(h_v; v | s) = \mathbb{E}_{P(s)} \mathbb{E}_{P(v)} \left[D_{\text{KL}}(g_\phi(h_v | s, v) || P(h_v | s)) \right], \quad (3)$$

256 where $P(s)$ and $P(v)$ represent the state and video distribution, respectively. $P(h_v | s)$ is a marginal
 257 distribution $P(h_v | s) = \mathbb{E}_{P(v)} g_\phi(h_v | s, v)$. As estimating this marginal distribution could be in-
 258 tractable in practice, previous work (Goyal et al., 2018; Eysenbach et al., 2021) commonly ap-
 259 proximates it with some prior $g(h_v | s)$. As we can see from Equation (3), the goal of this objec-
 260 tive is to minimize the distance between the video representation produced by the video encoder
 261 $h_v = g_\phi(h_v | s, v)$ and some prior $g(h_v | s)$ that does not consider video v at all. As shown later in
 262 the experiment, such a target for distance minimization is undesirable as it induces too much loss of
 263 video information. To solve this problem, we need to find a better alternative for $g(h_v | s)$ such that
 264 the loss of video information can be controlled at a proper level.

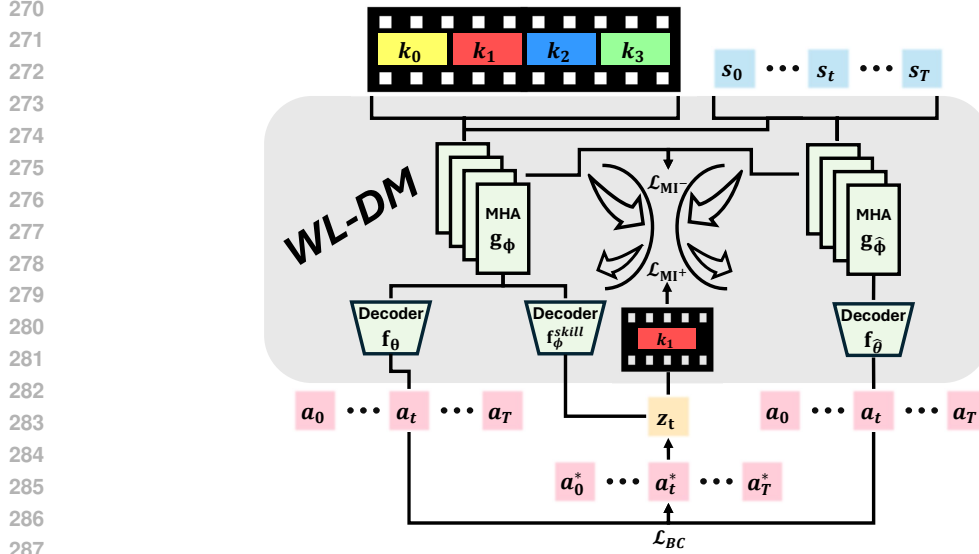


Figure 1: Overall Framework of WL-DM. We introduce an information bottleneck-based objective to achieve implicit video segmentation and skill discovery. Blocks with different colors represent different tasks. MHA stands for Multi-head Attention.

Recall the intuition in Section 4.1, we want the representation to be only related to the video segment of the current task v_{cur} . Although we cannot access the precise v_{cur} , we do have access to the video $v = (f_0, \dots, f_T)$, which allows us to approximate v_{cur} using a future video segment $\tilde{v}_{\text{cur}} = (f_t, f_{t+1}, \dots, f_{t+L})$ for state s_t , where L is a randomly sampled window size. Therefore, we can use \tilde{v}_{cur} as the input and get a prior video encoder $g_\phi^{\text{prior}}(h_v|s, \tilde{v}_{\text{cur}})$, which leads to the minimization of the following equation:

$$\mathbb{E}_{P(s)} \mathbb{E}_{P(v)} \left[D_{\text{KL}}(g_\phi(h_v|s, v) || g_\phi^{\text{prior}}(h_v|s, \tilde{v}_{\text{cur}})) \right]. \quad (4)$$

The relationship between Equation (3) and Equation (4) can be seen from the following inequality:

$$\begin{aligned} & \mathbb{E}_{P(s)} \mathbb{E}_{P(v)} \left[D_{\text{KL}}(g_\phi(h_v|s, v) || g_\phi^{\text{prior}}(h_v|s, \tilde{v}_{\text{cur}})) - D_{\text{KL}}(g_\phi(h_v|s, v) || P(h_v|s)) \right] \\ &= \mathbb{E}_{P(s)} \mathbb{E}_{P(v)} \mathbb{E}_{g_\phi(h_v|s, v)} \left[\log \frac{P(h_v|s)}{g_\phi^{\text{prior}}(h_v|s, \tilde{v}_{\text{cur}})} \right] \\ &= \mathbb{E}_{P(s)} \mathbb{E}_{P(h_v|s)} \left[\log \frac{P(h_v|s)}{g_\phi^{\text{prior}}(h_v|s, \tilde{v}_{\text{cur}})} \right] \\ &= \mathbb{E}_{P(s)} \left[D_{\text{KL}}(P(h_v|s) || g_\phi^{\text{prior}}(h_v|s, \tilde{v}_{\text{cur}})) \right] \\ &\geq 0, \end{aligned}$$

which indicates that, using prior encoder $g_\phi^{\text{prior}}(h_v|s, \tilde{v}_{\text{cur}})$, we construct an upper bound of Equation (3). With this prior encoder, we can get the final objective for mutual information minimization:

$$\mathcal{L}_{\text{MI}} = \mathbb{E}_{s, v \sim D} \left[D_{\text{KL}}(g_\phi(h_v|s, v) || g_\phi^{\text{prior}}(h_v|s, \tilde{v}_{\text{cur}})) \right]. \quad (5)$$

In addition to Equation (5), the prior encoder g_ϕ^{prior} is also trained via behavior cloning similar to Equation (1) with another action decoder $f_{\hat{\theta}}$ attached after it.

4.4 MAXIMIZING MUTUAL INFORMATION WITH SKILL APPROXIMATION

Another term in Equation (2) is to maximize the mutual information between the video representation h_v and some skill approximation z . We follow Yuan et al. (2024) and use the short-term

behavior $z = (a_t, a_{t+1}, \dots, a_{t+M})$ as the representation of skills for state s_t , where M is a randomly sampled window size. To enhance the level of abstraction of the skill representation, we further propose to first cluster all actions in the training dataset \mathcal{D} and then use the cluster id x_t of each action to improve the skill representation, such that $z = (x_t, x_{t+1}, \dots, x_{t+M})$. As the mutual information is to measure the dependency between two variables, to maximize $\text{MI}(h_v; z)$, we can simply maximize $\log P(z|h_v)$. As we have $z = (x_t, x_{t+1}, \dots, x_{t+M})$, similar to Yuan et al. (2024), we can decompose the above maximization into each timestep and get the final objective for mutual information maximization:

$$\mathcal{L}_{\text{MI}^+} = -\mathbb{E}_{s_t, x_t, v \sim \mathcal{D}} \left[\log f_{\psi}^{\text{skill}}(x_t | s_t, g_{\phi}(v)) \right], \quad (6)$$

where we introduce the skill decoder f_{ψ}^{skill} to enhance the dependency between h_t and z .

4.5 SUMMARY

Putting Equations (1), (5) and (6) together, we can now have the total loss for WL-DM:

$$\mathcal{L}_{\text{WL-DM}} = \mathcal{L}_{\text{BC}} + \alpha_1 \mathcal{L}_{\text{MI}^-} + \alpha_2 \mathcal{L}_{\text{MI}^+},$$

where we have two coefficients α_1 and α_2 to balance the scale of these three terms. The overall framework of our algorithm is illustrated in Figure 1. We use multiple self-attention layers as the encoder g_{ϕ} to process video tokens and state tokens and then use the action decoder f_{θ} to predict action labels a_t^* . The joint optimization of $\mathcal{L}_{\text{MI}^-}$ and $\mathcal{L}_{\text{MI}^+}$ ensures the video representation contains only information related to the current task. The pseudocode of our algorithm is summarized in Algorithm 1. It is worth noting that the skill decoder f_{ψ}^{skill} , the prior video encoder g_{ϕ}^{prior} , and the prior action decoder $f_{\bar{\theta}}$ will only be used during training, we will keep only the video encoder g_{ϕ} and the action decoder f_{θ} for execution.

Algorithm 1 WL-DM

- 1: Initialize video encoder g_{ϕ} , action decoder f_{θ} and skill decoder f_{ψ}^{skill}
 - 2: Initialize prior video encoder g_{ϕ}^{prior} and prior action decoder $f_{\bar{\theta}}$
 - 3: Initialize training dataset \mathcal{D}
 - 4: **for** $i = 1$ **to** I **do**
 - 5: Sample data (s_t, a_t^*, v) from \mathcal{D}
 - 6: Construct approximation of current video segment v_{cur}
 - 7: Construct approximation of current skill z
 - 8: Update g_{ϕ} and f_{θ} by Equation (1) with (s_t, a_t^*, v)
 - 9: Update g_{ϕ}^{prior} and $f_{\bar{\theta}}$ by Equation (1) with $(s_t, a_t^*, v_{\text{cur}})$
 - 10: Update g_{ϕ} and g_{ϕ}^{prior} by Equation (5) with (s_t, v, v_{cur})
 - 11: Update g_{ϕ} and f_{ψ}^{skill} by Equation (6) with (s_t, z, v)
 - 12: **end for**
-

5 EXPERIMENTS

5.1 EXPERIMENT SETUP

To validate our method, we conduct empirical evaluations on two different robotic environments, Franka Kitchen (Gupta et al., 2020) and Meta World (Yu et al., 2020). The visualization of these two environments is presented in Figure 2.

In Franka Kitchen (FK), we control a Franka Panda robot in the kitchen environment to perform seven possible tasks: microwave (M), kettle (K), bottom burner (B), top burner (T), light switch (L), slide cabinet (S) and hinge cabinet (H). The dataset from the original paper (Gupta et al., 2020) contains 566 trajectories corresponding to 24 different task combinations. To enable video-conditioned policy training, we train expert policy using the original dataset to collect trajectories and corresponding video demonstrations. To evaluate the one-shot imitation learning ability, we split

the dataset into a training dataset and a test dataset, where the training dataset contains 17 different task combinations and the test dataset contains 7 different task combinations, and there is no overlap of task combinations between the training set and the test set. During testing, we sample 3 different video demonstrations for each task combination, run the evaluation 10 times, and report the average performance.

In Meta world (MW), we modify the original environment (Yu et al., 2020) to perform multiple tasks within a single episode. In this newly devised environment, we control a Sawyer robot to perform four possible tasks: close drawer (D), open door (O), push button (B) and open window (W). We use expert policy to collect the dataset for all 24 different tasks. It is worth noting that, as the dataset contains all possible task combinations, the task orders presented in the video demonstration bring additional difficulties for policy learning. To evaluate the one-shot imitation learning ability, we split the dataset into a training dataset and a test dataset, where the training dataset contains 17 different task combinations and the test dataset contains 7 different task combinations, and there is no overlap of task combinations between the training set and the test set. During testing, we sample 3 different video demonstrations for each task combination, run the evaluation 10 times, and report the average performance.

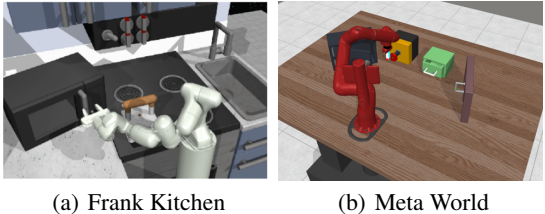


Figure 2: Visualization of Experiment Environments

We include several challenging imitation learning algorithms as our baselines: C-bet (Cui et al., 2022), decision transformer (Chen et al., 2021b), and VIMA (Jiang et al., 2023). As C-bet and decision transformer were not proposed for video-conditioned policy learning, we modify them to additionally take videos as input and get baselines **V-BET** and **V-DT**. For VIMA, it was originally proposed for multimodal prompts. However, as we do not assume data of other modalities, we train VIMA on our video-only dataset and serve as our baseline **VIMA**. More details of experiments can be found in Appendix A.

5.2 EXPERIMENT: MAIN

Table 1: The performance of WL-DM and baselines on all FK and MW tasks.

Env	Methods	Tasks							Avg
		MBTH	MBLS	MBTL	KBTH	MTLH	KBLS	KBTS	
FK	WL-DM	2.30	2.57	2.37	2.10	1.83	1.97	3.17	2.33 \pm 0.78
	V-BET	1.37	2.47	0.83	2.27	1.73	1.80	3.10	1.94 \pm 1.06
	V-DT	1.00	2.33	1.70	1.33	1.47	1.70	2.10	1.66 \pm 0.79
	VIMA	0.77	0.50	0.50	1.27	0.20	1.67	1.70	0.94 \pm 0.85
MW		ODWB	DOBW	DBWO	WBOD	BDOW	BDWO	BWDO	
	WL-DM	3.33	2.00	2.00	2.00	2.67	2.00	4.00	2.57 \pm 0.90
	V-BET	1.87	2.00	0.73	1.33	0.33	0.00	1.97	1.18 \pm 0.85
	V-DT	1.33	2.13	1.23	1.93	0.37	0.83	0.83	1.24 \pm 0.86
	VIMA	1.80	1.00	0.37	0.37	1.17	0.57	0.83	0.87 \pm 0.84

As shown in Table 1, our method achieves better average performance in both environments. In Franka Kitchen, our method achieves an improvement of approximately 20.1% compared to the best baseline (V-BET). In Meta World, the improvement is even more significant, with our method achieving a 100.1% improvement compared to the best baseline (V-DT). Although our method consistently outperforms all baselines, we note that there is a large gap in terms of the degree of improvement between the two environments. This is because, in Franka Kitchen, task combinations in the dataset are not diverse enough, and there is no variation in terms of the task orders (A, B vs B,

A). Therefore, it can be considered to have a strong correlation within the task combinations. Therefore, even without considering the segmentation of tasks in video demonstrations, we can still utilize this correlation to achieve a policy that performs well during testing. However, in Meta World, our dataset includes all task combinations and considers different task orders, making task segmentation in videos even more critical, which explains the gap in improvement of our algorithm in the two environments. More specifically, out of a total of 14 test tasks, our method achieves the best performance in 12 of them. Such overall performance validates the effectiveness of our algorithm.

For the performance of baselines, we found that V-BET and V-DT perform at a similar level. This is because the main difference between V-BET and V-DT in our implementation is whether or not action is used as part of the trajectory. For the robotic environments we used, this action information can often be directly inferred from changes in the state of the robot. Therefore, the advantage of using action information is not significant. VIMA, on the other hand, does not perform well in both environments. One potential reason is that VIMA was originally proposed for multimodal scenarios. Although its training process can be transferred to the pure video scenario, this direct transfer is clearly not effective. Moreover, in terms of one-shot video imitation, VIMA mainly considers object-level variations in a single task rather than variations in task combinations, so we believe its performance decline is acceptable.

5.3 EXPERIMENT: ABLATION

Table 2: The performance of WL-DM and ablation baselines on all MW tasks.

Env	Methods	Tasks						Avg	
		ODWB	DOBW	DBWO	WBOD	BDOW	BDWO		BWDO
MW	WL-DM	3.33	2.00	2.00	2.00	2.67	2.00	4.00	2.57 \pm 0.90
	WL-DM w/o \mathcal{L}_{MI+}	1.00	2.00	1.83	1.67	1.67	1.67	2.00	1.69 \pm 0.46
	WL-DM w/o \mathcal{L}_{MI-}	2.00	2.00	2.00	2.00	1.67	2.00	3.33	2.14 \pm 0.64
	WL-DM w \emptyset	2.00	2.13	1.93	2.00	2.00	1.33	2.00	1.91 \pm 0.37

As our objective function contains two different mutual information terms, we conduct ablation studies in this section to verify the contribution of these two components to our method. We construct two ablation baselines, WL-DM w/o \mathcal{L}_{MI+} and WL-DM w/o \mathcal{L}_{MI-} . The ablation baselines are identical to our algorithm in all aspects, except that WL-DM w/o \mathcal{L}_{MI+} does not use Equation (6) and WL-DM w/o \mathcal{L}_{MI-} does not use Equation (5). We validate these two ablation baselines in the Meta World environment and compare them with our method. As shown in Table 2, both ablation baselines achieve worse performance compared to WL-DM, thereby verifying that both Equation (5) and Equation (6) contribute to our algorithm. Notably, even with only Equation (5), the ablation baseline still achieves better performance than the baselines in Section 5.2, which further validates the effectiveness of Theorem 1 in practice.

As mentioned in Section 4.3, directly minimizing the mutual information using prior $g(h_v|s)$ can lead to excessive loss of video information in practice, thereby affecting the performance of the algorithm. To verify this point, we construct another ablation baseline WL-DM w \emptyset . This baseline is again identical to our method in all aspects, except that it does not use $g_{\phi}^{\text{prior}}(h_v|s, \tilde{v}_{\text{cur}})$ but instead uses prior $g(h_v|s)$. From Table 2, we can see that using $g(h_v|s)$ indeed leads to a decline in the performance, thus verifying the upper bound we constructed in Section 4.3. Additionally, it is worth noting that WL-DM w/o \mathcal{L}_{MI-} demonstrates that when we do not minimize mutual information at all, that is, do not control the information provided by the video, the algorithm cannot achieve its best performance. In contrast, WL-DM w \emptyset indicates that excessively reducing the information of the video also leads to a decline in the final performance. This phenomenon further demonstrates the importance of a proper approximation for v_{cur} and validates our statements in Section 4.3.

6 CONCLUSION

In this paper, we investigate the problem of one-shot video imitation for video-conditioned policies. To enhance the compositional generalization ability of the learned policy, we propose an imitation learning framework, **Watch-Less-Do-More (WL-DM)**. Our method introduces an information bottleneck-based objective, which leads to implicit skill discovery for video-conditioned policies. The intuition behind this method is that by segmenting the video into different tasks, the policy learns diverse skills corresponding to these tasks. When faced with unseen videos, the policy can also decompose them into combinations of previously encountered tasks, thereby completing these tasks through the learned skills. To better explain our method, we build a theoretical connection between our method and this intuition using information theory. We also present a practical implementation of our algorithm and evaluate it on a variety of tasks across multiple environments. The experimental results indicate that our algorithm outperforms baselines in terms of the compositional generalization ability, which verifies the effectiveness of our algorithm.

7 LIMITATIONS AND FUTURE WORK

The limitation of this work is that our approximations of the current task k_{curr} and the current skill z remain naive, which are just the information from the immediate future. Although similar approximations have been used in many previous studies (Pertsch et al., 2021; Liu et al., 2021; Xie et al., 2023; Yuan et al., 2024), one issue with this approach is the need for video data to be aligned with trajectory data in timesteps (for a state s_t , we can access its corresponding video frame f_t), which can be costly to collect in many cases.

Our future work will primarily focus on extending our algorithm to real-world scenarios, such as real-world robots, thereby broadening the application scope of our algorithm. Additionally, we will consider extending our algorithm to multimodal scenarios, utilizing multimodal information to obtain better approximations of tasks and skills, thereby not only enhancing the performance of the algorithm but also expanding the range of instruction formats it can process.

REFERENCES

- Bowen Baker, Ilge Akkaya, Peter Zhokhov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt) learning to act by watching unlabeled online videos. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, pp. 24639–24654, 2022.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pp. 1877–1901, 2020.
- Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves Oudeyer. Grounding large language models in interactive environments with online reinforcement learning. In *International Conference on Machine Learning*, pp. 3676–3713. PMLR, 2023.
- Wilka Carvalho, Angelos Filos, Richard L Lewis, Satinder Singh, et al. Composing task knowledge with modular successor feature approximators. *arXiv preprint arXiv:2301.12305*, 2023.
- Elliot Chane-Sane, Cordelia Schmid, and Ivan Laptev. Learning video-conditioned policies for unseen manipulation tasks. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 909–916. IEEE, 2023.
- Annie S Chen, Suraj Nair, and Chelsea Finn. Learning generalizable robotic reward functions from “in-the-wild” human videos. *arXiv preprint arXiv:2103.16817*, 2021a.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: reinforcement learning via sequence modeling. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, pp. 15084–15097, 2021b.

- 540 Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- 541
- 542 Zichen Jeff Cui, Yibin Wang, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. From
543 play to policy: Conditional behavior generation from uncurated robot data. *arXiv preprint*
544 *arXiv:2210.10047*, 2022.
- 545 Sudeep Dasari and Abhinav Gupta. Transformers for one-shot visual imitation. In *Conference on*
546 *Robot Learning*, pp. 2071–2084. PMLR, 2021.
- 547 Yan Duan, Marcin Andrychowicz, Bradly Stadie, Jonathan Ho, Jonas Schneider, Ilya Sutskever,
548 Pieter Abbeel, and Wojciech Zaremba. One-shot imitation learning. In *Proceedings of the 31st*
549 *International Conference on Neural Information Processing Systems*, pp. 1087–1098, 2017.
- 550 Alejandro Escontrela, Adcmi Adcniji, Wilson Yan, Ajay Jain, Xue Bin Peng, Ken Goldberg, Young-
551 woon Lee, Danijar Hafner, and Pieter Abbeel. Video prediction models as rewards for reinforce-
552 ment learning. In *Proceedings of the 37th International Conference on Neural Information Pro-*
553 *cessing Systems*, pp. 68760–68783, 2023.
- 554 Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. Robust predictable control. In
555 *Proceedings of the 35th International Conference on Neural Information Processing Systems*, pp.
556 27813–27825, 2021.
- 557
- 558 Chrisantus Eze and Christopher Crick. Learning by watching: A review of video-based learning
559 approaches for robot manipulation. *arXiv preprint arXiv:2402.07127*, 2024a.
- 560
- 561 Chrisantus Eze and Christopher Crick. Learning by watching: A review of video-based learning
562 approaches for robot manipulation. *arXiv preprint arXiv:2402.07127*, 2024b.
- 563
- 564 Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imi-
565 tation learning via meta-learning. In *Conference on robot learning*, pp. 357–368. PMLR, 2017.
- 566 Anirudh Goyal, Riashat Islam, DJ Strouse, Zafarali Ahmed, Hugo Larochelle, Matthew Botvinick,
567 Yoshua Bengio, and Sergey Levine. Infobot: Transfer and exploration via the information bottle-
568 neck. In *International Conference on Learning Representations*, 2018.
- 569 Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy
570 learning: Solving long-horizon tasks via imitation and reinforcement learning. In *Conference on*
571 *Robot Learning*, pp. 1025–1037. PMLR, 2020.
- 572
- 573 Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. Imitation learning: A
574 survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.
- 575 Vidhi Jain, Maria Attarian, Nikhil J Joshi, Ayzaan Wahid, Danny Driess, Quan Vuong, Pannag R
576 Sanketi, Pierre Sermanet, Stefan Welker, Christine Chan, et al. Vid2robot: End-to-end video-
577 conditioned policy learning with cross-attention transformers. *arXiv preprint arXiv:2403.12943*,
578 2024.
- 579 Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine,
580 and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In *Confer-*
581 *ence on Robot Learning*, pp. 991–1002. PMLR, 2022.
- 582
- 583 Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-
584 Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: robot manipulation with multimodal
585 prompts. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 14975–
586 15022, 2023.
- 587 Seungjae Lee, Yibin Wang, Haritheja Etukuru, H Jin Kim, Nur Muhammad Mahi Shafiullah, and
588 Lerrel Pinto. Behavior generation with latent actions. *arXiv preprint arXiv:2403.03181*, 2024.
- 589
- 590 Baihan Lin, Djallel Bouneffouf, and Irina Rish. A survey on compositional generalization in appli-
591 cations. *arXiv preprint arXiv:2302.01067*, 2023.
- 592 Bo Liu, Qiang Liu, Peter Stone, Animesh Garg, Yuke Zhu, and Anima Anandkumar. Coach-player
593 multi-agent reinforcement learning for dynamic team composition. In *International Conference*
on Machine Learning, pp. 6860–6870. PMLR, 2021.

- 594 Minghuan Liu, Menghui Zhu, and Weinan Zhang. Goal-conditioned reinforcement learning: Prob-
595 lems and solutions. *arXiv preprint arXiv:2201.08299*, 2022.
- 596
- 597 Lajanugen Logeswaran, Wilka Carvalho, and Honglak Lee. Learning compositional tasks from lan-
598 guage instructions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 13300–
599 13308, 2023.
- 600 Zhao Mandi, Fangchen Liu, Kimin Lee, and Pieter Abbeel. Towards more generalizable one-shot
601 visual imitation learning. In *2022 International Conference on Robotics and Automation (ICRA)*,
602 pp. 2434–2444. IEEE, 2022.
- 603
- 604 Robert McCarthy, Daniel CH Tan, Dominik Schmidt, Fernando Acero, Nathan Herr, Yilun Du,
605 Thomas G Thuruthel, and Zhibin Li. Towards generalist robot learning from internet video: A
606 survey. *arXiv preprint arXiv:2404.19664*, 2024.
- 607 Suraj Nair, Eric Mitchell, Kevin Chen, Silvio Savarese, Chelsea Finn, et al. Learning language-
608 conditioned robot behavior from offline data and crowd-sourced annotation. In *Conference on*
609 *Robot Learning*, pp. 1303–1315. PMLR, 2022.
- 610
- 611 Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A univer-
612 sal visual representation for robot manipulation. In *Conference on Robot Learning*, pp. 892–909.
613 PMLR, 2023.
- 614 Soroush Nasiriany, Vitchyr H Pong, Steven Lin, and Sergey Levine. Planning with goal-conditioned
615 policies. In *Proceedings of the 33rd International Conference on Neural Information Processing*
616 *Systems*, pp. 14843–14854, 2019.
- 617
- 618 Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generalization with
619 multi-task deep reinforcement learning. In *International Conference on Machine Learning*, pp.
620 2661–2670. PMLR, 2017.
- 621 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong
622 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
623 instructions with human feedback. In *Proceedings of the 36th International Conference on Neural*
624 *Information Processing Systems*, pp. 27730–27744, 2022.
- 625
- 626 Karl Pertsch, Youngwoon Lee, and Joseph Lim. Accelerating reinforcement learning with learned
627 skill priors. In *Conference on robot learning*, pp. 188–204. PMLR, 2021.
- 628 Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey
629 Levine, and Google Brain. Time-contrastive networks: Self-supervised learning from video. In
630 *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 1134–1141. IEEE,
631 2018.
- 632
- 633 Sangwoo Shin, Daehee Lee, Minjong Yoo, Woo Kyung Kim, and Honguk Woo. One-shot imitation
634 in a non-stationary environment via multi-modal skill. In *Proceedings of the 40th International*
635 *Conference on Machine Learning*, pp. 31562–31578, 2023.
- 636 Sangwoo Shin, Minjong Yoo, Jeongwoo Lee, and Honguk Woo. Semtra: A semantic skill translator
637 for cross-domain zero-shot policy adaptation. In *Proceedings of the AAAI Conference on Artificial*
638 *Intelligence*, pp. 15000–15008, 2024.
- 639
- 640 Sam Spilsbury and Alexander Ilin. Compositional generalization in grounded language learning via
641 induced model sparsity. In *Proceedings of the 2022 Conference of the North American Chap-*
642 *ter of the Association for Computational Linguistics: Human Language Technologies: Student*
643 *Research Workshop*, pp. 143–155, 2022.
- 644 Elias Stengel-Eskin, Andrew Hundt, Zhuohong He, Aditya Murali, Nakul Gopalan, Matthew Gom-
645 bolay, and Gregory Hager. Guiding multi-step rearrangement tasks with natural language instruc-
646 tions. In *Conference on Robot Learning*, pp. 1486–1501. PMLR, 2022.
- 647
- Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.

- 648 Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In
649 *2015 IEEE Information Theory Workshop (ITW)*, pp. 1–5. IEEE, 2015.
- 650
- 651 Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv*
652 *preprint physics/0004057*, 2000.
- 653 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
654 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
655 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 656
- 657 Chen Wang, Linxi Fan, Jiankai Sun, Ruohan Zhang, Li Fei-Fei, Danfei Xu, Yuke Zhu, and Anima
658 Anandkumar. Mimicplay: Long-horizon imitation learning by watching human play. In *Confer-*
659 *ence on Robot Learning*, pp. 201–221. PMLR, 2023.
- 660 Philipp Wu, Kourosh Hakhmaneshi, Yuqing Du, Igor Mordatch, Aravind Rajeswaran, and Pieter
661 Abbeel. Semi-supervised one-shot imitation learning. *arXiv preprint arXiv:2408.05285*, 2024.
- 662
- 663 Zhihui Xie, Zichuan Lin, Deheng Ye, Qiang Fu, Yang Wei, and Shuai Li. Future-conditioned unsu-
664 pervised pretraining for decision transformer. In *International Conference on Machine Learning*,
665 pp. 38187–38203. PMLR, 2023.
- 666 Mengda Xu, Zhenjia Xu, Cheng Chi, Manuela Veloso, and Shuran Song. Xskill: Cross embodiment
667 skill discovery. In *Conference on Robot Learning*, pp. 3536–3555. PMLR, 2023.
- 668
- 669 Karmesh Yadav, Arjun Majumdar, Ram Ramrakhya, Naoki Yokoyama, Alexei Baevski, Zsolt Kira,
670 Oleksandr Maksymets, and Dhruv Batra. Ovrl-v2: A simple state-of-art baseline for imagenav
671 and objectnav. *arXiv preprint arXiv:2303.07798*, 2023a.
- 672 Karmesh Yadav, Ram Ramrakhya, Arjun Majumdar, Vincent-Pierre Berges, Sachit Kuhar, Dhruv
673 Batra, Alexei Baevski, and Oleksandr Maksymets. Offline visual representation learning for em-
674 bodied navigation. In *Workshop on Reincarnating Reinforcement Learning at ICLR 2023*, 2023b.
- 675
- 676 Xue Ying. An overview of overfitting and its solutions. In *Journal of physics: Conference series*,
677 volume 1168, pp. 022022. IOP Publishing, 2019.
- 678 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
679 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
680 In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.
- 681 Haoqi Yuan, Zhancun Mu, Feiyang Xie, and Zongqing Lu. Pre-training goal-based models for
682 sample-efficient reinforcement learning. In *The Twelfth International Conference on Learning*
683 *Representations*, 2024.
- 684
- 685 Bohan Zhou, Jiangxing Wang, and Zongqing Lu. Nolo: Navigate only look once. *arXiv preprint*
686 *arXiv:2408.01384*, 2024.
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

A IMPLEMENTATION DETAILS

We implement our algorithm and all baselines based on the codebase of C-bet (Cui et al., 2022). For WL-DM and V-BET, we only consider only observations in the trajectory, and for V-DT and VIMA, we consider both observations and actions in the trajectory, which aligns with the implementation stated in the paper of C-bet (Cui et al., 2022), DT Chen et al. (2021b) and VIMA (Jiang et al., 2023). For WL-DM, V-BET, and V-DT, we use the same transformer model as stated in C-bet, which contains multiple self-attention layers to process video information and trajectory information at the same time. For VIMA, we use alternating cross-attention and self-attention layers as described in its paper (Jiang et al., 2023).

For all experiments, we set the learning rate to be 3×10^{-4} and set the window size for the trajectory to be 20 (for V-DT and VIMA, it means 20 observation-action pairs). For WL-DM, the window size of future video segments is sampled from $[20, 40]$. As we use the codebase of C-bet, all methods use the same action decoder, where we set the number of bins for action discretization to 32, and the id of each cluster will also be used for the representation of skills for WL-DM. For the Franka Kitchen environment (Gupta et al., 2020), we use decoders with 3 layers, and 3 heads and set the hidden dimension to be 60 (for VIMA, it means in total 3 self-attention layers and 3 cross-attention layers). We train all methods for 10 epochs. For WL-DM, α_1 is fixed to be 1×10^{-2} and α_2 is fixed to be 1×10^{-1} during the training process. For the Meta World environment (Yu et al., 2020), we use decoders with 6 layers, and 6 heads and set the hidden dimension to be 120 (for VIMA, it means in total 6 self-attention layers and 6 cross-attention layers). We train all methods for 30 epochs. For WL-DM, α_1 is set to be 0 in the beginning and fixed to be 1×10^{-3} after 10 epochs, and α_2 is fixed to be 10 during the training process.

B PROOF OF THEOREM 1

Theorem 1. *If we have $\text{MI}(h_v; v_{\text{other}} | s, v_{\text{cur}}) = 0$, then $D_{\text{KL}}(\pi(a|s, v) || \pi(a|s, v_{\text{cur}})) = 0$ for all state-video pairs $(s, v) \in \mathcal{S} \times \mathcal{V}$ with non-zero probability $P(s, v) > 0$.*

Proof. By expanding the mutual information $\text{MI}(v_{\text{other}}; h_v, a | s, v_{\text{cur}})$, we can have the following equality:

$$\begin{aligned}
 & \text{MI}(v_{\text{other}}; h_v, a | s, v_{\text{cur}}) \\
 &= \mathbb{E}_{P(s, v_{\text{cur}})} \mathbb{E}_{P(v_{\text{other}}, h_v, a | s, v_{\text{cur}})} \left[\log \frac{P(v_{\text{other}}, h_v, a | s, v_{\text{cur}})}{P(v_{\text{other}} | s, v_{\text{cur}}) P(h_v, a | s, v_{\text{cur}})} \right] \\
 &= \mathbb{E}_{P(s, v_{\text{cur}})} \mathbb{E}_{P(v_{\text{other}}, h_v, a | s, v_{\text{cur}})} \left[\log P(h_v, a | s, v_{\text{cur}}, v_{\text{other}}) - \log P(h_v, a | s, v_{\text{cur}}) \right] \\
 &= \mathbb{E}_{P(s, v_{\text{cur}})} \mathbb{E}_{P(v_{\text{other}}, h_v, a | s, v_{\text{cur}})} \left[\log P(h_v | s, v_{\text{cur}}, v_{\text{other}}) + P(a | h_v, s, v_{\text{cur}}, v_{\text{other}}) \right. \\
 & \quad \left. - \log P(h_v | s, v_{\text{cur}}) - \log P(a | h_v, s, v_{\text{cur}}) \right] \\
 &= \mathbb{E}_{P(s, v_{\text{cur}}, v_{\text{other}})} \left[D_{\text{KL}}(P(h_v | s, v_{\text{cur}}, v_{\text{other}}) || P(h_v | s, v_{\text{cur}})) \right] \\
 & \quad + \mathbb{E}_{P(h_v, s, v_{\text{cur}}, v_{\text{other}})} \left[D_{\text{KL}}(P(a | h_v, s, v_{\text{cur}}, v_{\text{other}}) || P(a | h_v, s, v_{\text{cur}})) \right] \\
 &= \text{MI}(h_v; v_{\text{other}} | s, v_{\text{cur}}) + \text{MI}(a; v_{\text{other}} | s, v_{\text{cur}}, h_v).
 \end{aligned}$$

Similarly, we can also have:

$$\text{MI}(v_{\text{other}}; h_v, a | s, v_{\text{cur}}) = \text{MI}(a; v_{\text{other}} | s, v_{\text{cur}}) + \text{MI}(h_v; v_{\text{other}} | s, v_{\text{cur}}, a).$$

Combining these two equality, we can have:

$$\begin{aligned}
 & \text{MI}(h_v; v_{\text{other}} | s, v_{\text{cur}}) + \text{MI}(a; v_{\text{other}} | s, v_{\text{cur}}, h_v) \\
 &= \text{MI}(a; v_{\text{other}} | s, v_{\text{cur}}) + \text{MI}(h_v; v_{\text{other}} | s, v_{\text{cur}}, a).
 \end{aligned}$$

As a and v_{other} become independent with each other when h_v is given, we have $\text{MI}(a; v_{\text{other}} | s, v_{\text{cur}}, h_v) = 0$. As we also have $\text{MI}(h_v; v_{\text{other}} | s, v_{\text{cur}}, a) \geq 0$, we can have the

following inequality, which basically gives us the conditional version of data processing inequality (Cover, 1999):

$$MI(h_v; v_{\text{other}} | s, v_{\text{cur}}) \geq MI(a; v_{\text{other}} | s, v_{\text{cur}}).$$

Since $MI(a; v_{\text{other}} | s, v_{\text{cur}}) \geq 0$, if we can also have $MI(h_v; v_{\text{other}} | s, v_{\text{cur}}) = 0$, then we can conclude that:

$$MI(a; v_{\text{other}} | s, v_{\text{cur}}) = 0.$$

By expanding this mutual information term, we have:

$$\begin{aligned} & MI(a; v_{\text{other}} | s, v_{\text{cur}}) \\ &= \mathbb{E}_{P(s, v_{\text{cur}}, v_{\text{other}})} \left[D_{\text{KL}}(\pi(a | s, v_{\text{cur}}, v_{\text{other}}) || \pi(a | s, v_{\text{cur}})) \right] \\ &= \mathbb{E}_{P(s, v)} \left[D_{\text{KL}}(\pi(a | s, v) || \pi(a | s, v_{\text{cur}})) \right] \\ &= 0. \end{aligned}$$

Since the KL divergence is non-negative, for the above expectation to be zero, there must be for all state-video pairs $(s, v) \in \mathcal{S} \times \mathcal{V}$ with non-zero probability $P(s, v) > 0$, we have the KL divergence to be zero, $D_{\text{KL}}(\pi(a | s, v) || \pi(a | s, v_{\text{cur}})) = 0$, and conclude our proof. \square

C VISUALIZATION

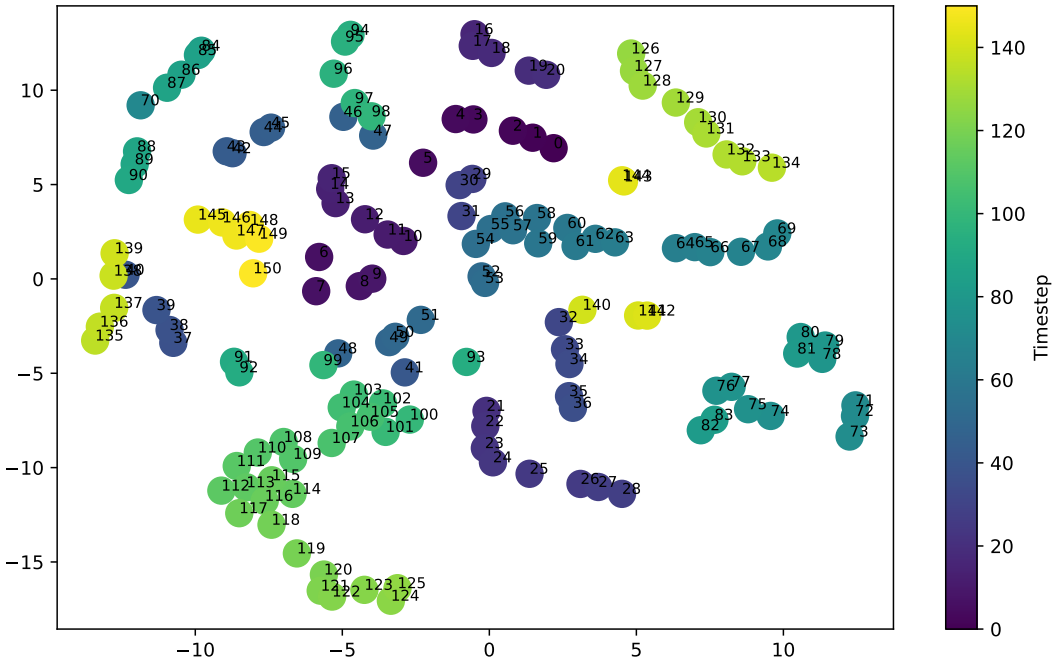


Figure 3: Visualization of h_v over timesteps.

In this section, we present the visualization result of our method. We visualize how h_v of WL-DL changes over timesteps. As shown in Figure 3, we can observe that h_v of WL-DM tends to converge at adjacent timesteps. It is worth noting that since we use a GPT-like transformer architecture as the encoder, the information of video tokens and obs tokens are mixed together in h_v . Furthermore, we do not introduce any task-level information (such as task-level video segmentation annotations), so the clustering results of h_v do not fully correspond to the task.