

# Is Imitation All You Need? Generalized Decision-Making with Dual-Phase Training

Yao Wei<sup>1</sup>, Yanchao Sun<sup>2</sup>, Ruijie Zheng<sup>2</sup>, Sai Vemprala<sup>4</sup>, Rogerio Bonatti<sup>3</sup>, Shuhang Chen<sup>4</sup>  
 Ratnesh Madaan<sup>3</sup>, Zhongjie Ba<sup>1</sup>, Ashish Kapoor<sup>4</sup> and Shuang Ma<sup>3\*</sup>

<sup>1</sup>Jiaxing Research Institute, Zhejiang University <sup>2</sup>University of Maryland, College Park

<sup>3</sup>Microsoft <sup>4</sup>Scaled Foundations

\* Project lead yunyikristy@gmail.com

## Abstract

We introduce *DualMind*, a generalist agent designed to tackle various decision-making tasks that addresses challenges posed by current methods, such as overfitting behaviors and dependence on task-specific fine-tuning. *DualMind* uses a novel “Dual-phase” training strategy that emulates how humans learn to act in the world. The model first learns fundamental common knowledge through a self-supervised objective tailored for control tasks and then learns how to make decisions based on different contexts through imitating behaviors conditioned on given prompts. *DualMind* can handle tasks across domains, scenes, and embodiments using just a single set of model weights and can execute zero-shot prompting without requiring task-specific fine-tuning. We evaluate *DualMind* on *MetaWorld* [55] and *Habitat* [39] through extensive experiments and demonstrate its superior generalizability compared to previous techniques, outperforming other generalist agents by over 50% and 70% on *Habitat* and *MetaWorld*, respectively. On the 45 tasks in *MetaWorld*, *DualMind* achieves over 30 tasks at a 90% success rate $\beta$ . Our source code is available at <https://github.com/yunyikristy/DualMind>.

## 1. Introduction

Transformer-based models, combined with large-scale data, have shown success in generalizing across various tasks in both language and vision. Notable examples include BERT [13], GPT [36], MAE [19], CLIP [35] and Flamingo [1], etc. Recently, there has been a significant focus on developing such general-purpose models for sequential decision-making and control tasks, such as GATO [41]. The prominent approach is to train a decoder-only Transformer with Imitation Learning (IL) on massive datasets from all targeted tasks. By training with prompts, the model can perform zero-shot inference with just task prompts.

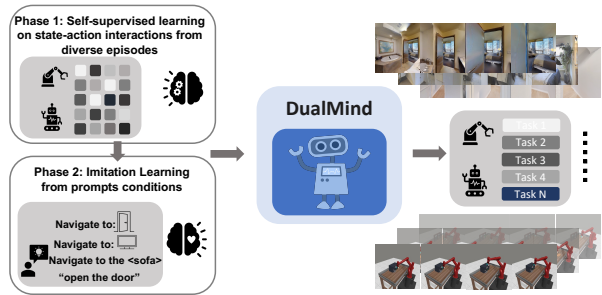


Figure 1: A high-level overview of *DualMind*’s Dual-phase training scheme.

However, such IL-based approaches to general-purpose models face limitations when it comes to sequential control tasks, as highlighted below: (1) *Memorizing behaviors hinders generalization to diverse tasks*: Imitating expert behaviors can lead to memorization and overfitting of specific behaviors that may not be applicable to new situations or variations of tasks, thus limiting the model’s ability to generalize. This limitation is particularly challenging when dealing with a wide range of decision-making tasks that have vastly different configurations, transition functions, and state and action spaces. (2) *Dependence on high-quality data impedes practical application*: IL methods rely heavily on the availability of high-quality expert demonstrations, which can be difficult and expensive to obtain. When the available data is of low quality or not representative of the target task, the performance of the model may suffer.

In light of the aforementioned limitations, self-supervised pretraining has emerged as a viable solution. By focusing on learning common underlying information, a pretrained model can be better equipped to handle diverse tasks. Recently, a study known as SMART [49] has demonstrated the potential of self-supervised pretraining for multi-task decision-making.

Although SMART has shown promising results in promoting generalization, it still requires additional fine-tuning

to adapt to each task. Furthermore, it has only been demonstrated on a small set of tasks on DMC [50]. For decision-making problems that involve numerous tasks with different configurations, finetuning the model for each task can become time-consuming and resource-intensive.

Given the limitations of both IL and self-supervised pre-training discussed earlier, a natural question arises: *How can we develop a decision-making approach that achieves a high degree of generalization without requiring task-specific fine-tuning?* In this paper, we propose DualMind, a generalist agent, to address this question, which stands for our proposed Dual-phase training scheme. The name ‘Dual-Mind’ is derived from our main idea of Dual-phase training for generalized decision-making. Our approach introduces an Encoder-Decoder Control Transformer (Enc-Dec Control Transformer) that models state-action interactions from complex high-dimensional observations. To further improve computational efficiency, DualMind uses Token-Learner [45] as an attention-based Information Bottleneck (IB) [51] to compress the number of tokens so that to speed up training and inference. Building upon Enc-Dec Control Transformer, we propose a Dual-phase training scheme that initially prioritizes policy-independent transition probabilities and encourages the model to capture both short- and long-term temporal granularities. To facilitate zero-shot prompting, we train a second phase on a small fraction of model parameters to learn a generic policy by conditioning on various prompts (such as images, annotations, and language instructions) using a cross-attention mechanism (XAtten.). The Dual-phase training scheme parallels how humans learn to act in the world by first learning underlying common knowledge and subsequently making decisions based on different contexts. Our contributions are summarized below:

1. We propose DualMind, a solution for general-purpose decision-making that can handle various tasks using a single set of weights without task-specific fine-tuning.
2. We introduce a Dual-phase training scheme that overcomes limitations of IL and self-supervised learning.
3. We propose an Encoder-Decoder Transformer (Enc-Dec Control Transformer) that efficiently learns state-action transitions from high-dimensional observation spaces.
4. We conduct extensive experiments on Metaworld [55] and Habitat [39] and show that DualMind outperforms other generalist agents by over 50% and 70% on Habitat and MetaWorld, respectively. We also analyze and ablate different design choices to demonstrate the superior generalizability of DualMind.

## 2. Related work

*Pretraining Visual Representations for Policy Learning:* Recent studies such as R3M [31], APV [47] VPT [4], NRNS [18], PVR [33] and MVP[37] have shown that pre-

trained visual representations can significantly enhance the efficiency of downstream policy learning. However, these works mainly focus on learning object-centric semantics, potentially losing essential control-relevant information. To address this issue, VIP [29] formulates the problem as an offline goal-conditioned RL problem and proposes a visual representation algorithm capable of generating dense reward functions for downstream robotics tasks. On the other hand, COMPASS [28] introduces a general-purpose pretraining pipeline that effectively integrates multimodal signals for autonomous systems.

*Transformer-Based Foundational Model:* The use of high-capacity transformer architectures trained on large-scale datasets has led to significant breakthroughs in various domains. Examples include language models such as BERT [13], GPT-3 [7], T5 [38], and PaLM [11], as well as vision and vision-language models such as MAE [19], Multi-MAE [3], BiT [24], MuST [16], Flamingo [15], and CLIP [35]. For decision-making tasks, recent work such as SMART [49] has proposed a self-supervised pretraining framework tailored for control tasks. For robotics control problems, PACT [5] has shown that a pretrained representation could speed up various downstream tasks of mobile agents, such as navigation and localization.

*A General-Purpose Model for Control:* Since the groundbreaking success of GPT [36], recent research has focused on using Transformer decoder-based models to tackle control tasks in an auto-regressive manner. Decision Transformer (DT) [10, 56] builds on the architecture of GPT to create a generalist agent for sequential decision-making tasks. This has been followed by Multi-game DT [26] and Online-DT [57], which demonstrate the potential of DTs for multi-task and online learning. GATO [41] imitates expert demonstrations from a vast dataset and showcases its ability to handle a large number of tasks. VIMA [21] is an agent that can accept multi-modal prompts for solving various robotics manipulation tasks. In real-life applications, RT-1 [6] has demonstrated the efficacy of this approach in robotic control.

## 3. Preliminary and Overview of DualMind

### 3.1. Problem formulation

We focus on a set of tasks, denoted as  $\mathcal{T}$ , from two representative benchmarks, namely Metaworld [55] and Habitat [39], which cover the *Manipulation* and *Navigation* domains, respectively. As shown in Table 2, our selection of these two benchmarks allows us to conduct a comprehensive study on tasks with a wide variety of characteristics. Here, we define a task as a partially observable Markov decision process (POMDP). The tasks we consider span across several factors, as defined below:

- *Domain:* This refers to tasks with different state/action

	Self-superv.	IL-prompt	Dual-phase (ours)
Learning	Pre.: generic info. FT: task-specific policy	Cond. generic policy	I: generic info. II: cond. generic policy
Data	Pre: Multi-task large set FT: Single-task small set	Multi-task large set+prompts	I: Multi-task large set II: +prompts
Optim. weights	Pre: whole model FT: Entire/freeze+Task heads	Entire model	I: Entire model II: Partial/freeze+XAtten.
Inference task	Single	Multiple	Multiple
No need FT	✗	✓	✓
Zero-shot promp.	✗	✓	✓
Final utilization	Many models for each task	Single model	Single model

Table 1: Comparisons of different training approaches.

Bench.	Dom.	Sc.	Emb.	Prom.	Tasks	Epis.
Meta.	Man.	1	1	inst.	50	50K
Habit.	Nav.	933	1	Obj. / Img	27	50K
Total	2	934	2	3	77	100K

Table 2: Dataset summerization Dom.: domains, Sc.: number scenes, Emb.: number of embodiments, Prom.: types of prompts, Epis.: number of episodes.

spaces and application scenarios. In our study, *Manipulation* and *Navigation* are two domains we focus on.

- *Embodiment*: This factor is used to differentiate tasks that have different physics and action spaces. For instance, a robot arm and an embodied agent in MetaWorld and Habitat are considered as different embodiments. Differences can also exist in the same domain, such as arms with distinct joint torques and/or hardware configurations.
- *Scene*: This refers to tasks that are performed in different observation spaces, state spaces, and world structures. For example, in Habitat, agents that navigate in different rooms should adapt to various visual appearances and geometry structures.
- *Prompt*: This factor captures different forms of prompt conditions. In MetaWorld, prompts are natural language instructions, while in Habitat, we use a single RGB image or an object annotation as the navigation goal to prompt our model.

### 3.2. Overview of Dual-phase training scheme

In this section, we provide a brief overview of DualMind and compare it with two other prominent approaches: self-supervised pretraining (Self-superv.) and Imitation Learning with prompt conditions (IL-prompt). We also provide insights into the central idea behind our proposed approach. A summarized comparison of these approaches is shown in Table 1.

As shown in Fig. 2, In Phase I, we train the entire Enc-Dec Control Transformer (Sec. 4.1) with a self-supervised

training objective to capture generic information of state-action transitions. In Phase II, we train only a small part of Enc-Dec Control Transformer attached with XAtten. on a diverse set of prompts for a conditional generic policy. After the Dual-phase training, we can obtain one model with a single set of weights that can be directly applied to a large number of tasks with corresponding prompts.

Compared to other generalist agents like GATO [41], which trains an imitating policy directly, DualMind demonstrates superior generalization capability. Moreover, our Phase II requires training only a small fraction of model weights while freezing the remaining parts, resulting in faster learning and reduced training cost when optimizing the model with the same number of iterations. Additionally, compared to self-supervised learning approaches such as SMART [49], DualMind is simple and effective, making it suitable for a wide range of application scenarios.

### 3.3. Insights

The central idea behind our Dual-phase is to mimic how humans learn to act in the world, first by learning underlying common knowledge and then by learning to make decisions based on different contexts. Our approach relates to InstructGPT [32], which aims to align language models with user intent by fine-tuning them with human feedback. In analogy to InstructGPT, our Phase I can be considered as learning a general model that captures the common essential information. However, as stated in InstructGPT, this is different from the objective of “following task instructions (i.e. prompt conditions),” and thus such a model is *misaligned*. Therefore, in the second phase, we leverage conditional IL to align the model so that it can perform well for any given prompts.

## 4. Approach

In this section, we introduce our proposed DualMind. We present the model architecture in Section 4.1, and illustrate the training objective for DualMind in Section 4.2.

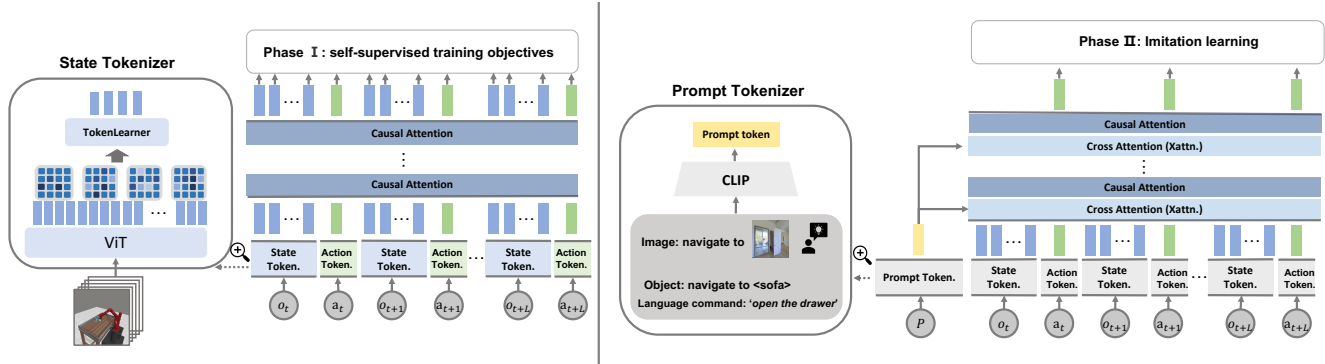


Figure 2: The architecture diagram of DualMind. **Left: Phase I.** Agent is trained with self-supervised learning objectives. During this phase, Transformer encoder and decoder jointly trained. **Right: Phase II.** Agent is trained with prompt conditional imitation learning. We tokenize task prompts with a pretrained CLIP encoder, and condition the Transformer decoder on the prompt through XAttn. layers. The gray color indicates frozen modules. (Detailed training objectives are in Sec. 4.2.)

### 4.1. Model Architecture

We propose an Encoder-Decoder Control Transformer to process state-action interaction sequences, as illustrated in Figure 2. The implementation details of each component in the Enc-Dec Control Transformer are outlined below.

**State tokenizer.** We utilize a ViT model [14] to tokenize raw pixel states. To reduce the computational burden of dealing with sequential decision-making tasks, we leverage an attention-based Information Bottleneck (IB) to further compress the number of tokens so as to speed up training and inference (Fig. 2-left). Specifically, we use TokenLearner [45] which is an element-wise attention module that learns to soft-select image tokens, passing only the important ones to subsequent layers. The inclusion of TokenLearner sub-samples the 196 state tokens that come out of ViT to just 8 tokens that are then passed to the Transformer decoder layers.

**Action tokenizer.** To handle both continuous and discrete action spaces in our two domains, we adapt a strategy similar to GATO [41] by discretizing continuous actions into bins. We first flatten the actions into sequences of floating point values in row-major order, and then mu-law encode them to the range [-1, 1] before discretizing them into 256 uniform bins. Discrete actions are tokenized into 256 bins in the same way.

**Transformer decoder.** Our transformer decoder architecture is similar to Control Transformer [49], but with a modification. In our approach, we encode each state into 8 tokens, which is different from SMART’s single token representation. This modification enables richer representation learning, making it suitable for more complex visual control environments.

**Prompt tokenizer.** We tokenize prompts using a pretrained CLIP encoder [35]. For “image goal” prompts in Habitat, we use the CLIP image encoder, while for “ob-

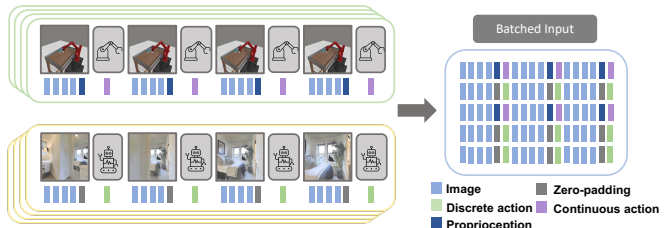


Figure 3: Batch input when training on multiple domains. For “image goal” prompts in Habitat and “language instruction” prompts in MetaWorld, we use the CLIP text encoder. A learnable linear layer is added on top of the CLIP encoders to map all prompts to prompt tokens with the same dimensions. During training in both phases, we freeze the CLIP encoders.

**XAttn. layer.** We condition the Transformer decoder by training it to learn from the prompt sequence through a series of cross-attention layers. The output sequence from each cross-attention layer is computed by  $\text{softmax}(\frac{q_H k_P^T}{\sqrt{d}})v_P$ , where  $H$  is the sequence of episodes,  $P$  is prompt, and  $d$  is the embedding dimension. This design builds a stronger connection between the prompts and the demonstrations, which is an improvement over prefix-style prompting approaches [41]. We will show the benefits of this design in Sec. 5.4.

### 4.2. Training objectives

**Phase I: self-supervised SMART training.** The goal of this phase is to learn a good representation that captures control-relevant information shared across tasks. In this phase, we jointly train the encoder and the decoder following the self-supervised training objectives of SMART [49]. We use  $F_\theta$  to denote the learned model with parametrization  $\theta$ , such that  $F_\theta(o_{i:j}, a_{i:j})$  refers to the output tokens of the decoder corresponding to raw inputs  $o_{i:j}$  and  $a_{i:j}$ , the observation and action sequence from step  $i$  to step  $j$ . For a sequence of observations and actions denoted as

$\{o_t, a_t, \dots, o_{t+L}, a_{t+L}\}$  with context length  $L$ , we minimize the following objective.

$$\mathcal{L}_{P1} := \mathcal{L}_1 + \mathcal{L}_2 + \mathcal{L}_3, \text{ where} \quad (1)$$

$$\mathcal{L}_1 := \sum_{i=0}^{L-1} l(f_1(F_\theta(o_{t:t+i}, a_{t:t+i})), \bar{\phi}(o_{t+i+1})), \quad (2)$$

$$\mathcal{L}_2 := \sum_{i=1}^L l(f_2(F_\theta(o_{t:t+i}, a_{t:t+i-1}), a_{t+i}), \quad (3)$$

$$\mathcal{L}_3 := \sum_{i=1}^{L-1} l(f_3(F_\theta(\text{Mask}(o_{t:t+L}, a_{t:t+L})), a_{t+i}). \quad (4)$$

Here  $l$  is a loss function that is selected by the variable type. For latent states, we use a mean squared error, while for discrete actions, we use the cross-entropy loss.  $\mathcal{L}_1$  is to learn a forward prediction head  $f_1$  that can predict the next state representation based on the historical interactions. Since the groundtruth state representation is unknown, we use the learned state embedding from the ViT model to encode the next observation, denoted as  $\bar{\phi}$  where the overline stands for gradient stopping.  $\mathcal{L}_2$  aims to recover the action token in each step conditioning on the history and the next state.  $\mathcal{L}_3$  masks a proportion of input tokens and learns to recover the masked actions, which can extract long-term temporal dependence for control.

### Phase II: Imitation learning with prompt conditions.

In this phase, we train the model to follow prompt conditions. We formulate various tasks as a conditional generation problems, where the conditions can be goals, commands, prompts, etc. During Phase II, we let the agent learn a conditional policy, using expert trajectories with associated prompts. Let  $\psi$  be the prompt tokenizer, and  $\pi$  be the learned policy whose inputs are the representation tokens given by the decoder. For an expert sequence  $\{o_t, a_t, \dots, o_{t+L}, a_{t+L}\}$  with prompt  $P$ , we minimize loss

$$\mathcal{L}_{P2} := \sum_{i=0}^{L-1} l(\pi(F_\theta(o_{t:t+i}, a_{t:t+i}; \psi(P))), a_{t+i+1}). \quad (5)$$

Note that in this phase, we do not train the entire model  $F_\theta$ , and instead only re-train a small fraction of it. More discussion is in Sec. 5.4.

## 5. Experiments

### 5.1. Experimental setup

**Data.** We evaluate and train DualMind on two benchmarks, Habitat [39] and MetaWorld [55]. Habitat is a photorealistic simulation platform for research in Embodied AI, emphasizing active perception and long-term planning, while MetaWorld is a simulated benchmark for multi-task learning and meta-reinforcement learning, comprising 50 distinct robotic manipulation environments. Training on datasets collected from both these benchmarks allows us to demonstrate the model’s generalizability across domains, embodiments, scenes, and prompts. We provide a detailed introduction to these factors in Sec. 3.1 and summarize them in Table 2. Additionally, we use 10 tasks as an out-of-distribution testbed to showcase the model’s generalization

capability. More details about our data collection process can be found in Appendix A.

**Comparing baselines.** We compare DualMind with existing transformer-based approaches and present results from two versions of our model: a generalist agent trained on the full dataset (DualMind) and a single-domain specialist trained only on data from either MetaWorld or Habitat (DualMind/single). To ensure fair comparisons, we implemented related works ourselves and trained and evaluated them on the same data and model architecture. We provide information on each baseline below:

- **IL-only** is a model trained only with prompt-conditioned imitation learning, which is related to GATO but uses a different prompting conditioning method.
- **SMART-only** is a model trained only using SMART training objectives (purely self-supervised).
- **Jointly** is a model jointly trained with both SMART objectives and prompt-conditioned IL loss.
- **GATO\*** is the model described in the original paper. We include its reported performance on the Metaworld benchmark for reference. Notably, this model has 1.18 billion parameters and was trained on massive datasets, including 94.6k episodes from Metaworld. In comparison, DualMind has 175 million parameters and was trained on a smaller dataset consisting of 100k episodes, of which 50k are from MetaWorld.
- **GATO** is a model we implemented ourselves, reproducing the main technical approaches presented in the original paper. For a fair comparison, we used the same base model architecture (Enc-Dec Control Transformer), but replaced our XAtten.-based prompting approach with their proposed prefix prompting approach.

Specifically, we train **IL-only**, **SMART-only**, and **Joint-only** on Enc-Dec Control Transformer +XAtten., which has the same model architecture as DualMind. For **GATO**, we train it using Enc-Dec Control Transformer but insert prompts in a prefix manner since it uses a different prompting method. Moreover, we provide the performance of **GATO** reported in their paper for reference, denoted as **GATO\***. More details about our baselines can be found in Appendix A.

**Implementation details.** Our implementation of DualMind uses a Transformer-based architecture consisting of a ViT-B [17] model, a TokenLearner [45], and a GPT model [36] as the encoder and decoder, respectively. The decoder consists of 8 layers and 8 attention heads, with a context length of  $L=6$  and an embedding size of  $d=512$ <sup>1</sup>. We trained our model with the AdamW optimizer and a learning rate of  $5e-5$  for both training phases. In Phase i, we trained the model for around 40 hours with  $BS=16$  on

<sup>1</sup>We found that longer context lengths can produce better performance, particularly on tasks that rely on long-range temporal dependencies. See Appendix B for more ablations.

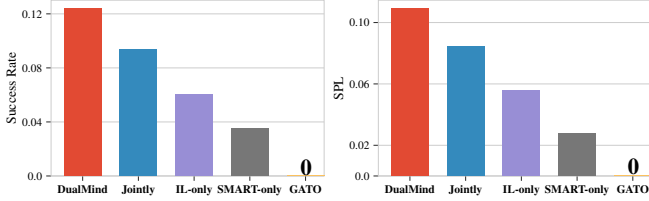


Figure 4: Comparisons of **generalist agents** on *Habitat 4* scenes with 3 difficulty levels per scene. We roll out the agents 3 times on each scene and average the defined scores, and compare agents by Success Rate (SR) (left) and Success weighted by Path Length (SPL) (right).

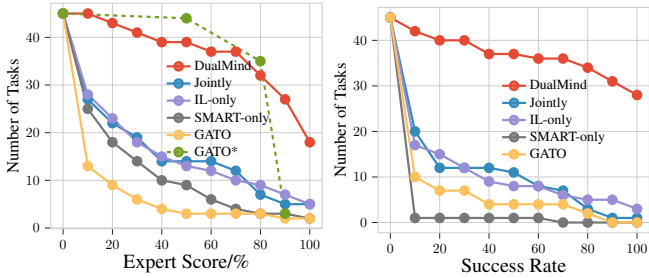


Figure 5: Comparisons of **generalist agents** on *MetaWorld 45* tasks on Percentage of Expert Score (PES) (left) and Success Rate (SR) (right).

5x8xV100 GPUs. In Phase ii, the model was trained for about 12 hours with BS=128 on 2x8xV100 GPUs. Further implementation details are provided in Appendix A.

## 5.2. Capabilities of DualMind

In this section, we aim to demonstrate the capabilities of DualMind on all tasks. Note that, as a generalist agent, the performance on both MetaWorld and Habitat are achieved by a single model. The performance is shown in Fig. 4 and Fig. 5. To provide a reference for readers, we follow GATO’s evaluation protocol and report the Percentage Expert Score (PES), which measures the number of distinct tasks for which each model performs above a given score threshold relative to the expert performance. For each task, we roll out the model 10 times and average the defined scores. As shown in Fig. 5, DualMind achieves over 90% expert score threshold across more than 27 tasks, outperforming GATO\* by a large margin, which only has three tasks above the threshold. On lower expert score thresholds, for example, 80% and 50%, DualMind can also achieve comparable performance. However, it should be noted that GATO’s performance was achieved by their 1.18B model trained on massive datasets. Therefore, this is just a reference for readers, and a fully fair comparison with GATO cannot be performed without access to both the model and data. To provide a more fair comparison, we compare DualMind with a self-implemented GATO, which will be discussed in more detail in Sec. 5.1. We also report the num-

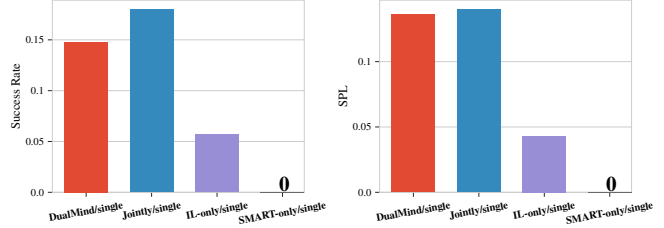


Figure 6: Comparisons of **single-domain specialist** on *Habitat 12* scenes by SR (left) and SPL (right).

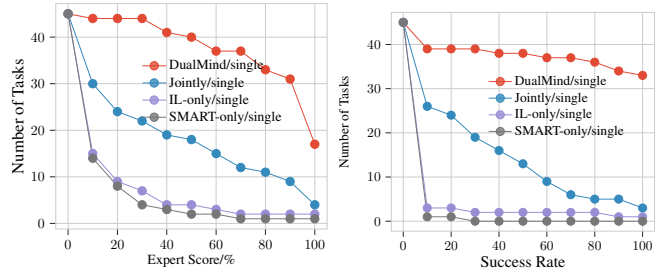


Figure 7: Comparisons of **single-domain specialist** on *MetaWorld 45* tasks by PES (left) and SR (right).

ber of tasks for which our model performs above a given Success Rate (SR). DualMind achieves 39 tasks at over 0.5 SR and can maintain good performance on higher SRs, with 34 tasks at over 0.8 SR and 28 tasks at over 1 SR. We present the performance of DualMind on Habitat by averaging across all 12 testing scenes and reporting the success rate (SR) and success weighted by path length (SPL) evaluation metrics. As shown Fig. 4, DualMind outperforms the other baseline models by a large margin under both evaluation metrics. (See performance on each task in Appendix B.)

## 5.3. Analysis

### 5.3.1 Different training regimes

#### Is imitation learning all you need for a generalist agent?

To answer the question, in this experiment, we compare DualMind with its counterpart trained only with Imitation Learning objective, i.e., IL-only. In Fig. 4 and Fig. 5, we present the comparison results between the generalist multi-domain agents. As shown in the figures, DualMind outperforms its IL-only counterpart by over 50% and 70% on Habitat and MetaWorld, respectively. Specifically, DualMind performs well on 39 out of 45 tasks over the 50% expert score threshold, while IL-only only performs well on 13 tasks. As the difficulty of the tasks increases, DualMind still maintains good performance, achieving 18 tasks and 28 tasks at the 100% expert score and SR, respectively, while IL-only only achieves 5 tasks. Similar observations can also be made when comparing the single-domain specialist agents (Fig. 7 and Fig. 6). (See performance on each tasks in Appendix B.)

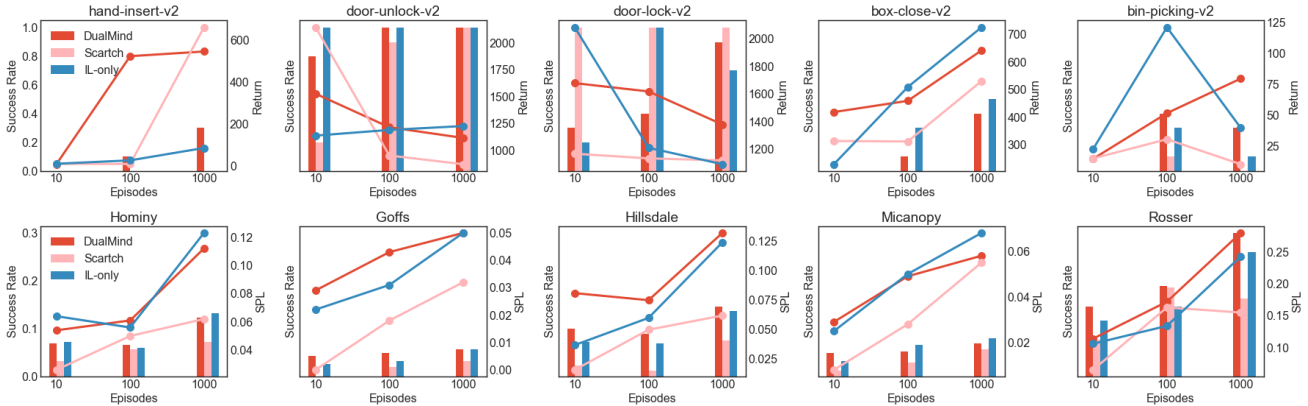


Figure 8: Few-shot comparisons of generalist agents on out-of-distribution tasks. The performance of the success rate (left axis, bar-chart) or Return/SPL (right axis, line-chart) on different tasks after we performed 10-, 100-, and 1000-shot learning on DualMind (red), IL-only (pink), and Scratch (blue).

We can infer from this that Imitation Learning alone may not suffice to build a truly general-purpose model, particularly when aiming to tackle tasks that span a broad range of domains. Even within a single domain, variations in embodiments, scenes, and instructions can pose significant challenges. We conducted additional investigations into the generalization capabilities by comparing different approaches on out-of-distribution tasks, as demonstrated in Section 5.3.2.

### Can self-supervised learning well-align with instructions without FT?

To address this inquiry, we compare DualMind with its self-supervised equivalent, SMART-only, while also evaluating both single- and multi-domain agents. As depicted in Fig. 4 and Fig. 5, DualMind exhibits superior performance compared to SMART-only, with over 75% and 78% better results on Habitat and MetaWorld, respectively. Notably, SMART-only is unable to succeed on any tasks when applied to single-domain agents, whereas DualMind maintains a significant advantage, particularly on MetaWorld.

Our hypothesis is that SMART, being a pretrain-finetune pipeline, is unlikely to attain the desired performance without post-finetuning. Even when training SMART-only by providing prompts in the same manner as DualMind, zero-shot prompting may not be achievable due to limitations in the self-supervised training objective not being well-aligned with task instructions, as detailed in Section 3. Additionally, we noted that SMART-only surpasses its single-domain equivalent, suggesting its effectiveness in capturing shared knowledge across diverse data.

### Do we need to train them in two phases?

As DualMind is trained using different objectives in two phases, one may question the necessity of such an approach. Firstly, from an optimization standpoint, training all four losses jointly may present more challenges in terms of steady optimization. Different optimization directions could potentially conflict with each other, and vary-

ing convergence rates could hinder all objectives from being trained to reach optimality. Furthermore, in terms of computational costs, DualMind only needs to optimize a small portion of the model weights in phase 2 (as demonstrated in the ablations presented in Section 5.4). This makes the training process more efficient and cost-effective compared to its jointly trained counterpart. In this experiment, we provide further empirical evidence to support this claim.

As illustrated in Fig. 4, Fig. 5, Fig. 6, and Fig. 7, Jointly outperforms IL-only and SMART-only, thereby confirming the necessity of utilizing all training objectives. However, it lags behind DualMind by a considerable margin in both multi- and single-domain comparisons. Interestingly, Jointly slightly outperforms DualMind in single-domain comparisons. We hypothesize that the optimization challenges may not be as significant as those encountered when training on data from the same domain.

### 5.3.2 Out-of-distribution tasks

The objective of this experiment is to assess the ability of our model to solve novel tasks. To achieve this, we evaluate our models on 10 held-out tasks from two domains, namely MetaWorld and Habitat. The MetaWorld tasks consist of "hand-insert-v2", "door-unlock-v2", "door-lock-v2", "box-close-v2", and "bin-picking-v2", whereas the Habitat tasks include "Goffs", "Hominy", "Hillsdale", "Micanopy", and "Rosser". To evaluate the performance of our models, we follow the evaluation protocol with GATO, which involves finetuning each agent on a limited number of demonstrations. Specifically, we conduct 10-, 100-, and 1000-shot learning. Further details on the evaluation protocol can be found in Appendix A.

We compare the performance of three models, namely DualMind, IL-only, and Scratch. Scratch refers to the model that is trained on few-shot demonstrations from randomly initialized model weights. As demonstrated in Fig. 8, Scratch performs the worst among the three mod-

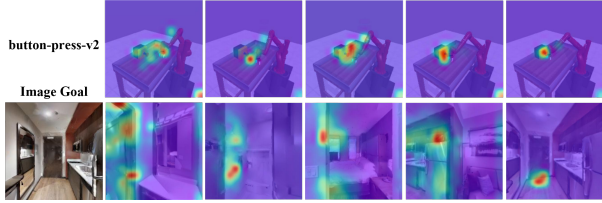


Figure 9: Attention map visualization.

els in most cases.

Upon comparing DualMind with `IL-only`, we observe that DualMind exhibits superior performance across various shot settings. Specifically, in terms of the SR metric, DualMind outperforms `IL-only` on 8 out of 10 tasks at 10-shot and on 7 tasks at 100- and 1000-shot demonstrations. Furthermore, with respect to the SPL and PES metrics, DualMind achieves better results than `IL-only` on 9 tasks in the 10-shot experiment. These results provide further evidence that the proposed Dual-phase training approach can enhance the generalization ability of models even when dealing with novel tasks and limited demonstrations.

### 5.3.3 Attention visualization

To gain insight into how DualMind is able to perform diverse tasks, we conduct attention visualization. We present attention maps for tasks from both Habitat and MetaWorld, where we display a sequence of frames from the episode for each task.

The attention maps reveal that when performing manipulation tasks in MetaWorld, such as and “button-press-v2”, the model initially focuses on the execution context and then shifts its attention to the targeting instance, such as the “button”, until the task is completed. Notably, for navigation tasks in Habitat, DualMind learns to explore the scene to locate the goal. For example, as shown in Fig. 9, given an image goal, the agent first attends to the entrance to navigate into the restroom. Upon realizing that the goal is not there, it steps out and searches for another room to enter. After spotting the refrigerator, which appears in the image goal, the agent quickly locks onto the goal and completes the task. These attention maps provide insight into how DualMind leverages its generalization ability to solve new tasks.

## 5.4. Ablation study

### Training parts in Phase II

In this section, we ablate DualMind by varying model weights that been trained in Phase II, as listed below:

- ①: freeze the entire Enc-Dec Control Transformer arch by only train the cross-attention layers.
- ②: freeze the Transformer Encoder (State tokenizer) and the first 4 layers of Transformer Decoder.
- ③: freeze the Transformer Encoder.

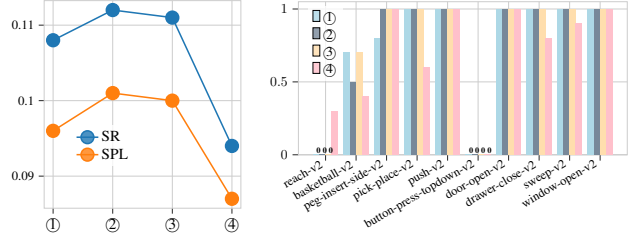


Figure 10: Comparisons of frozen parts in Phase II.

- ④: no frozen part, optimize the entire model in Phase II. As shown in Fig. 10, ② and ③ perform the best in most cases. For our experiments, we use ③. However, for future scaled-up models and data, we would recommend using ② since it saves more computational cost. When training each setting with the same number of iterations, ④ performs poorly, which may be due to slow convergence with more model weights. This result also suggests that after training in phase I, our model has learned useful information, but insufficient re-training in phase II may lead to performance deterioration due to potential forgetting issues.

### Prompt conditioning

We conducted an ablation study on DualMind by comparing two prompt conditioning approaches: prefix and XAtten. prompting. We used the average success rate of ML10 training tasks as the comparison metric for Metaworld. Results show that XAtten. prompting achieves a 0.76 SR on Metaworld and an 0.11 SR on Habitat, while prefix prompting only achieves 0.29 SR and 0 SR, respectively. The cross-attention mechanism in XAtten. prompting allows the agent to establish a strong connection between prompts and demonstrations, which is particularly useful for goal-conditioned tasks. (See more details and discussion in Appendix B.)

## 6. Conclusion

This paper presents a new training approach for generalist agents called DualMind, which consists of two phases: self-supervised learning of basic and generic knowledge across various tasks, followed by imitation of expert behaviors with different types of prompt conditioning. By utilizing a carefully designed Transformer Encoder-Decoder architecture and a dual-phase training scheme, DualMind is scalable, versatile, and generalizable. Empirical evaluation on two challenging domains, Habitat and MetaWorld, shows that DualMind outperforms previous generalist learning methods and pretraining approaches. Further analysis and ablations demonstrate the effectiveness of the dual-phase design.

Future work includes expanding DualMind to more domains and tasks, finding efficient solutions for handling longer context lengths in demonstrations, and enabling practical training in online interactive scenarios.



## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*, 2022. [1](#)
- [2] Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, et al. On evaluation of embodied navigation agents. *arXiv preprint arXiv:1807.06757*, 2018. [16](#)
- [3] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. Multima: Multi-modal multi-task masked autoencoders. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII*, pages 348–367. Springer, 2022. [2](#), [13](#)
- [4] 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. *arXiv preprint arXiv:2206.11795*, 2022. [2](#)
- [5] Rogerio Bonatti, Sai Vemprala, Shuang Ma, Felipe Frujeri, Shuhang Chen, and Ashish Kapoor. Pact: Perception-action causal transformer for autoregressive robotics pre-training. *arXiv preprint arXiv:2209.11133*, 2022. [2](#)
- [6] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. [2](#)
- [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020. [2](#)
- [8] Arthur Buckner, Luis Figueredo, Sami Haddadin, Ashish Kapoor, Shuang Ma, Sai Vemprala, and Rogerio Bonatti. Latte: Language trajectory transformer, 2022. [13](#)
- [9] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*, 2017. [14](#)
- [10] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021. [2](#)
- [11] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. [2](#)
- [12] Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Jordi Salvador, Kiana Ehsani, Winson Han, Eric Kolve, Ali Farhadi, Aniruddha Kembhavi, and Roozbeh Mottaghi. Proctor: Large-scale embodied ai using procedural generation, 2022. [13](#)
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics, 2019. [1](#), [2](#)
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [4](#), [12](#)
- [15] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. [2](#)
- [16] Golnaz Ghiasi, Barret Zoph, Ekin D. Cubuk, Quoc V. Le, and Tsung-Yi Lin. Multi-task self-training for learning general representations. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8836–8845, 2021. [2](#)
- [17] Agrim Gupta, Stephen Tian, Yunzhi Zhang, Jiajun Wu, Roberto Martín-Martín, and Li Fei-Fei. Maskvit: Masked visual pre-training for video prediction. *arXiv preprint arXiv:2206.11894*, 2022. [5](#)
- [18] Meera Hahn, Devendra Singh Chaplot, Shubham Tulsiani, Mustafa Mukadam, James M Rehg, and Abhinav Gupta. No rl, no simulation: Learning to navigate without navigating. *Advances in Neural Information Processing Systems*, 34:26661–26673, 2021. [2](#)
- [19] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. [1](#), [2](#)
- [20] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Fredrik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning, 2022. [13](#)
- [21] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: General robot

- manipulation with multimodal prompts. *arXiv preprint arXiv:2210.03094*, 2022. [2](#), [17](#)
- [22] Michael Kelly, Chelsea Sidrane, Katherine Driggs-Campbell, and Mykel J. Kochenderfer. Hg-dagger: Interactive imitation learning with human experts, 2019. [13](#)
- [23] Apoorv Khandelwal, Luca Weihs, Roozbeh Mottaghi, and Aniruddha Kembhavi. Simple but effective: Clip embeddings for embodied ai. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14829–14838, 2022. [13](#)
- [24] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V*, page 491–507, Berlin, Heidelberg, 2020. Springer-Verlag. [2](#)
- [25] Michael Laskey, Jonathan Lee, Roy Fox, Anca Dragan, and Ken Goldberg. Dart: Noise injection for robust imitation learning, 2017. [13](#)
- [26] Kuang-Huei Lee, Ofir Nachum, Sherry Yang, Lisa Lee, C. Daniel Freeman, Sergio Guadarrama, Ian Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, and Igor Mordatch. Multi-game decision transformers. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. [2](#)
- [27] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*. [13](#)
- [28] Shuang Ma, Sai Vemprala, Wenshan Wang, Jayesh K Gupta, Yale Song, Daniel McDufft, and Ashish Kapoor. Compass: Contrastive multimodal pretraining for autonomous systems. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1000–1007. IEEE, 2022. [2](#)
- [29] Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy Zhang. VIP: Towards universal visual reward and representation via value-implicit pre-training. In *The Eleventh International Conference on Learning Representations*, 2023. [2](#)
- [30] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 220–229, 2019. [13](#)
- [31] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022. [2](#)
- [32] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022. [3](#)
- [33] Simone Parisi, Aravind Rajeswaran, Senthil Purushwalkam, and Abhinav Gupta. The unsurprising effectiveness of pre-trained vision models for control. In *International Conference on Machine Learning*, pages 17359–17371. PMLR, 2022. [2](#)
- [34] Peter Pastor, Heiko Hoffmann, Tamim Asfour, and Stefan Schaal. Learning and generalization of motor skills by learning from demonstration. In *2009 IEEE International Conference on Robotics and Automation*, pages 763–768, 2009. [13](#)
- [35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. [1](#), [2](#), [4](#), [13](#), [15](#)
- [36] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018. [1](#), [2](#), [5](#)
- [37] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell. Real-world robot learning with masked visual pre-training. *arXiv preprint arXiv:2210.03109*, 2022. [2](#)
- [38] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. [2](#)
- [39] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexander Clegg, John M Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. [1](#), [2](#), [5](#), [14](#)
- [40] Ram Ramrakhya, Eric Undersander, Dhruv Batra, and Abhishek Das. Habitat-web: Learning embodied object-search strategies from human demonstrations at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5173–5183, 2022. [13](#)
- [41] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022. [1](#), [2](#), [3](#), [4](#), [12](#), [13](#), [14](#), [16](#)
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. [17](#)
- [43] Stephane Ross and J. Andrew Bagnell. Reinforcement and imitation learning via interactive no-regret learning, 2014. [13](#)
- [44] Stéphane Ross, Geoffrey J Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *AISTATS*, volume 1, page 6, 2011. [13](#)
- [45] Michael Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova. Tokenlearner: Adaptive space-time tokenization for videos. In M. Ranzato,

- A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 12786–12797. Curran Associates, Inc., 2021. [2](#), [4](#), [5](#), [12](#)
- [46] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9339–9347, 2019. [13](#)
- [47] Younggyo Seo, Kimin Lee, Stephen L James, and Pieter Abbeel. Reinforcement learning with action-free pre-training from videos. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 19561–19579. PMLR, 17–23 Jul 2022. [2](#)
- [48] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10740–10749, 2020. [13](#)
- [49] Yanchao Sun, Shuang Ma, Ratnesh Madaan, Rogerio Bonatti, Furong Huang, and Ashish Kapoor. Smart: Self-supervised multi-task pretraining with control transformers. *arXiv preprint arXiv:2301.09816*, 2023. [1](#), [2](#), [3](#), [4](#), [12](#)
- [50] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018. [2](#)
- [51] Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000. [2](#)
- [52] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio, 2016. [14](#)
- [53] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models, 2023. [13](#)
- [54] Fei Xia, Amir R Zamir, Zhiyang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: Real-world perception for embodied agents. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9068–9079, 2018. [14](#)
- [55] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pages 1094–1100. PMLR, 2020. [1](#), [2](#), [5](#), [13](#), [15](#), [16](#)
- [56] Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 27042–27059. PMLR, 17–23 Jul 2022. [2](#)
- [57] Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. In *International Conference on Machine Learning*, pages 27042–27059. PMLR, 2022. [2](#)

# Appendix

## Model Card of DualMind

<b>Model Details</b>	
Model Type	Encoder-Decoder Transformer (Enc-Dec Control Transformer), built upon a ViT encoder [14], a TokenLearner [45], and a Control Transformer [49].
Training Process	The training process is divided into two phases. In Phase I, the entire Enc-Dec Control Transformer is trained with a self-supervised training objective. In Phase II, a small part of the model is trained with imitation learning conditioned on prompts. Detailed training objective is in Sec. 4.2.
Model Version	Initial release.
<b>Intended Uses</b>	
Primary Intended Uses	The proposed model aims to perform a wide range of control tasks spanning multiple domains, visual scenes and embodiments. Our intention is to create a general-purpose decision-making solution capable of handling various tasks using a single set of weights, without requiring task-specific fine-tuning.
<b>Factors</b>	
Relevant Factors	Multiple factors can influence the performance of the model. First, the quality of training dataset has influence on the results, including task diversity, behavior policy performance, data volume, etc. Second, model implementation hyperparameter setting, and training objectives will also alter the final performance.
Evaluation Factors	We report the performance of the model in multiple sets of tasks, and conducted ablation study in Sec. 5.3.2.
<b>Metrics</b>	
Model Performance Measures	Our downstream task performance is measured using success rate, SPL, and expert score, as detailed in Sec. A.6. The expert score is calculated in the same manner as GATO [41], while using a different dataset.
Decision thresholds	N/A
Approaches to Uncertainty and Variability	The model evaluation process inevitably involves uncertainties. In order to reduce the variance introduced during the evaluation, we employed 3 random seeds for the Habitat evaluation and 10 random seeds for the Metaworld evaluation.
<b>Evaluation Data</b>	
Datasets	Our DualMind is evaluated on multiple control tasks from Habitat and Metaworld. Both in-distribution and out-of-distribution tasks are considered. <b>Habitat:</b> Our experiments are focused on the ImageNav task, we chose 4 scenes. We hold out 5 Gibson scenes for the experiments of out-of-distribution tasks, as detailed in Sec. A.6. <b>Meta-world:</b> We select 45 training tasks on ML45 for the experiments of evaluation, and hold out 5 test tasks for the experiments of out-of-distribution, as detailed in Sec. A.6.
Motivation	Our evaluation of DualMind consists of two components. First, we evaluated its performance on in-distribution tasks to understand how well it handles tasks across domains, scenes, and embodiments using a single set of model weights. Second, we evaluated DualMind on out-of-distribution tasks to assess its ability to adapt to entirely new tasks.
Preprocessing	Observations are tokenized into the same embedding sequence before being input to transformer decoder, as detailed in Sec. 4.1.
<b>Training Data</b>	
Datasets	The model is trained using 100K episodes collected from Habitat and Metaworld, with 50k episodes (~3.26M interaction steps) on Habitat and 50K episodes (~3.82M interaction steps) on Metaworld, respectively.

Motivation	In order to ensure that DualMind can handle tasks across domains, scenes, and embodiments, we collected data for all tasks in Metaworld and all scenes in Habitat. The data collection process is detailed in Sec. A.4.
Preprocessing	The multi-domain data is tokenized into the same embedding sequence before being fed to the transformer decoder, as detailed in Sec. 4.1.
<b>Quantitative Analyses</b>	
Unitary Results	We evaluated the performance of DualMind on the Metaworld and Habitat benchmarks. In Sections 5.2 and 5.3.1, we demonstrate the general capabilities of DualMind across both Metaworld and Habitat tasks. Additionally, in Section 5.3.2, we analyze its performance on out-of-distribution tasks.
<b>Ethical Considerations</b>	
Data	Our data is collected from simulators of navigation and manipulation, and thus it does not include any unethical data.
Risks and Harms	Our current training and evaluation are conducted in simulators, and do not involve physical robots where model malfunctioning could lead to safety issues.
Mitigations	N/A
<b>Caveats and Recommendation</b>	
Future work	Our future work includes expanding DualMind to more domains and tasks, finding efficient solutions for handling longer context lengths in demonstrations, and enabling practical training in online interactive scenarios.

Table 3: Model card of DualMind, following the framework proposed by [30].

## A. Implementation details

### A.1. Additional related work

*Embodied AI*: Embodied AI is an interdisciplinary research field that intersects with robotics, computer vision, machine learning, artificial intelligence, and simulation [46, 40, 48, 55]. It distinguishes itself from other areas through its emphasis on enabling AI systems to understand and interact with the physical world. Embodied AI often employs simulation to abstract away lower-level robotic control, emphasizing high-level task planning. Robotic manipulation and navigation tasks are addressed by LaTTe [8] and Embodied-CLIP [23] using the frozen visual and textual representations of CLIP [35]. Proctor [12] illustrates how leveraging massive datasets for training can significantly enhance the performance of embodied AI tasks. Voyager [53] proposes an autonomous Minecraft agent, using GPT-4 to learn and excel in gameplay, solving novel tasks in new worlds.

*Imitation Learning (IL)*: Imitation learning operates by reducing the discrepancy between the expert’s action and the action that the policy predicts, using the expert’s trajectory for guidance without reward for learning. Imitation Learning (IL) [34] has typically been employed for tasks with low-dimensional state spaces. DAgger [44] and its variants [43, 25, 22] employ active learning strategies to mitigate the distribution shift encountered by the learner, thereby enhancing its ability to generalize effectively. Recent works [20, 41] leverage the power of scaling and diversifying data collection, thereby enhancing generalization capabilities.

### A.2. Model and hyperparameters

In this section, we provide a summary of the architecture and hyperparameters used in the Encoder-Decoder Control Transformer. Our model consists of a ViT encoder, a TokenLearner, and a Control Transformer. The Control Transformer we use is composed of 8 causal attention layers with 8 attention heads, 8 cross-attention layers with 8 attention heads, and an embedding dimension of 512. The ViT encoder is ViT-B/16 and we load pretrained weights from MultiMAE [3]. Instead of using mean pooling and a linear projection layer, we employ a TokenLearner that subsamples the 196 patch tokens output by the ViT encoder to 8 tokens, which are then passed to the Transformer decoder layers.

For both Phase I and Phase II, we utilize the default AdamW optimizer [27]. For Phase I, the learning rate and batch size are set to 5e-5 and 16, respectively, while for Phase II, they are set to 1e-4 and 128, respectively. Additionally, a context length of 6 is used in all models for both training and execution. Phase I has 175M trainable parameters while Phase II has

	Training objectives	Model structure	Dual-phase
DualMind	Phase I: Self-superv.	Phase I.:Enc-Dec Control Transformer	✓
	Phase II: IL-prompt	Phase II: +XAtten.	
IL-only	IL-prompt	Enc-Dec Control Transformer +XAtten.	✗
SMART-only	Self-superv. prompt	Enc-Dec Control Transformer +XAtten.	✗
Jointly	Self-superv. + IL-prompt	Enc-Dec Control Transformer +XAtten.	✗
GATO	IL-prompt	Enc-Dec Control Transformer	✗
GATO*	IL-prompt	GATO [41]	✗

Table 4: Comparisons of different baselines.

51.1M trainable parameters. All models are trained for 10 epochs in Phase I, and 10 epochs for Phase II, with additional training details provided in Sec. A.5.

### Action embedding

**Action encoding.** To facilitate the model’s adaptation to a broader spectrum of robot action spaces, which encompass both continuous (Metaworld) and discrete (Habitat) action spaces. For each continuous action spaces, we first compand each action independently into the [-1,1] range using mu-law [41, 52]:

$$F(x) = \text{sgn}(x) \frac{\log(|x|\mu + 1.0)}{\log(M\mu + 1.0)}$$

Then we discrete them into 256 uniform bins. Discrete actions expand to 256 bins with 0 in the same way. Each action is encoded independently so we map the action space  $\mathbb{A}$  to  $\mathbb{A} \times 256$ .

**Action decoding.** After obtaining the predicted action token, we utilize a set of action heads to anticipate the encoded actions. Following this, we decode these predicted encodings into corresponding actions for the relevant tasks.

### A.3. Baselines architecture

We summary the differences between DualMind and Baselines in Table 4. The details of Baselines are listed below.

#### A.3.1 IL-only, SMART-only and Jointly.

The `IL-only`, `SMART-only`, and `Jointly` models utilize the same architecture as our structure in Phase II, with alterations made solely to the training objectives and phase. The `IL-only` model focuses solely on prompt-conditioned imitation learning during its training process. In contrast, the `SMART-only` model leverages SMART training objectives in a purely self-supervised learning context with prompt-conditioning. The `Jointly` model synthesizes these methods, employing both SMART objectives and prompt-conditioned imitation learning loss in its comprehensive training strategy.

#### A.3.2 GATO and GATO\*

**GATO\*:** GATO [41] (In this paper, we use **GATO\*** to denote) is decode-only model, which imitates expert demonstrations from a vast dataset by prompting the model with the state and action subsequence. This model has 1.18 billion parameters and was trained on massive datasets, including 94.6k episodes from Metaworld. We include its reported performance on the Metaworld benchmark for reference.

**GATO:** For a fair comparison, we used the same base model architecture (Enc-Dec Control Transformer), but replaced our XAtten.-based prompting approach with their proposed prefix prompting approach, denote as **GATO**. Similar to `IL-only`, only imitation learning loss is used to predict future actions, But replace the XAtten. module with one that prefixes the model with prompt token. Details are provided in Fig. 11.

### A.4. Data collection

**Habitat.** We collect shortest path episodes sampled from each of the 72 Gibson [54], 61 mp3d [9] and 800 hm3d [39] training scenes. These demonstrations are generated by greedily fitting actions to follow the geodesic shortest path to the nearest navigable goal object viewpoint. We hold out 5 Gibson scenes (hominy, Goffs, Hillsdale, Micanopy, and Rosser) for

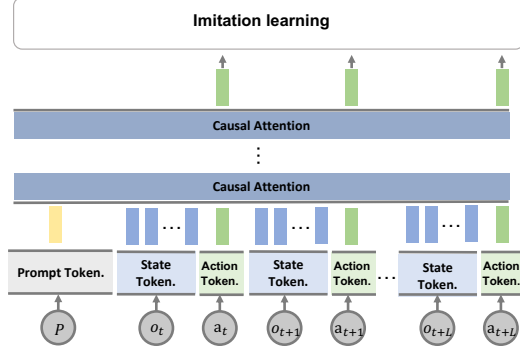


Figure 11: The architecture diagram of GATO.

the experiments of out-of-distribution. The data we collected included RGB images ( $3 \times 224 \times 224$ ), goal, and actions. We collected about 1000 episodes for each scene. Then, we divided the dataset and created the following dataset based on their intended purposes.

- **Habitat 50k.** We select all scenes of the Habitat dataset, and randomly sample about 50 episodes pre scene from the Habitat dataset. This data has 50K episodes and about  $\sim 3.26\text{M}$  interaction steps. The dataset is used in Phase I and Phase II training.
- **Habitat 10k.** We randomly select 10 scenes of Habitat scenes, and randomly sample 1000 episodes pre scene from the Habitat dataset. This data has 10K episodes and about  $\sim 0.54\text{M}$  interaction steps. The dataset is used to train the model in Phase II of the ablation study.
- **Out-of-distribution tasks.** We select 5 Gibson scenes (“Goffs”, “Hominy”, “Hillsdale”, “Micanopy”, and “Rosser”) held out, and randomly sample 10, 100, and 1000 episodes pre scenes from the Habitat dataset.

**Metaworld.** We collected data for all tasks in the MT50 [55] using scripted policies, which allowed us to generate expert demonstrations across an unlimited number of environment seeds. The data we collected included RGB images ( $3 \times 224 \times 224$ ) rendered by the physical simulator, physics engine states, and actions. We collected 2000 episodes for each tasks. We use 45 tasks in the ML45 for Phase I, and hold out other 5 tasks (hand-insert-v2, door-lock-v2, door-unlock-v2, box-close-v2 and bin-picking-v2) for the experiments of out-of-distribution. Then, we divided the dataset and created the following dataset based on their intended purposes.

- **ML45.** We select 45 training tasks of ML45 in Metaworld, and randomly sample 1000 episodes pre task from the Metaworld dataset. This data has 45K episodes and about  $\sim 3.40\text{M}$  interaction steps. The dataset is used in Phase I and Phase II training.
- **ML10.** We select 10 training tasks of ML10 in Metaworld, and randomly sample 1000 episodes pre task from the Metaworld dataset. This data has 10K episodes and about  $\sim 0.79\text{M}$  interaction steps. The dataset is used to train the model in Phase II of the ablation study.
- **Out-of-distribution tasks.** We select 5 test tasks of ML45 in Metaworld (“hand-insert-v2”, “door-unlock-v2”, “door-lock-v2”, “box-close-v2”, and “bin-picking-v2”), and randomly sample 10, 100, and 1000 episodes pre task from the Metaworld dataset.

### A.5. Training detail

**Phase I.** In Phase I, the entire model, except for the cross-attention layers (XAtten.), is trained using a self-supervised training objective on the ML45 dataset.

**Phase II.** In Phase II, we freeze the model encoder and only finetune a small part of the model, namely the Control Transformer, using imitation learning based on prompts. To encode the prompts, we use the CLIP encoder (CLIP/ViT-B/16) [35] and denote the resulting prompt sequence as  $P$ . The output sequence from each cross-attention layer is computed by  $\text{softmax}(\frac{q_H k_P^T}{\sqrt{d}})v_P$ , where  $H$  is the sequence of episodes and  $d$  is the embedding dimension. In ablation study, we use

Habitat 10K and ML10 datasets for Phase II training dataset, while for the other experiments we use the Habitat 50K and ML45 datasets as training data for Phase II.

**Out-of-distribution tasks.** In Sec. 5.3.2, we use `DualMind`, `IL-only`, and `Scratch` for out-of-distribution tasks. `DualMind` and `IL-only` model are trained beforehand and further finetuned with few-shot demonstrations. `Scratch` refers to the model that is trained on few-shot demonstrations from randomly initialized model weights. We randomly select 10, 100 and 1000 episodes for few-shot learning. We use batch size  $bs = 64$  and  $lr = 1e-4$ . We train all models for 10000 gradient steps. The data for the out-of-distribution tasks are generated in the same way as we did in Sec. A.4.

**Ablation study.** In Sec. 5.4, we use Phase I model pretrained on Habitat 50k and ML45 datasets. And the training parameters were the same as in Phase II except for the change in ablation condition and datasets.

## A.6. Evaluation detail

**Habitat.** Habitat is an immersive navigation task that provides a visually realistic environment. Our experiments are focused on the ImageNav task, in which the agent navigates towards a target position based on a goal image. The agent should stop within 1000 steps and reach a distance of 1m from the target image. To conduct our evaluation, we chose 4 scenes (Convoy, Beach, Cooperstown and Eagerville). We hold out 5 Gibson scenes (hominy, Goffs, Hillsdale, Micanopy, and Rosser) for the experiments of out-of-distribution tasks. For each scene, we randomly select three difficulty levels based on path length (EASY: 1.5-3m, MEDIUM: 3-5m, and HARD: 5-10m), resulting in a total of 300 episodes per scene. The metrics of the Habitat benchmark are listed below:

- **Success Rate(SR) and Success weighted by Path Length(SPL).** The success rate(SR) and success weighted by Path Length(SPL), proposed by [2], are estimated over 100 episodes on 4 scenes with 3 difficulty levels per scene, for a total of 1200 episodes per seed.

**Metaworld.** Metaworld is a benchmark of 50 diverse simulated manipulation tasks. We select 45 training tasks on ML45 for the experiments of evaluation, and hold out 5 test tasks ("hand-insert-v2", "door-unlock-v2", "door-lock-v2", "box-close-v2", and "bin-picking-v2") for the experiments of out-of-distribution. The metrics of the Metaworld benchmark are listed below.

- **Success Rate(SR).** We refer to the evaluation method in Metaworld [55]. The success rate is estimated over 10 seeds per task.
- **Expert Score.** The expert score is a measure of the difference between the performance of agents and experts, and is calculated as the ratio of the return obtained by agents to the expert return. We use the same expert return calculation method as GATO [41].

$$\max_{j \in [0, 1, \dots, N-W]} \left( \frac{\sum_{i=j}^{j+W-1} R_i}{W} \right)$$

where  $N$  is the total number of collected episodes for the task,  $W$  is the window size, and  $R_i$  is the total return for episode  $i$ .

## B. More experiments

### B.1. Comparisons of varying context length

We conducted experiments on different context lengths, as illustrated in Fig. 12 and Fig. 13. On the navigation tasks in Habitat, long-range temporal dependencies are important for decision-making. As a result, the model’s performance is improved progressively as the length of the context increases, as shown in Fig. 12. On the other hand, we observed that setting the context length to 6 leads to better performance on the Metaworld dataset, as demonstrated in Fig. 13. Therefore, we choose a context length of 6 as means of balancing performance and compute cost. However, if one seeks to capture long-term temporal dependence, increasing the context length may be necessary.

### B.2. Comparison with vision tokenization

We conducted a set of experiments to demonstrate the effectiveness of the multi-state tokens module (i.e., TokenLearner), by comparing it with a single-state tokens module that uses mean pooling and a linear projection layer to convert patch tokens to one token. In contrast to other ablation studies, we trained both Phase I and Phase II on the ML45 and Habitat 50k datasets.



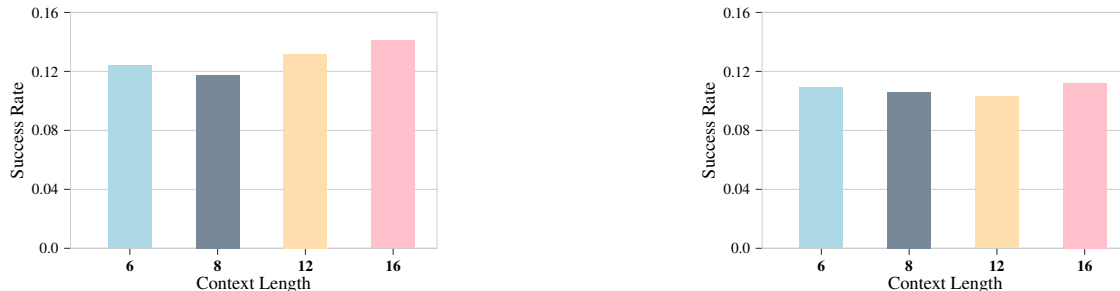


Figure 12: Comparison of varying **context length** on *Habitat*, and compare agents by Success Rate (SR) (left) and Success weighted by Path Length (SPL) (right).

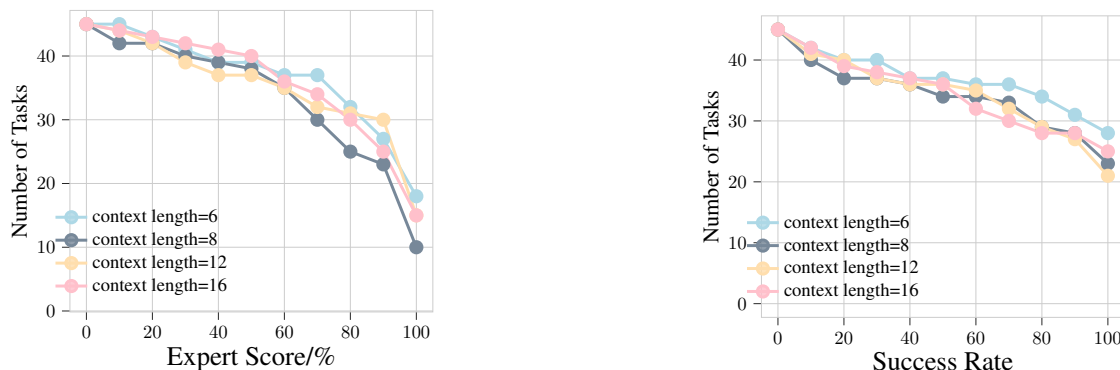


Figure 13: Comparisons of varying **context length** on *MetaWorld 45 tasks* on Percentage of Expert Score (PES) (left) and Success Rate (SR) (right).

The results in Table 5 indicate a minor difference between the two on habitat, but an average success rate difference of approximately 0.09 on ML45. These findings support our expectation that multi-state tokens can extract additional information from the encoder to improve decision-making and enhance overall learning performance.

tokenization	Habitat	Metaworld
multi-state token.	0.1239	0.802
single-state token.	0.1217	0.713

Table 5: Comparisons of tokenizations.

### B.3. Prompt conditioning discussion

**Implementation details.** In Sec. 5.4, we conducted an ablation study by comparing two prompt conditioning approaches: prefix and XAtten. prompting. The prefix approach is a conventional prompting method that splices the prompt sequences in front of the token sequences, which are directly fed into the Transformer decoder layers. In contrast, XAtten. prompting uses a cross-attention layer to fuse the prompt sequences and token sequences together. We utilized the base model that was pretrained on Habitat 50 and ML45 after Phase I. We use Habitat 10K and ML10 datasets for Phase II training dataset.

**Discussion.** In the experiments discussed in Sec. 5.4, it was found that XAtten. prompting outperforms prefix prompting. This suggests that the cross attention mechanism is effective in establishing a strong connection between prompt and token sequences, which has also been demonstrated in other recent works, such as Vima[21] and Stable Diffusion [42]. One potential limitation of prefix prompting is that the prompt token sequence may be too short to attract sufficient attention from the attention mechanism, leading to suboptimal performance. To address this, future research may explore alternative encoding methods for prompts that can better capture the information necessary for guiding the model’s output.

### B.4. Attention visualization

In Figure 14, we provide additional attention maps that reveal how DualMind tends to focus on the object being manipulated, as well as its surrounding context and relevant visual cues (such as "plate-slide-v2", "push-v2", and "hammer-v2") when performing manipulation tasks in MetaWorld. Furthermore, the attention maps show that the model focuses on the location of the item being manipulated, and then interacts with the corresponding item to complete the task. In Habitat, our model (DualMind) focuses more on exploring the scene and then locating the goal, as illustrated in Figure 14. The attention maps demonstrate that DualMind quickly identified the location of the goal image at the outset. Despite that there are obstacles blocking the shortest path, DualMind was able to eventually reach the goal.

### B.5. Performance on each tasks

We show the detailed results of all models on Metaworld and Habitat in Table 6, Table 7, and Table 8. Specifically, Table 6 presents the Habitat results, while Table 7 and Table 8 present the Metaworld results.

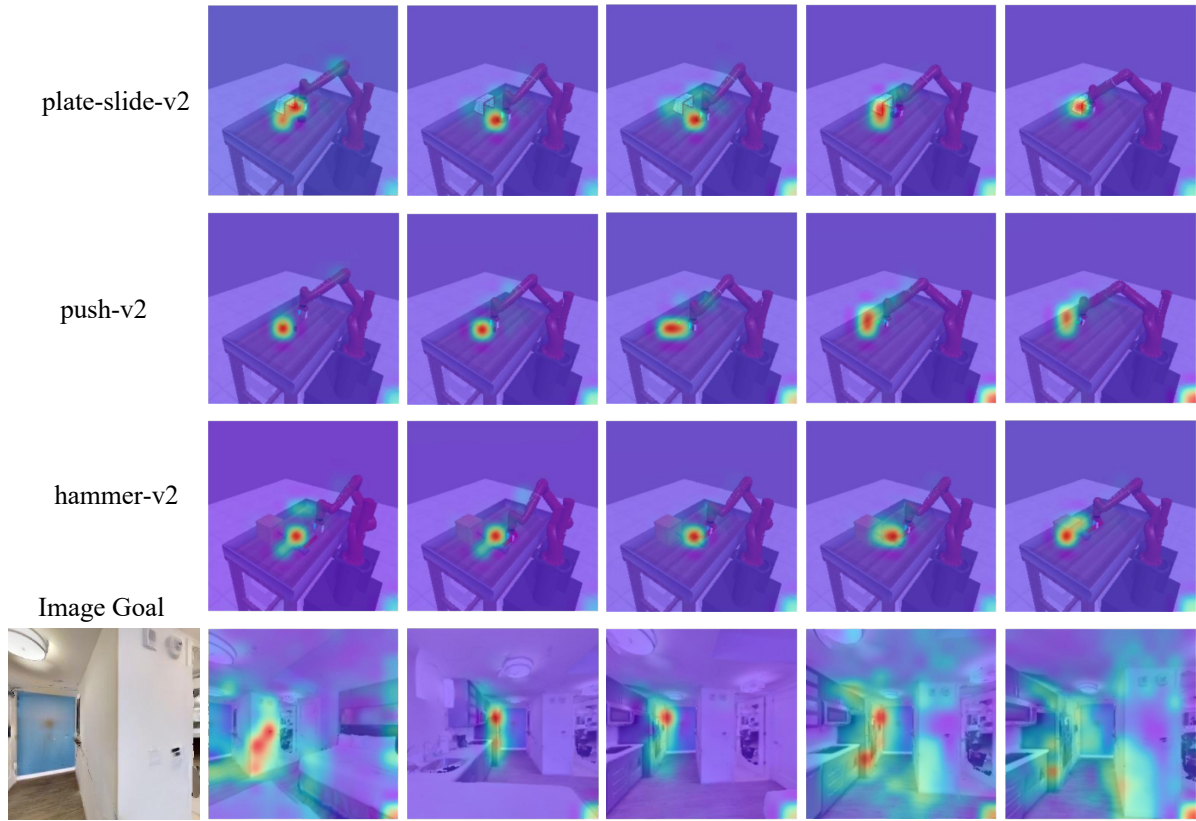


Figure 14: More attention map visualization. On Metaworld, the attention maps show that the model focuses on the location of the item being manipulated, and then interacts with the corresponding item to complete the task. On Habitat, DualMind focuses more on exploring the scene and then locating the goal.

scene	DualMind		DualMind/single	
	SR	SPL	SR	SPL
Convoy(easy)	0.143±0.021	0.124±0.017	0.160±0.026	0.115±0.027
Convoy(medium)	0.150±0.017	0.144±0.020	0.157±0.006	0.126±0.009
Convoy(hard)	0.067±0.025	0.062±0.029	0.173±0.042	0.162±0.039
Beach(easy)	0.170±0.044	0.144±0.045	0.170±0.053	0.133±0.036
Beach(medium)	0.157±0.015	0.146±0.012	0.200±0.050	0.165±0.050
Beach(hard)	0.133±0.015	0.128±0.014	0.227±0.045	0.215±0.047
Cooperstown(easy)	0.147±0.021	0.122±0.030	0.180±0.070	0.141±0.058
Cooperstown(medium)	0.110±0.035	0.103±0.027	0.190±0.046	0.175±0.044
Cooperstown(hard)	0.073±0.021	0.070±0.020	0.140±0.046	0.132±0.043
Eagerville(easy)	0.113±0.025	0.079±0.027	0.087±0.015	0.054±0.003
Eagerville(medium)	0.127±0.029	0.106±0.020	0.107±0.025	0.087±0.019
Eagerville(hard)	0.097±0.046	0.081±0.043	0.147±0.021	0.130±0.018
scene	Jointly		Jointly/single	
	SR	SPL	SR	SPL
Convoy(easy)	0.130±0.020	0.122±0.023	0.173±0.031	0.120±0.020
Convoy(medium)	0.080±0.000	0.076±0.002	0.143±0.031	0.111±0.035
Convoy(hard)	0.050±0.010	0.049±0.011	0.177±0.015	0.160±0.017
Beach(easy)	0.137±0.021	0.115±0.013	0.220±0.056	0.156±0.026
Beach(medium)	0.093±0.021	0.086±0.019	0.230±0.030	0.172±0.015
Beach(hard)	0.057±0.021	0.054±0.021	0.147±0.042	0.115±0.018
Cooperstown(easy)	0.117±0.035	0.106±0.034	0.173±0.021	0.139±0.027
Cooperstown(medium)	0.077±0.006	0.069±0.007	0.240±0.026	0.219±0.027
Cooperstown(hard)	0.080±0.020	0.076±0.022	0.183±0.032	0.167±0.032
Eagerville(easy)	0.137±0.035	0.116±0.035	0.100±0.010	0.058±0.002
Eagerville(medium)	0.097±0.015	0.085±0.016	0.180±0.030	0.112±0.021
Eagerville(hard)	0.067±0.023	0.061±0.022	0.193±0.040	0.150±0.020
scene	IL-only		IL-only/single	
	SR	SPL	SR	SPL
Convoy(easy)	0.113±0.006	0.110±0.006	0.110±0.010	0.09±0.011
Convoy(medium)	0.040±0.026	0.039±0.026	0.037±0.025	0.030±0.021
Convoy(hard)	0.027±0.006	0.027±0.006	0.053±0.012	0.043±0.011
Beach(easy)	0.103±0.055	0.095±0.051	0.070±0.010	0.052±0.010
Beach(medium)	0.053±0.035	0.050±0.033	0.050±0.010	0.035±0.0135
Beach(hard)	0.030±0.017	0.027±0.015	0.033±0.012	0.027±0.010
Cooperstown(easy)	0.080±0.010	0.074±0.010	0.103±0.015	0.076±0.01
Cooperstown(medium)	0.043±0.012	0.041±0.023	0.053±0.012	0.038±0.013
Cooperstown(hard)	0.047±0.020	0.045±0.009	0.050±0.000	0.039±0.006
Eagerville(easy)	0.067±0.025	0.064±0.025	0.063±0.021	0.042±0.023
Eagerville(medium)	0.070±0.030	0.054±0.018	0.033±0.021	0.022 ±0.01
Eagerville(hard)	0.047±0.021	0.040±0.024	0.033±0.021	0.026±0.014
scene	SMART-only		SMART-only/single	
	SR	SPL	SR	SPL
Convoy(easy)	0.113±0.006	0.088±0.003	0.007±0.006	0.007±0.006
Convoy(medium)	0.003±0.006	0.003±0.005	0.0±0.0	0.0±0.0
Convoy(hard)	0.007±0.006	0.005±0.004	0.0±0.0	0.0±0.0
Beach(easy)	0.063±0.025	0.044±0.014	0.007±0.012	0.007±0.012
Beach(medium)	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
Beach(hard)	0.010±0.010	0.008±0.009	0.0±0.0	0.0±0.0
Cooperstown(easy)	0.090±0.020	0.079±0.018	0.013±0.006	0.013±0.006
Cooperstown(medium)	0.010±0.010	0.007±0.006	0.0±0.0	0.0±0.0
Cooperstown(hard)	0.003±0.006	0.003±0.006	0.0±0.0	0.0±0.0
Eagerville(easy)	0.087±0.006	0.076±0.003	0.017±0.006	0.016±0.005
Eagerville(medium)	0.023±0.006	0.017±0.004	0.0±0.0	0.0±0.0
Eagerville(hard)	0.010±0.010	0.008±0.008	0.0±0.0	0.0±0.0

Table 6: performance on each tasks on Habitat

task	DualMind		DualMind/single		Jointly		Jointly/single	
	SR	return	SR	return	SR	return	SR	return
assembly-v2	0.9	1039.1	1.0	1276.0	0.0	198.8	0.0	477.0
basketball-v2	0.3	400.0	0.5	621.7	0.0	10.6	0.0	42.4
button-press-topdown-v2	1.0	364.9	1.0	1175.4	0.6	113.9	0.6	45.2
button-press-topdown-wall-v2)	1.0	0.0	1.0	0.0	0.7	55.2	1.0	39.5
button-press-v2	1.0	347.3	1.0	357.6	0.7	318.2	0.5	326.2
button-press-wall-v2	0.3	1373.7	1.0	1195.2	0.0	128.8	0.0	0.5
coffee-button-v2	1.0	301.0	1.0	299.2	1.0	304.3	0.6	268.6
coffee-pull-v2	1.0	429.5	1.0	407.3	0.8	303.8	0.9	324.2
coffee-push-v2	1.0	443.6	1.0	509.0	0.0	30.3	0.2	66.9
dial-turn-v2	1.0	1220.0	1.0	1203.3	0.0	61.6	0.3	17.7
disassemble-v2	1.0	615.3	1.0	553.4	0.0	220.9	0.0	213.2
door-close-v2	1.0	946.0	1.0	754.9	0.1	2956.1	1.0	1286.4
door-open-v2	1.0	1756.7	1.0	1775.4	0.1	696.5	0.9	1457.9
drawer-close-v2	1.0	61.4	1.0	81.2	0.1	4.9	0.4	24.6
drawer-open-v2	1.0	1965.0	1.0	1989.5	0.7	1517.8	0.5	1103.6
faucet-open-v2	0.3	1693.4	1.0	2192.2	0.0	1388.9	0.0	1276.4
faucet-close-v2	0.1	1624.0	0.0	1695.1	0.0	1856.1	0.2	2075.1
hammer-v2	1.0	951.9	1.0	929.5	0.0	675.8	0.1	869.2
handle-press-side-v2	1.0	839.9	1.0	808.8	0.7	307.2	1.0	877.9
handle-press-v2	1.0	671.1	1.0	782.7	0.4	195.6	0.7	427.7
handle-pull-side-v2	0.8	580.4	0.0	29.1	0.0	12.4	0.0	11.4
handle-pull-v2	0.7	182.9	0.2	177.5	0.1	70.4	0.5	139.9
lever-pull-v2	0.0	291.1	0.8	946.7	0.1	420.3	0.0	283.8
peg-insert-side-v2	0.9	990.8	0.7	909.5	0.5	1413.7	0.3	1096.3
pick-place-wall-v2	1.0	656.9	1.0	1698.2	0.0	0.1	0.0	14.0
pick-out-of-hole-v2	0.0	261.1	0.0	362.4	0.0	25.9	0.0	12.3
reach-v2	0.1	2507.2	0.0	2374.8	0.0	241.8	0.0	598.6
push-back-v2	1.0	193.7	1.0	283.7	0.0	6.5	0.0	6.4
push-v2	1.0	1264.1	1.0	1446.0	0.0	22.7	0.0	23.4
pick-place-v2	0.7	608.4	1.0	303.0	0.0	7.2	0.0	19.2
plate-slide-v2	1.0	1255.4	1.0	1214.2	0.0	269.7	0.2	417.6
plate-slide-side-v2	1.0	1281.6	1.0	1278.6	0.0	143.1	0.2	397.8
plate-slide-back-v2	1.0	1207.5	1.0	1170.6	0.8	909.0	0.4	618.6
plate-slide-back-side-v2	1.0	1340.7	1.0	1321.9	0.1	499.6	0.6	782.1
peg-unplug-side-v2	0.8	343.9	1.0	344.0	0.1	109.2	0.3	122.0
soccer-v2	0.0	336.4	0.0	329.6	0.1	180.3	0.0	164.3
stick-push-v2	1.0	1328.5	1.0	1316.9	0.0	24.2	0.5	475.5
stick-pull-v2	0.9	190.4	1.0	648.5	0.0	6.3	0.1	144.8
push-wall-v2	1.0	1387.4	1.0	1782.3	0.0	51.3	0.0	49.4
reach-wall-v2	0.5	3151.8	0.0	4190.0	0.0	341.4	0.0	714.7
shelf-place-v2	0.8	752.2	0.9	864.4	0.0	0.0	0.0	0.1
sweep-into-v2	1.0	880.6	0.8	783.1	0.0	47.5	0.0	52.0
sweep-v2	1.0	1346.6	1.0	1108.6	0.0	93.3	0.0	91.8
window-open-v2	1.0	438.6	1.0	494.6	0.5	379.8	0.2	436.8
window-close-v2	1.0	784.2	1.0	806.3	0.5	799.5	0.4	541.1

Table 7: The detailed Metaworld ML45 results of the DualMind, DualMind/single, Jointly and Jointly/single on each tasks.

	IL-only		IL-only/single		SMART-only		SMART-only/single	
	SR	return	SR	return	SR	return	SR	return
assembly-v2	0.0	252.1	0.0	169.1	0.0	197.2	0.0	189.0
basketball-v2	0.0	13.9	0.0	5.1	0.0	2.0	0.0	1.3
button-press-topdown-v2	0.0	194.3	0.0	33.3	0.0	116.4	0.0	0.1
button-press-topdown-wall-v2	0.6	88.6	0.0	1.3	0.0	35.7	0.0	14.5
button-press-v2	0.3	366.0	0.0	43.4	0.0	53.3	0.0	45.3
button-press-wall-v2	0.0	75.9	0.0	19.2	0.0	59.1	0.0	25.2
coffee-button-v2	1.0	301.0	0.0	38.6	0.6	297.3	0.0	65.1
coffee-pull-v2	0.9	365.8	0.0	13.5	0.0	11.6	0.0	11.9
coffee-push-v2	0.0	83.6	0.0	13.1	0.0	10.9	0.0	5.1
dial-turn-v2	0.0	19.3	0.0	6.7	0.0	4.4	0.0	8.3
disassemble-v2	0.0	210.9	0.0	206.6	0.0	204.8	0.0	206.1
door-close-v2	0.2	2742.7	0.0	30.4	0.0	659.8	0.2	327.4
door-open-v2	0.4	1169.6	0.0	169.6	0.0	212.3	0.0	383.7
drawer-close-v2	0.0	0.0	1.0	71.3	0.0	2.3	0.0	0.0
drawer-open-v2	1.0	1976.5	0.0	389.6	0.0	493.4	0.0	389.9
faucet-open-v2	0.1	1547.4	0.0	427.1	0.0	490.4	0.0	302.0
faucet-close-v2	0.0	1062.5	0.0	453.2	0.0	863.6	0.0	552.3
hammer-v2	0.0	588.5	0.0	528.0	0.0	263.3	0.0	585.6
handle-press-side-v2	0.6	235.4	0.8	493.9	0.0	26.3	0.0	28.7
handle-press-v2	0.0	86.1	0.0	18.5	0.0	34.3	0.0	22.9
handle-pull-side-v2	0.3	12.1	0.0	10.7	0.0	2.1	0.0	2.5
handle-pull-v2	0.2	82.5	0.0	6.4	0.0	14.4	0.0	4.4
lever-pull-v2	0.0	350.2	0.0	24.2	0.0	125.0	0.0	90.0
peg-insert-side-v2	0.9	1289.4	0.0	2.2	0.0	2.4	0.0	1.7
pick-place-wall-v2	0.0	27.4	0.0	0.0	0.0	0.0	0.0	0.0
pick-out-of-hole-v2	0.0	19.1	0.0	3.4	0.0	3.1	0.0	1.2
reach-v2	0.0	195.0	0.0	122.3	0.0	144.5	0.0	195.3
push-back-v2	0.0	5.0	0.0	1.7	0.0	2.3	0.0	1.7
push-v2	0.0	21.8	0.0	10.2	0.0	3.8	0.0	4.6
pick-place-v2	0.0	6.7	0.0	3.0	0.0	2.3	0.0	3.2
plate-slide-v2	0.3	393.9	0.0	72.1	0.0	97.2	0.0	44.2
plate-slide-side-v2	0.0	45.8	0.2	419.0	0.0	20.6	0.0	5.0
plate-slide-back-v2	0.7	1027.7	0.0	43.3	0.0	48.3	0.0	21.8
plate-slide-back-side-v2	0.0	200.1	0.0	1185.7	0.0	25.5	0.0	25.6
peg-unplug-side-v2	0.2	131.4	0.0	3.6	0.0	2.7	0.0	2.8
soccer-v2	0.0	21.0	0.0	38.1	0.0	3.4	0.0	6.9
stick-push-v2	0.0	16.7	0.0	5.7	0.0	1.9	0.0	3.1
stick-pull-v2	0.0	6.6	0.0	5.8	0.0	2.2	0.0	6.9
push-wall-v2	0.0	24.4	0.0	18.0	0.0	3.9	0.0	5.0
reach-wall-v2	0.0	558.7	0.0	305.2	0.0	159.9	0.0	435.3
shelf-place-v2	0.0	214.2	0.0	0.0	0.0	0.0	0.0	0.0
sweep-into-v2	0.0	55.3	0.0	8.3	0.0	10.7	0.0	9.1
sweep-v2	0.0	83.6	0.0	15.9	0.0	17.1	0.0	13.7
window-open-v2	1.0	449.5	0.0	101.4	0.0	91.4	0.0	92.7
window-close-v2	0.1	462.9	0.0	10.9	0.0	374.4	0.0	216.0

Table 8: The detailed Metaworld ML45 results of the IL-only, IL-only/single, SMART-only and SMART-only/single on each tasks.