

# EGOPLAN: TOWARDS EFFECTIVE EMBODIED AGENTS VIA EGOCENTRIC PLANNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We explore leveraging large multi-modal models (LMMs) and Text2image models to build a more general embodied agent. LMMs excel in planning long-horizon tasks over symbolic abstractions but struggle with grounding in the physical world, often failing to accurately identify object positions in images. A bridge is needed to connect LMMs to the physical world. The paper proposes a novel approach, egocentric vision language planning (EgoPlan), to handle long-horizon tasks from an egocentric perspective in varying household scenarios. This pipeline leverages a diffusion model to simulate the fundamental dynamics between states and actions, discusses how to integrate computer vision related techniques like style transfer and optical flow to enhance ability of modeling spatial states and generalization across different environmental dynamics. The LMM serves as a planner, breaking down instructions into sub-goals and selecting actions based on their alignment with these sub-goals, thus enabling more generalized and effective decision-making. By using LMM, we can output text actions, using a series of mechanisms such as reflection to perform high-level task decomposition and low-level action output end-to-end. Experiments show that EgoPlan improves long-horizon task success rates from the egocentric view compared to baselines across household scenarios.

## 1 INTRODUCTION

The advent of large language models (LLMs) (et al., 2024b; Touvron et al., 2023) and large multi-modal models (LMMs) (202, 2023; Girdhar et al., 2023; Zhang et al., 2023a; Zhu et al., 2023) has revolutionized the field of artificial intelligence. Their strong reasoning (Wang et al., 2023b; Wei et al., 2023) and powerful generalization capabilities allow them to be directly applied in various scenarios. In the next step toward artificial general intelligence (AGI), researchers are considering enabling large models (LMs), especially LMMs, to break through the world expressed by text and images to interact with the physical world. They aim to build a general embodied agent that intelligently interacts with the physical world.

LMMs exhibit impressive long-horizon planning over symbolic abstractions (Wake et al., 2024), yet struggle with grounding text in the physical world, often failing in precise object localization. While LMMs understand what to do, they lack understanding of how the world functions, necessitating a world model to bridge this gap. Two potential solutions exist: implicitly integrating dynamics via extensive fine-tuning on state-action sequences (Driess et al., 2023; et al., 2023), which demands substantial resources, or explicitly employing pre-trained generative world models (e.g., Text2image/video) as auxiliary tools (Radford et al., 2021; Saharia et al., 2022). Prior work (Black et al., 2023; Du et al., 2023b) suggests that these models can inject world knowledge by predicting future observations or trajectories. This work investigates the latter approach, exploring the potential of leveraging LMMs and Text2image models for more general embodied agents.

Existing approaches (Du et al., 2023a; Zhou et al., 2024) using Text2image/video models as world models for decision-making face limitations. First, their focus on fully observable object manipulation tasks is atypical of real-world scenarios and their adaptability to partially observable settings is unclear. For instance, methods requiring multi-step image generation (Black et al., 2023; Du et al., 2023b) suffer from error accumulation in partially observed environments like autonomous driving. Second, their frameworks exhibit limited capability in: (i) task-specific low-level policies with potential for collapse upon new dynamics; (ii) coarse-grained text action representations hindering the mapping

to fine-grained state transitions, especially in complex, partially observable tasks; and (iii) the lack of individual entity motion patterns, limiting generalization to novel environments with different dynamics within the same task category. We aim for generalization to varying dynamics within fixed household scenarios.

To address these limitations, we propose Egocentric Vision Language Planning (**EgoPlan**), a general embodied agent for long-horizon egocentric tasks. Recognizing the rich action and state transition information in optical flow (Ko et al., 2023; Yang & Ramanan, 2020), we integrate it into our world model for enhanced spatial understanding in navigation and object motion prediction in manipulation, contrasting with traditional text-based actions. Furthermore, to handle visual style variations across different simulated home environments, we employ LoRA fine-tuning to enable adaptation to diverse visual distributions while preserving learned motion patterns. This enhances fine-grained texture modeling and generation across scenes, allowing for transfer to new environments with limited samples, achieving a style transfer-like effect.

We conduct a comprehensive evaluation and analysis of each module of the embodied agent. Empirically, we demonstrate the high quality of image generation by the dynamics model and the high accuracy of optical flow prediction. Subsequently, we verify the dynamics model’s effectiveness in aiding decision-making in more complex tasks. Lastly, we confirm the method’s generalization capabilities in a different environment. Our major contributions are summarized as follows:

- We have collected a dataset on Virtualhome, which views high-level manipulation/navigation actions of the agent in Virtualhome as trajectories and provides egocentric observations each time-step and fine-grained action information, which will provide data support for navigation and manipulation tasks in the embodied environment. See Section 3 for details.
- We propose **EgoPlan**, a framework for complex task planning that combines LMM and a dynamics model that predicts an egocentric view of the next time step and the subgoal is completed. We also introduce optical flow into the dynamics model and borrow the idea of style transfer in computer vision and adopt the LoRA (Hu et al., 2021) model to achieve few-shot generalization in different embodied scenarios.
- For the action selection and decision-making module, we employ the LMM as the execution module in both the high-level task decomposition and low-level action selection components. We utilize a series of reflection and summarization mechanisms to accomplish tasks, while also ensuring the agent inherits this ability of generalizing the downstream policies to new dynamics. Experiments on comprehensive tasks demonstrate the effectiveness of our framework through LMM+dynamics model planning.

## 2 RELATED WORK

In this section, we present a brief overview of related work. More discussions are in Appendix B.

### 2.1 DYNAMIC MODEL AND WORLD MODEL FOR DECISION-MAKING

The world model is used to model the dynamics of the environment. It is crucial for building autonomous agents and enabling intelligent interactions in various scenarios. However, developing a precise world model remains a significant challenge in model-based decision-making. The advancements in diffusion-based world models are reshaping how we model physical motion laws in real-world settings, particularly in robotics. UniPi (Du et al., 2023a) frames the decision-making problem in robotics as a Text2video task. The generated video is fed into an inverse dynamics model (IDM) that extracts underlying low-level control actions, which are executed in simulation or by a real robot agent. Video Language Planning (VLP) (Du et al., 2023b) introduces a novel method for task planning that integrates video generation with tree search algorithms. This methodology lets robots plan over longer horizons by visualizing future actions and outcomes. Unlike previous works, SuSIE (Black et al., 2023) leverages pre-trained image-editing models to predict the hypothetical future frame. A low-level goal-reaching policy is trained on robot data to reach this hypothetical future frame. Since one goal frame prediction does not require the model to understand the intricacies of the robot’s low-level precisely dynamics, it should facilitate transfer from other data sources such as human videos. RoboDreamer (Zhou et al., 2024) advances the field by utilizing video



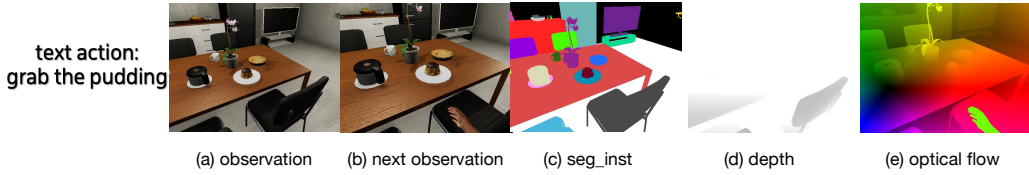


Figure 1: An illustration sample in VH-1.5M, which includes current image observation, next image observation given the text action, semantic segmentation map, depth map, and optical flow map.

diffusion to formulate plans combining actions and objects, solving novel tasks in unexplored robotic environments.

## 2.2 EMBODIED AGENT WITH LMMs

Recent methods use LMMs to assist planning and reasoning in simulation environments (Fan et al., 2022; Wang et al., 2023a; Yao et al., 2023) and robot learning (Ahn et al., 2022; Liang et al., 2023; Zeng et al., 2022). LMMs are also applied to help robot navigation (Parisi et al., 2022; Majumdar et al., 2020) and manipulation (Jiang et al., 2022; Ren et al., 2023; Khandelwal et al., 2022). Among them, ReAct (Yao et al., 2023) uses chain-of-thought prompting by generating both reasoning traces and action plans with LMMs. SayCan (Ahn et al., 2022) leverages the ability of LLMs to understand human instructions to make plans for completing tasks without finetuning LLMs. Voyager (Wang et al., 2023a) leverages GPT-4 to learn and continually discover skills during learning. While these studies demonstrate encouraging outcomes, they depend significantly on the inherent capabilities of powerful large language models (LLMs), which poses challenges for their application to smaller language and multimodal models (LMMs) with limited reasoning abilities.

## 3 VH-1.5M DATASET

Existing vision-language-action datasets, such as RT-X (et al., 2024a) and RH20T (Fang et al., 2023), often utilize static views to mitigate perspective change issues, providing "fixed camera" observations suitable for fully observable task planning where all manipulable objects are within a constant field of view. In contrast, datasets like ALFRED (Shridhar et al., 2020) and ProcTHOR (Deitke et al., 2022) employ egocentric views, introducing significant perspective changes and necessitating embodied task planning under partial observability, where current observations may be insufficient for task completion without viewpoint adjustments for information gathering (e.g., navigating to an unseen object). Our dataset distinguishes itself from other embodied navigation datasets by incorporating agent motion trajectories (e.g., grasp, put) and the coupled perspective and hand position changes. To address the aforementioned limitations, we introduce the VH-1.5M dataset, built upon the VirtualHome environment (Puig et al., 2018; 2020).

We construct our dataset VH-1.5M in the VirtualHome environment, which comprises 50 distinct houses. Each house contains approximately 300 interactive objects, and the embodied agent can perform more than 10 actions. Note that the VirtualHome environment is a simulator tailored for embodied agents, offering a detailed simulation of a residential living scenario. It enables a range of household tasks, such as navigation and object manipulation.

The VH-1.5M dataset is organized in a structured manner, encapsulating the relationship between actions, houses, agents, and trajectories. Each task sequence entry follows a hierarchical structure, for example: "/open/house\_0/Female4/2\_fridge" (female4 open the fridge2 in house0).

**Dataset Details:** The VH-1.5M dataset consists of:

- 13 Actions: Various physical actions and interactions for agents within the houses. These action instructions are high level and can be completed in a sequence of time steps, such as walk to microwave.
- 50 Houses: Uniquely designed houses with diverse layouts and object placements.
- 4 Agents: Four distinct agents (simulated humans), each capable of performing the full range of actions.
- 1.5M Samples: Dataset has numerous detailed sequences, each executing one action instruction. Information from each step in the sequence is stored as one sample. One example is shown in Figure 1.

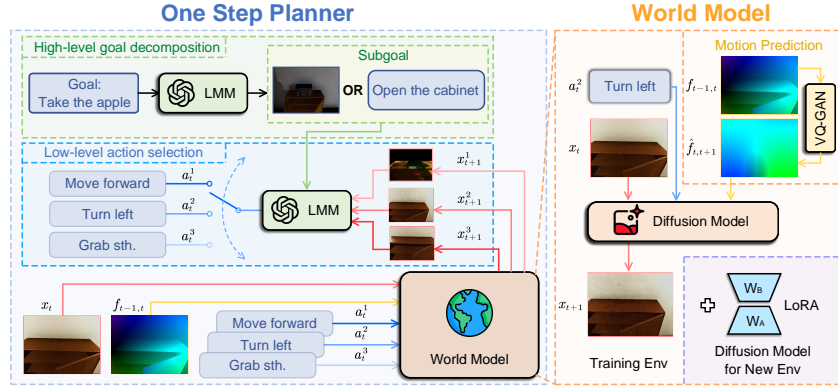


Figure 2: Overview of EgoPlan. The left side features a one-step planner that provides the agent with decision-making capabilities, while the right side includes a world model (dynamics model) that provides the agent with an understanding of the current environment.

More details of the dataset can be found in the Appendix D, and *we will open-source the dataset*.

## 4 METHOD

Our embodied agent, EgoPlan, takes visual observation  $x_t$  of the scene at the current timestep  $t$  and a natural language goal  $g$  as inputs and outputs an action  $a_t$  to interact with the environment. Note that the  $x_t$  only partially represents the current environment state. In addition, the agent uses encapsulated skills as actions, such as moving forward, turning, and grabbing objects. For problem settings, the decision-making environment is typically characterized as a Partially Observable Markov Decision Process (POMDP) (Smallwood & Sondik, 1973), defined as a tuple  $(\mathcal{O}, \mathcal{A}, p, r, \gamma)$ . In our pipeline, we define egocentric observation  $x_t$  as partial observation  $\mathcal{O}$ , and we will train a text2image model as dynamic model to model dynamic processes  $p(o_t|o_{t-1}, a_{t-1})$ . For EgoPlan agent, we model the actions in a Markov process using either textual descriptions as  $a_t = l_t$  or a more fine-grained description of actions: optical flow, which can be denoted as  $a_t = f_{t,t+1}$ .

EgoPlan consists of two parts, as illustrated in Figure 2. The first is a dynamics model that gives the agent the concept of the current environment, and the other is the planner that endows the agent with decision-making capabilities. Intuitively, we humans first envision the outcomes of each action in our minds, and then, by comparing the results, we make the best decision. In the same way, we use a dynamic model to create an egocentric scenario where different actions can be taken, which is then fed into LLM to determine which action is more reasonable.

### 4.1 DIFFUSION-BASED DYNAMICS MODEL

#### 4.1.1 LEARNING DYNAMICS

From a first-person perspective, the view after two or more steps may be completely different, making it difficult to model. Therefore, we aim to model the fundamental dynamics model,  $p_\theta(x_{t+1}|x_t, a_t)$ , for one-step planning usage. In more detail, we want to generate a new image  $x_{t+1}$ , representing the next state given the current visual observation  $x_t$  and the text of the action  $a_t$ . Then, we cast our eyes on the Text2image model and resort to the diffusion model for modeling specifically. It has an irreplaceable advantage in easily incorporating other modalities as a condition.

Although the open-sourced diffusion model (Ho et al., 2022; Luo et al., 2023),  $p_\theta(x_{\text{tar}}|x_{\text{src}}, l)$ , trained on a wealth of online videos, has demonstrated the ability to predict the future, their generated results are hard to control, and most are only semantically reasonable. Moreover, most of the text in the pre-trained dataset consists of image descriptions  $l$  rather than action instructions  $a$ . Therefore, supervised fine-tuning is adopted based on our VH-1.5M dataset to better model the dynamics,  $p_{\theta_{\text{ft}}}(x_{t+1}|x_t, a_t)$ . Formally, the training objective is given by:

$$\mathcal{L}_{\text{MSE}} = \left\| \epsilon - \epsilon_\theta \left( q \left( x_{t+1}^{(k)} | x_t, a_t \right), k \right) \right\|^2 = \left\| \epsilon - \epsilon_\theta \left( \sqrt{\alpha_t} x_t + \sqrt{1 - \alpha_t} \epsilon(a_t) \right) \right\|^2 \quad (1)$$

where  $\epsilon_\theta$  is a learnable denoising model for reverse process,  $k$  is denoising steps, and  $\alpha_t$  are a set of  $K$  different noise levels for each  $k \in [1, K]$ , and  $x_t, a_t$  separately represent the current observation image and action description text. In practice, we'll use Instructpix2pix as the backbone network; see the Appendix I for training details. However, we find it difficult to generalize directly to other

environments since our dataset only includes VirtualHome scenes. The difference between two environments, such as Habitat 2.0 (Savva et al., 2019; Szot et al., 2022) and VirtualHome, primarily lies in their different motion patterns for the same action and distinct visual styles. Especially for the former, the motion pattern, such as the amplitude of the same action, performed by agents in a different environment can be unpredictable.

#### 4.1.2 GENERALIZATION

We want to improve the model’s generalization ability from a different perspective. In other words, instead of enhancing generalization through big data and large models, we aim to explicitly address the differences between environments such as the visual style of indoor environments and the definition of action amplitudes at the methodological level.

**Motion Regularization.** Firstly, we must combine the motion information into the diffusion model to distinguish the different motion patterns. Optical flow has thus caught our attention. It refers to the pattern of apparent motion of image objects between two consecutive frames caused by objects or camera movement. In optical flow maps, colors represent the direction of motion, and the depth or intensity of the colors indicates the magnitude of the motion, which is a general feature across different environments.

However, in practice, in the absence of the next observation, we cannot obtain the current optical flow,  $f_{t,t+1}$ . Inspired by other motion estimation works (Chen & Koltun, 2016; Zach et al., 2007), we assume motion consistency holds over short intervals, meaning abrupt changes do not occur. Consequently, the consecutive optical flow maps are highly correlated, allowing us to predict the current optical flow map using the previous map. The previous map is calculated from the previous two frames and reflects the actual motion pattern in the current environment.

We notice that optical flow generation does not require complex texture generation, and it is expected not to cause a significant delay in the pipeline. Therefore, we adopt a less powerful but lightweight generative model, VQ-GAN (Esser et al., 2021), and train it on our dataset to predict the optical flow map. Empirically, the generalization ability to predict optical flow is much better than predicting actual images. Formally, the training objective is given by:

$$\min \mathcal{L}_{VQ}(E, G, Z) = \|x - \hat{x}\|_2^2 + \|\text{sg}[E(x)] - z_q\|_2^2 + \beta \|\text{sg}[z_q] - E(x)\|_2^2,$$

where  $E$  is the encoder,  $G$  is the generator,  $Z$  represents the latent space,  $x$  is the input image,  $\hat{x}$  is the reconstructed image,  $z_q$  is the quantized latent vector,  $\text{sg}$  denotes the stop-gradient operator, and  $\beta$  is a hyperparameter that balances the commitment loss.

*In summary, we use a simple model to predict motion patterns and then a more complex model to reconstruct real textures based on motion patterns.* Therefore, we adopt ControlNet (Zhang et al., 2023b) to incorporate the optical flow map,  $f_{t,t+1}$ , into the default diffusion model,  $p_{\theta_{\text{sf}}}(x_{t+1}|x_t, a_t, f_{t,t+1})$ . Only the ControlNet part needs to be trained on VH-1.5M at this stage. The training details of VQ-GAN and ControlNet can be found in Appendix I. Formally, the training objective is given by:

$$\mathcal{L}_{\text{MSE}} = \left\| \epsilon - \epsilon_{\theta} \left( q \left( x_{t+1}^{(k)} | x_t, a_t, f_{t,t+1} \right), k \right) \right\|^2 \quad (2)$$

$$= \left\| \epsilon - \epsilon_{\theta} \left( \sqrt{\bar{\alpha}_t} x_t + \sqrt{1 - \bar{\alpha}_t} \epsilon | a_t, f_{t,t+1} \right) \right\|^2. \quad (3)$$

**Style Transfer.** Secondly, we use LoRA to fine-tune the diffusion model for visual style transfer. Note that LoRA requires very little data, just about 20 of samples. Normally, it is convenient to collect data on such a scale in new environments. We expect the model to achieve generalization with as little effort as possible. In Section 5.2, we can find the role of LoRA method in maintaining the action pattern of the model between different environments, while flexibly transferring the style of fine-grained observation images.

## 4.2 PLANNING WITH DYNAMICS MODEL

To avoid further training in new environments, we prompt the LMM GPT-4V, as the planner. The LMM needs to be responsible for high-level goal decomposition as well as low-level action selection. Meanwhile, the pre-trained dynamics model can help LMM better understand the world.

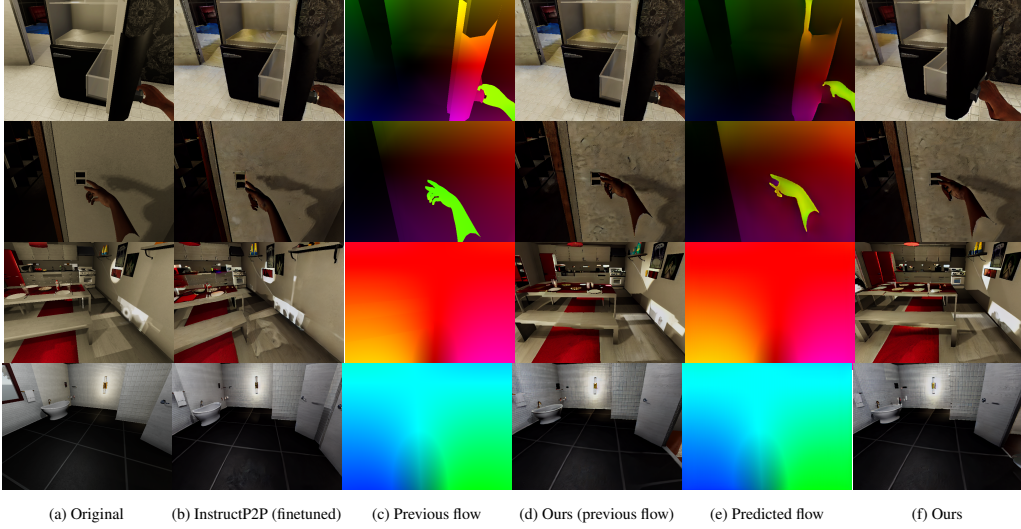


Figure 3: Examples of the generated image of the next observation in VirtualHome. The tasks from rows 1 to 4 are: close the fridge, switch off the light, turn left, and turn right.

#### 4.2.1 GOAL DECOMPOSITION

For long-term complex tasks, subgoal decomposition is crucial. Subgoals can be represented as text ( $g_{\text{tar}}$ ) or images ( $x_{\text{tar}}$ ). For text-based subgoals, we prompt the LMM for a plausible one. Additionally, we train a diffusion model,  $p_{\theta_{\text{sft}}}(x_{\text{tar}}|x_t, g_{\text{tar}})$ , to generate image-based subgoals conditioned on the text subgoal and current observation. While prior work (Black et al., 2023; Zhou et al., 2024) uses diffusion models serially to predict subgoal state images for long-horizon manipulation tasks, generating subgoal scene images for composite manipulation and navigation tasks, particularly navigation, presents a greater challenge. This is due to the substantial changes in the entire image scene and the joint positions of numerous objects required, demanding a strong understanding of spatial attributes beyond simple object-centric image editing. Consequently, predicting subgoal images can be less precise than predicting the next observation. We plan to investigate the impact of different subgoal types on task performance (Section 5.4).

#### 4.2.2 ONE-STEP PLANNER

Since we can only ensure that the prediction for the next step is relatively accurate, we adopt a one-step planning method. In more detail, we utilize the pre-trained dynamics model to predict the visual outcomes of all the actions in the next state. Once the text/image-based subgoal is obtained, we send the subgoal and all the visual outcomes to the LMM. Then, we prompt it to compare all the potential outcomes with the subgoal and determine which action can bring the agent closer to the goal. So the process of goal decomposition and one-step planner is equivalent to the following formula.

$$\{G_0, G_1, \dots, G_n\} = \text{LMM}(s_0, \text{task}) \quad (4)$$

$$a^* = \arg \min_{a \in A} d(f(s_t, a), G \in \{G_0, G_1, \dots, G_n\}) \quad (5)$$

In the aforementioned equations,  $\{G_0, G_1, \dots, G_n\}$  refers to a series of subgoals that are decomposed from the task using LMM.  $f$  is the dynamic model and  $d$  is the distance metric function, in our pipeline, GPT4V judges how far the target is from the dynamic model prediction. It is noteworthy that, in selecting the optimal action for one-step planning process, inspired by Tan et al. (2024); Zhai et al. (2024), we utilize LMM to generate low-level actions in contrast to reinforcement learning or imitation learning algorithms. In this context, we leverage the comprehension capabilities of LMM to ensure the generalization of the low-level action in cross-environment decision-making. We also employing mechanisms like React (Yao et al., 2023) and Reflexion (Shinn et al., 2023) to enhance the agent’s performance, which are shown in Appendix H. The prompt of task-decomposition and low-level action selection has been listed in Appendix G. Black et al. (2023) has discussed the generalization of objects concerning various operational targets; however, the generalization of underlying policy networks based on reinforcement learning or imitation learning algorithms, particularly in response to changes in the entire environmental scene—especially in navigation tasks, the ability of the pipeline still requires improvement. We will further discuss the experimental outcomes related to this in Sections 5.2 and 5.4.





(a) Enclose the fridge (b) Go through door (c) Shut off the PC (d) Take hold of pillow (e) Switch off the light (f) Shut the stove (g) Open the cabinet

Figure 4: Examples of the generated image subgoals. The first row is the original image, and the second row is the image subgoal generated based on the text subgoal.

## 5 EXPERIMENT

In this section, we comprehensively evaluate and analyze each module of the embodied agent. We first evaluate the quality of image generation using the world model and the quality of optical flow prediction. Secondly, we evaluate whether our world model can assist task planners in completing more complex tasks. Finally, we assess the generalization of our method. In addition to the below experiments, we also do a series of works to discuss the **complexity of the system** to explain why we do one-step planning. See the Appendix J for detailed analysis.

### 5.1 VISUAL QUALITY

We adopt two metrics, FID (Heusel et al., 2018) and user score, to evaluate the visual quality of the generated image of the world model. For models, **InstructP2P (pre-trained)** is the default model of InstructP2P. **InstructP2P (fine-tuned)** is the model fine-tuned on our dataset. **Ours (previous flow)** is the world model that conditions on the previous optical flow map, while **Ours** is conditioned on the predicted optical flow map. Note that the validation set of VH-1.5M has around 5k samples.

**FID Score.** FID is a standard metric measuring the distance of two image distributions using the inception model. The smaller the FID is, the more similar the two images are. Table 1 shows the FID score of our model and baselines. We can see that using existing diffusion models as world models is ineffective because their training data often lacks state transition-related data. Meanwhile, introducing an optical flow map, which serves as motion pattern information, significantly enhances the generation results. In addition, world models based on predicted optical flow are slightly better than those based on the optical flow of the previous frame.

**User Study.** We also conduct a user study on the accuracy of world models for image generation. For the criterion, users judge the correctness of the direction and amplitude of the executed action. Each user investigates a total of 1000 samples from the validation set. There are 8 users participating in the survey in total. Our user study, shown in Table 2, again verifies our predicted optical flow can help generate higher-quality images.

**Analysis.** As illustrated in Figure 3, InstructP2P (fine-tuned) generates the scene of steering in the wrong direction. However, this flaw can be greatly improved by incorporating optical flow information. Moreover, it is observed that the dynamics of closing the refrigerator can be more accurately predicted if the prediction of the motion pattern is considered. More examples can be seen in Appendix E.

### 5.2 VIRTUALHOME TASKS

**Results.** To demonstrate that our world model can well assist the LMM in task planning, we evaluate various methods on 12 tasks in VirtualHome environment, each task described by an instruction and can be broken down into a number of subtasks. Each task instruction, subtasks and some experiment

Table 1: FID score comparison with other models on the validation set. It is calculated between the predicted observation and ground truth. The lower the number, the better the quality of the image.

Model	Mean	Variance
InstructP2P (pre-trained)	13.65	0.10
InstructP2P (fine-tuned)	1.06	0.05
Ours (previous flow)	0.83	0.03
Ours	<b>0.82</b>	0.03

Table 2: User score of the user study. The user score is the percentage of images that users consider to meet the criteria out of the total 1000 images. The higher the number, the better the quality of the image. The evaluated images are from the validation set.

Model	Mean	Variance
InstructP2P (fine-tuned)	54.10%	1.53%
Ours (previous flow)	69.35%	1.34%
Ours	<b>74.93%</b>	2.57%

Table 3: The average length of completed subtasks on 12 tasks for all the methods. Tasks 1-6 occur inside one room, while tasks 7-12 take place in two rooms. This metric measures the average number of subtasks completed per execution after 100 executions of each task. **We reported the task completion rate in the Appendix K.**

	GPT4+React	GPT4V	React	Reflexion	GPT4V+P2P	GPT4V+OF	SuSIE	GR-MG	Ours(text goal)	Ours(image goal)
take and place	0.11	0.26	0.57	0.87	1.21	1.64	1.42	1.61	1.63	<b>1.88</b>
take and put1	0.12	0.34	0.65	0.80	1.22	1.34	0.98	1.68	1.75	<b>2.02</b>
take and put2	0.21	0.34	0.59	0.76	1.32	1.47	1.41	1.63	1.69	<b>1.91</b>
take and drink	0.08	0.46	0.81	0.79	1.19	1.47	1.39	1.77	1.99	<b>2.11</b>
turn on sit	0.10	0.31	0.75	0.81	1.29	1.51	1.31	1.68	1.71	<b>2.00</b>
put apple	0.09	0.35	0.69	0.97	1.18	1.61	1.69	1.86	<b>1.93</b>	<b>1.97</b>
take and place2	0.16	0.45	0.66	0.96	1.28	1.50	1.23	1.75	1.81	<b>1.88</b>
take and place3	0.17	0.34	0.63	0.86	1.14	1.57	1.1	1.61	2.05	<b>2.12</b>
take and put3	0.15	0.33	0.74	0.96	1.11	1.55	1.17	1.83	1.81	<b>2.01</b>
take open and put	0.12	0.38	0.64	0.84	1.22	1.46	1.15	1.93	1.77	<b>1.99</b>
take put and open	0.12	0.29	0.66	0.87	1.30	1.58	1.62	1.74	1.89	<b>1.96</b>
take and put4	0.21	0.35	0.71	0.86	1.28	1.68	1.56	1.64	1.69	<b>1.81</b>

details can be found in Appendix C. Each task is tested 100 times, and the maximum step in one episode is 80. For each of the 12 tasks, we abbreviated the task names for convenience. For example, the instruction of task 1, "take the bread from the toaster and place it on the plate on the table," consists of four subtasks: a) walk to the toaster, b) grab the bread, c) walk to the plate, and d) place the bread on the plate. We use "take and place" to refer to task 1.

These 12 instructional tasks are comprised of multiple sequential sub-tasks. For baselines, we use GPT4 combined with React (Yao et al., 2023) as the task planner and policy, denoted as **GPT4+React**, and it takes input as the JSON format text environment description. We also directly use GPT-4V to make decisions, denoted as **GPT4V**, and we also combined GPT4V with React (Yao et al., 2023) and Reflexion (Shinn et al., 2023) as the task planner and policy. When employing the Reflexion algorithm, its actor component is based on the React algorithm. These two baselines are denoted as **React** and **Reflexion**. For ablation baselines, we use the fine-tuned InstructP2P as the world model, denoted as **GPT4V+P2P**. The world model that conditions on the previous optical flow map is denoted as **GPT4V+OF**.

As shown in Table 3, the dynamic model significantly improves the GPT-4V ability on various long-horizon tasks. Moreover, the inclusion of optical flow information enhances the accuracy of image generation and further improves task planning performance. The results also demonstrate the effectiveness of the predicted optical flow map.

**Image Subgoal vs. Text Subgoal.** In this part, we analyze the impact of different types of subgoals on tasks. During the goal decomposition process, the text subgoal directly outputted by the LLM task planner represents a high-level, coarse-grained description. If our method can generate images of the scene at the completion time of the subgoal, a more detailed, fine-grained description can be obtained. This might enhance the action selection ability that relies on the quality of the subgoal.

When employing images as subgoals, our approach contrasts with methods like SuSIE (Black et al., 2023) and GR-MG (Li et al., 2025), which generate actions based on these subgoals using a one-step planning world model. Instead, we leverage an LMM for an end-to-end pipeline encompassing both task decomposition and action selection, diverging from SuSIE’s goal-conditioned behavioral cloning (GCBC) for low-level policy and GR-MG’s goal-conditioned vision-language-action model. As shown in Table 3, our method (denoted as **Ours**) outperforms SuSIE (**SuSIE**) and GR-MG (**GR-MG**), particularly in long-horizon composite task planning involving substantial perspective shifts and the necessity for subgoal reasoning. The robustness and efficacy of our pipeline for extended tasks are further evaluated through a comparison of completion rates against several baselines on VirtualHome tasks, detailed in Appendix K.

Specifically, we have trained an InstructP2P model based on VH-1.5M to generate the image when the subgoal is completed, with the generation results illustrated in Figure 4. The decision-making results in Table 3 show that fine-grained subgoal description is better than coarse-grained description, even if the generated image is not that accurate.

We also conduct a user study to evaluate the visual quality of the generated image-based subgoals. More details can be found in the Appendix F.

**Real-world Experiments.** We also conduct real-world experiments. We use the **Qwen2.5-VL** models (Bai et al., 2025) and GPT4V as the LMMs for the experiments, and compared the results of **Cosmos** (Agarwal et al., 2025) as the world model. At the same time, we also compared the success rates of a series of reinforcement learning and imitation learning methods in terms of tasks. Detailed in Appendix L.

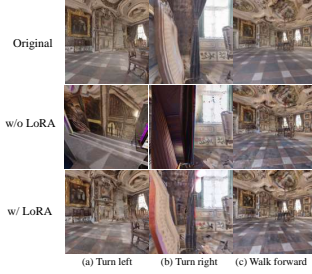


Figure 5: Examples of the generated images of the next observation in Habitat 2.0.

### 5.3 MOTION PATTERN

As mentioned before, we cannot obtain the optical flow from the current timestep to the next timestep. Therefore, we adopt the VQ-GAN model to predict the current optical flow map. The examples of prediction can be found in Appendix E. We can also found that the VQ-GAN trained on the VH-1.5M dataset can easily generalize to other environments, this is because the optical flow map is a universal feature and does not require the prediction of complex textures.

The average endpoint error (AEE) specifically measures the average distance between two motion vectors at the pixel level. As illustrated in Table 5, the gap between the predicted optical flow map and ground truth is narrower than that between the previous flow map and ground truth (current optical flow map). In addition, the model trained on VirtualHome can still predict optical flow maps in Habitat 2.0 and AI2-THOR (Kolve et al., 2017). This confirms the effectiveness and generalization of the VQ-GAN.

### 5.4 GENERALIZATION

To assess the generalization of our method, we also evaluate its performance in a new household environment. In more detail, we choose Habitat 2.0 due to its high-fidelity scenes compared with other simulators, such as AI2-THOR. However, Habitat 2.0 does not provide any inter-frame regarding manipulation skills, which is unrealistic. Therefore, we only carry out experiments on navigation tasks.

To enhance usability, we use the pre-trained optical flow model, RAFT (Teed & Deng, 2020), to calculate the optical flow for the previous step since the optical flow cannot be directly obtained. The RAFT results are shown in the last 2 columns of Figure 7. Since VQ-GAN has demonstrated some degree of generalization ability to Habitat 2.0 in Section 5.3, we can predict the motion pattern of the new environment. The remaining task is to transfer the visual style to a new environment, and we adopt LoRA to fine-tune the dynamic model. As shown in Figure 5, we successfully perform style transfer with a small amount of data (tens of samples), and the results with LoRA are closer to real scene images compared to those without LoRA visually.

Table 4 presents the success rate (SR) of LLM-based methods on the HM3D ObjectNav task (Yadav et al., 2023), where our method demonstrates strong generalization with a high SR. Notably, our approach surpasses existing LLM-based methods for the first time, achieving a +4.5% improvement in SR compared to PixelNav (Cai et al., 2023), which navigates to LLM-deduced points. Furthermore, when compared to mapping-based methods employing LLM-guided frontier exploration, our method shows improvements of +6.0% against L3MVN (Yu et al., 2023) and +2.0% against ESC (Zhou et al., 2023).

## 6 CONCLUSION AND LIMITATIONS

This paper introduces EgoPlan, an embodied agent that utilizes an LMM as a one-step planner and a Text2image model as a dynamic model for long-horizon tasks. We demonstrate EgoPlan’s capacity for high-quality image generation, precise optical flow prediction, and promising decision-making. Notably, we have shown its generalization capabilities across diverse environments. It is important to acknowledge a current limitation: EgoPlan employs encapsulated skills as actions, precluding direct low-level control (e.g., joint positions), which remains a subject for future research.

Table 4: We report the zero-shot evaluation results on the HM3D ObjectNav task. Comparison with state-of-the-art methods on the ObjectNav task.

Method	with Mapping	LLM	Extra Sensors	SR	SPL
L3MVN (Yu et al., 2023)	with	GPT-2	Depth, GPS	35.2	16.5
PixelNav (Cai et al., 2023)	without	GPT-4	-	37.9	20.5
ESC (Zhou et al., 2023)	with	GPT-3.5	Depth, GPS	39.2	22.3
Egoplan (Ours)	without	GPT-4	-	<b>41.2</b>	<b>22.5</b>

Table 5: Average endpoint error (AEE) results. The lower the number, the closer the image is to the ground truth.

	Previous flow	Prediction flow
Habitat 2.0	3.30	<b>3.09</b>
AI2-THOR	5.00	<b>4.08</b>
VirtualHome	21.22	<b>15.71</b>

## 7 ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. In this study, no human subjects or animal experimentation was involved. All datasets used, including VH-1.5M Dataset, were sourced in compliance with relevant usage guidelines, ensuring no violation of privacy. We have taken care to avoid any biases or discriminatory outcomes in our research process. No personally identifiable information was used, and no experiments were conducted that could raise privacy or security concerns. We are committed to maintaining transparency and integrity throughout the research process.

## 8 REPRODUCIBILITY STATEMENT

We have made every effort to ensure that the results presented in this paper are reproducible. The experimental setup, including training steps, model configurations, and hardware details, is described in detail in the paper.

Additionally, benchmarks such as Virtualhome and Habitat2.0, are publicly available, ensuring consistent and reproducible evaluation results.

We believe these measures will enable other researchers to reproduce our work and further advance the field.

## REFERENCES

- Gpt-4v(ision) system card. 2023. URL <https://api.semanticscholar.org/CorpusID:263218031>.
- Josh Abramson, Arun Ahuja, Iain Barr, Arthur Brussee, Federico Carnevale, Mary Cassin, Rachita Chhaparia, Stephen Clark, Bogdan Damoc, Andrew Dudzik, et al. Imitating interactive intelligence. *arXiv preprint arXiv:2012.05672*, 2020.
- Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. *arXiv preprint arXiv:2501.03575*, 2025.
- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibor Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Kevin Black, Mitsuhiko Nakamoto, Pranav Atreya, Homer Walke, Chelsea Finn, Aviral Kumar, and Sergey Levine. Zero-shot robotic manipulation with pretrained image-editing diffusion models, 2023.
- Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions, 2023.
- Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024.
- Wenzhe Cai, Siyuan Huang, Guangran Cheng, Yuxing Long, Peng Gao, Changyin Sun, and Hao Dong. Bridging zero-shot object navigation and foundation models through pixel-guided navigation skill. *CoRR*, abs/2309.10309, 2023.
- 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.



- Qifeng Chen and Vladlen Koltun. Full flow: Optical flow estimation by global optimization over regular grids. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4706–4714, 2016. doi: 10.1109/CVPR.2016.509.
- Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample efficiency of grounded language learning. *arXiv preprint arXiv:1810.08272*, 2018.
- Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, pp. 02783649241273668, 2023.
- 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. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- Ishita Dasgupta, Christine Kaeser-Chen, Kenneth Marino, Arun Ahuja, Sheila Babayan, Felix Hill, and Rob Fergus. Collaborating with language models for embodied reasoning. *arXiv preprint arXiv:2302.00763*, 2023.
- 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.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021.
- Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model, 2023.
- Yilun Du, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Joshua B. Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation, 2023a.
- Yilun Du, Mengjiao Yang, Pete Florence, Fei Xia, Ayzaan Wahid, Brian Ichter, Pierre Sermanet, Tianhe Yu, Pieter Abbeel, Joshua B. Tenenbaum, Leslie Kaelbling, Andy Zeng, and Jonathan Tompson. Video language planning, 2023b.
- Alejandro Escontrela, Ademi Adeniji, Wilson Yan, Ajay Jain, Xue Bin Peng, Ken Goldberg, Youngwoon Lee, Danijar Hafner, and Pieter Abbeel. Video prediction models as rewards for reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis, 2021.
- Anthony Brohan et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 2023.
- Embodiment Collaboration et al. Open x-embodiment: Robotic learning datasets and rt-x models, 2024a.
- OpenAI et al. Gpt-4 technical report, 2024b.
- 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. *Advances in Neural Information Processing Systems*, 35: 18343–18362, 2022.
- Hao-Shu Fang, Hongjie Fang, Zhenyu Tang, Jirong Liu, Chenxi Wang, Junbo Wang, Haoyi Zhu, and Cewu Lu. Rh20t: A comprehensive robotic dataset for learning diverse skills in one-shot, 2023.

- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022.
- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all, 2023.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination, 2020.
- Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models, 2022.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models, 2024.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control, 2022.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020.
- Jonathan Ho, Chitwan Saharia, William Chan, David J. Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation, 2021.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models, 2022.
- Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. Vln bert: A recurrent vision-and-language bert for navigation. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pp. 1643–1653, 2021.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp, 2019.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International conference on machine learning*, pp. 9118–9147. PMLR, 2022a.
- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022b.
- 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*, 2(3):6, 2022.
- Siddharth Karamcheti, Megha Srivastava, Percy Liang, and Dorsa Sadigh. Lila: Language-informed latent actions. In *Conference on Robot Learning*, pp. 1379–1390. PMLR, 2022.
- 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*, pp. 14829–14838, 2022.
- Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Po-Chen Ko, Jiayuan Mao, Yilun Du, Shao-Hua Sun, and Joshua B Tenenbaum. Learning to act from actionless videos through dense correspondences. *arXiv preprint arXiv:2310.08576*, 2023.

- Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. Ai2-thor: An interactive 3d environment for visual ai. *arXiv*, 2017.
- Peiyan Li, Hongtao Wu, Yan Huang, Chilam Cheang, Liang Wang, and Tao Kong. Gr-mg: Leveraging partially-annotated data via multi-modal goal-conditioned policy. *IEEE Robotics and Automation Letters*, 2025.
- Shuang Li, Xavier Puig, Chris Paxton, Yilun Du, Clinton Wang, Linxi Fan, Tao Chen, De-An Huang, Ekin Akyürek, Anima Anandkumar, et al. Pre-trained language models for interactive decision-making. *Advances in Neural Information Processing Systems*, 35:31199–31212, 2022.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 9493–9500. IEEE, 2023.
- Jessy Lin, Yuqing Du, Olivia Watkins, Danijar Hafner, Pieter Abbeel, Dan Klein, and Anca Dragan. Learning to model the world with language. *arXiv preprint arXiv:2308.01399*, 2023.
- Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with blockwise ringattention. *arXiv preprint arXiv:2402.08268*, 2024.
- Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao, Jingren Zhou, and Tieniu Tan. Videofusion: Decomposed diffusion models for high-quality video generation, 2023.
- Arjun Majumdar, Ayush Shrivastava, Stefan Lee, Peter Anderson, Devi Parikh, and Dhruv Batra. Improving vision-and-language navigation with image-text pairs from the web. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*, pp. 259–274. Springer, 2020.
- Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and Automation Letters*, 7(3):7327–7334, 2022.
- Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations, 2022.
- 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*, pp. 17359–17371. PMLR, 2022.
- Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. Virtualhome: Simulating household activities via programs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8494–8502, 2018.
- Xavier Puig, Tianmin Shu, Shuang Li, Zilin Wang, Joshua B. Tenenbaum, Sanja Fidler, and Antonio Torralba. Watch-and-help: A challenge for social perception and human-ai collaboration, 2020.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.
- Allen Z Ren, Bharat Govil, Tsung-Yen Yang, Karthik R Narasimhan, and Anirudha Majumdar. Leveraging language for accelerated learning of tool manipulation. In *Conference on Robot Learning*, pp. 1531–1541. PMLR, 2023.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022.

- Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A platform for embodied ai research, 2019.
- Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023. URL <https://arxiv.org/abs/2303.11366>.
- 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, 2020.
- Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation. In *Conference on robot learning*, pp. 894–906. PMLR, 2022.
- Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 11523–11530. IEEE, 2023.
- Richard D Smallwood and Edward J Sondik. The optimal control of partially observable markov processes over a finite horizon. *Operations research*, 21(5):1071–1088, 1973.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022.
- Andrew Szot, Alex Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Chaplot, Oleksandr Maksymets, Aaron Gokaslan, Vladimir Vondrus, Sameer Dharur, Franziska Meier, Wojciech Galuba, Angel Chang, Zsolt Kira, Vladlen Koltun, Jitendra Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0: Training home assistants to rearrange their habitat, 2022.
- Weihao Tan, Wentao Zhang, Shanqi Liu, Longtao Zheng, Xinrun Wang, and Bo An. True knowledge comes from practice: Aligning llms with embodied environments via reinforcement learning. *arXiv preprint arXiv:2401.14151*, 2024.
- Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow, 2020.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- Naoki Wake, Atsushi Kanehira, Kazuhiro Sasabuchi, Jun Takamatsu, and Katsushi Ikeuchi. Gpt-4v(ision) for robotics: Multimodal task planning from human demonstration, 2024.
- 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. *arXiv preprint arXiv:2305.16291*, 2023a.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models, 2023b.
- Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560*, 2023c.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
- Karmesh Yadav, Ram Ramrakhya, Santhosh Kumar Ramakrishnan, Theo Gervet, John Turner, Aaron Gokaslan, Noah Maestre, Angel Xuan Chang, Dhruv Batra, Manolis Savva, et al. Habitat-matterport 3d semantics dataset. In *CVPR*, pp. 4927–4936, 2023.

- Gengshan Yang and Deva Ramanan. Upgrading optical flow to 3d scene flow through optical expansion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1334–1343, 2020.
- Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. Unisim: A neural closed-loop sensor simulator, 2023.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models, 2023.
- Bangguo Yu, Hamidreza Kasaei, and Ming Cao. L3mvn: Leveraging large language models for visual target navigation. *arXiv preprint arXiv:2304.05501*, 2023.
- C. Zach, T. Pock, and H. Bischof. A duality based approach for realtime tv-l1 optical flow. In *Proceedings of the 29th DAGM Conference on Pattern Recognition*, pp. 214–223, Berlin, Heidelberg, 2007. Springer-Verlag. ISBN 9783540749332.
- Andy Zeng, Maria Attarian, Brian Ichter, Krzysztof Choromanski, Adrian Wong, Stefan Welker, Federico Tombari, Aveek Purohit, Michael Ryoo, Vikas Sindhwani, et al. Socratic models: Composing zero-shot multimodal reasoning with language. *arXiv preprint arXiv:2204.00598*, 2022.
- Yuexiang Zhai, Hao Bai, Zipeng Lin, Jiayi Pan, Shengbang Tong, Yifei Zhou, Alane Suhr, Saining Xie, Yann LeCun, Yi Ma, et al. Fine-tuning large vision-language models as decision-making agents via reinforcement learning. *arXiv preprint arXiv:2405.10292*, 2024.
- Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding, 2023a.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023b.
- Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. In *The Twelfth International Conference on Learning Representations*, 2023.
- Kaiwen Zhou, Kaizhi Zheng, Connor Pryor, Yilin Shen, Hongxia Jin, Lise Getoor, and Xin Eric Wang. ESC: exploration with soft commonsense constraints for zero-shot object navigation. volume 202 of *Proceedings of Machine Learning Research*, pp. 42829–42842. PMLR, 2023.
- Siyuan Zhou, Yilun Du, Jiaben Chen, Yandong Li, Dit-Yan Yeung, and Chuang Gan. Robodreamer: Learning compositional world models for robot imagination, 2024.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models, 2023.

## APPENDIX

### A LLM USAGE

Large Language Models (LLMs) were used to aid in the writing and polishing of the manuscript. Specifically, we used an LLM to assist in refining the language, improving readability, and ensuring clarity in various sections of the paper. The model helped with tasks such as sentence rephrasing, grammar checking, and enhancing the overall flow of the text.

It is important to note that the LLM was not involved in the ideation, research methodology, or experimental design. All research concepts, ideas, and analyses were developed and conducted by the authors. The contributions of the LLM were solely focused on improving the linguistic quality of the paper, with no involvement in the scientific content or data analysis.

The authors take full responsibility for the content of the manuscript, including any text generated or polished by the LLM. We have ensured that the LLM-generated text adheres to ethical guidelines and does not contribute to plagiarism or scientific misconduct.

### B RELATED WORK

#### B.1 DIFFUSION MODEL

The diffusion model (Ho et al., 2020; Song et al., 2022) has been extensively studied in the field of image generation (Dhariwal & Nichol, 2021; Ho et al., 2021; Rombach et al., 2022) and image editing (Gal et al., 2022; Hertz et al., 2022; Meng et al., 2022). Diffusion models can achieve a high degree of control during the image generation. In more detail, InstructPix2Pix (InstructP2P) (Brooks et al., 2023) trains a conditional diffusion model that, given an input image and text instruction for how to edit it, generates the edited image. ControlNet (Zhang et al., 2023b) is widely used to control the style of the generated image by using various forms of prior information, such as edge information and segmentation. By adding LoRA or adapter (Houlsby et al., 2019) modules to the network, the model trained on one data distribution can also be transferred to other data distributions (different visual styles) through a few picture examples. The images produced by current diffusion models are of very high quality, highly realistic, and easily controllable. It prompts various fields to consider using these generated images to assist in accomplishing other tasks. Our paper adopts the diffusion model to generate task subgoals and predict the image of the next state for decision-making.

#### B.2 DYNAMIC MODEL AND WORLD MODEL FOR DECISION-MAKING

In the works of using world model for long-range mission planning, the Dreamer series (Hafner et al., 2020; 2022; 2024) models environmental dynamics in latent space to predict future states within gaming contexts, enabling agents to learn tasks through imagination and reducing the number of interactions needed for effective learning. However, as these world models are developed in latent space rather than pixel space, they often struggle to generalize to unseen tasks and environments. A world model constructed in pixel space may offer improved generalization capabilities. Recent studies have sought to address how to learn world models from large-scale video datasets (Liu et al., 2024). In Genie (Bruce et al., 2024), researchers utilize a latent action representation, though their focus primarily revolves around 2D platform video games or simple robotic actions. By meticulously orchestrating rich data across various dimensions, UniSim (Yang et al., 2023) simulates realistic visual experiences in response to actions performed by humans, robots, and other interactive agents. Overall, the applications of world models extend beyond gaming and robotics. For instance, in Escontrela et al. (2024), frame-by-frame video prediction is employed as a mechanism for providing rewards in reinforcement learning. DynaLang (Lin et al., 2023) explores the integration of language prediction as an element of the world model, enabling the training of multimodal world models using datasets that lack explicit actions or rewards. In DynaLang, the representation is shared between vision and language within the world model.

### B.3 EMBODIED AGENT WITH LMMs

The successful integration of language as a semantically rich input for interactive decision-making underscores the pivotal role of LMMs in facilitating interaction and decision-making processes (Abramson et al., 2020; Karamcheti et al., 2022; Li et al., 2022). LMMs have also been employed across various environments to support robot navigation (Parisi et al., 2022; Hong et al., 2021; Majumdar et al., 2020) and manipulation tasks (Jiang et al., 2022; Ren et al., 2023; Karamcheti et al., 2022). Recently, numerous approaches have emerged that leverage LMMs to enhance the planning and reasoning capabilities of embodied agents. For instance, SayCan (Ahn et al., 2022) evaluates the affordance of potential actions by combining their probabilities derived from LMMs with a value function. (Zeng et al., 2022) integrate a language and multimodal model (LMM) with a visual-language model and a pre-trained language-conditioned policy (Shridhar et al., 2022) to facilitate open vocabulary robotic tasks. Similarly, Huang et al. (2022a) illustrate that LMMs can be effectively utilized for planning and executing simple household tasks, grounding LMM-generated actions by comparing their embeddings with a predefined list of acceptable actions. To incorporate environmental feedback, Inner Monologue (Huang et al., 2022b) enhances SayCan through a closed-loop principle. This principle is further employed in related works such as (Yao et al., 2023; Huang et al., 2022b; Kim et al., 2024; Singh et al., 2023; Liang et al., 2023; Shinn et al., 2023; Wang et al., 2023c) to continuously monitor agent behaviors and refine plans accordingly for tasks in domains like computer automation and Minecraft. Furthermore, there are methods that prompt language and multimodal models (LMMs) to generate temporally abstracted actions (Zheng et al., 2023). Dasgupta et al. (2023) utilize the LMM as both a planner and a success detector for an agent, with their actor module requiring pre-training using reinforcement learning to enable the agent to adhere to natural language instructions. While these studies yield impressive results, they are heavily dependent on the inherent capabilities of powerful LMMs, such as GPT-4 and PaLM (Chowdhery et al., 2023), which presents challenges when attempting to apply these approaches to smaller LMMs with limited reasoning abilities, such as LLaMA-7B. GLAM (Carta et al., 2023) employs RL fine-tuning to achieve functional grounding of LLMs and LMMs. However, their focus is primarily on simple primitive actions (e.g., turn left, turn right, go forward) evaluated within toy environments, such as BabyAI (Chevalier-Boisvert et al., 2018), using a significantly smaller encoder-decoder LMM, Flan-T5-780M. These primitive actions possess a similar token count and lack substantial semantic meaning, which leads to an underutilization of LMM capabilities. Consequently, they fail to adequately explore the effects of prompt design and address the imbalance within the action space, resulting in additional instability and reduced robustness.

## C DETAILS OF VIRTUALHOME TASKS

We conducted experiments to evaluate the decision-making ability of all methods in the VirtualHome environment. In total, we investigated 12 complex tasks, with detailed instructions and reference subtasks steps for each task as follows:

Listing 1: Instructions and subtasks.

```
<$one-house instructions$>

1. take and place: take the bread from the toaster and place it on the
   plate on the table
steps: (a). walk to the toaster
       (b). grab the bread
       (c). walk to the table
       (d). place the bread on the plate
2. take and put1: take the apple from the table and put it in the
   microwave
steps: (a). walk to the table
       (b). grab the apple
       (c). walk to the microwave
       (d). open the microwave (if the microwave is closed)
       (e). put the apple in the microwave
3. take and put2: take the book from the table and put it on the
   bookshelf
steps: (a). walk to the table
```

```

918         (b). take the book
919         (c). grab the book
920         (d). walk to the bookshelf
921         (e). put the book on the bookshelf
922 4. take and drink: take the water glass from the table and drink from it
923 steps: (a). walk to the table
924         (b). take the water glass
925         (c). drink the water glass
926 5. turn on sit: turn on the TV and sit down
927 steps: (a). walk to the TV
928         (b). turn on the TV
929         (c). walk to the chair
930         (d). sit down
931 6. put apple: Put an apple that is on the table into the bookshelf
932 steps: (a). walk to the table
933         (b). grab the apple
934         (c). walk to the bookshelf
935         (d). put the apple on the bookshelf
936
937 <$two-houses instructions$>
938
939 7. take and place2: take the frying pan from the counter and place it in
940 the sink
941 steps: (a). walk to the counter
942         (b). grab the frying pan
943         (c). walk through the door
944         (d). walk to the sink
945         (e). place frying pan in the sink
946 8. take and place3: take the condiment shaker from the bookshelf and
947 place it on the table
948 steps: (a). walk to the bookshelf
949         (b). grab the condiment shaker
950         (c). walk through the door
951         (d). walk to the table
952         (e). place condiment shaker on the table
953 9. take and put3: take the salmon on top of the microwave and put it in
954 the fridge
955 steps: (a). walk to the microwave
956         (b). grab the salmon
957         (c). walk through the door
958         (d). walk to the fridge
959         (e). open the fridge (if the fridge is closed)
960         (f). put salmon in the fridge
961 10. take open and put: take the pie on the table and warm it using the
962 stove
963 steps: (a). walk to the table
964         (b). grab the pie
965         (c). walk through the door
966         (d). walk to the stove
967         (e). put pie on the stove
968         (f). switch on the stove
969 11. take put and open: put the sponge in the sink and wet it by switching
970 on the faucet
971 steps: (a). walk to the sponge
972         (b). grab the sponge
973         (c). walk through the door
974         (d). walk to the sink
975         (e). put sponge in the sink
976         (f). switching on the faucet
977 12. take and put4: take the condiment bottle from the kitchen table and
978 put it on the plate
979 steps: (a). walk through the door
980         (b). walk to the kitchen table
981         (c). grab the condiment bottle
982         (d). walk to the plate

```



```
(e). put pie on the stove
(f). switch on the stove
```

In terms of the average task completion length, since we want to prove the effectiveness of our pipeline in long-term planning, we draw on the metric of calvin benchmark (Mees et al., 2022), where the next subtask is executed after the completion of the previous subtask, that is, the completion of the next subtask is conditional on the completion of the previous subtask. This index represents the average number of subtasks that each pipeline can complete after 100 repeated experiments of each task, which measures the long-term planning ability of the pipeline. One repeated experiment represents an **episode**. Since the virtualhome emulator can return instructions on whether the task was successfully executed, our evaluation is automated to calculate the success rate.

## D DETAILS OF VH-1.5M’S TEXT ACTIONS

We automatically collected the dataset in the order of action category to object. Firstly, 50 different indoor environments are randomly initialized as House1-50, and then the action types (such as put and walk to) are specified. Under each action type, items in the house are randomly selected as the imposed objects of the action. Such commands are executed in the VirtualHome simulator to form a trajectory in the dataset.

The dataset includes a wide range of action sequences, each meticulously annotated with corresponding text actions. These text actions are crucial for providing contextual information that aligns visual actions with natural language descriptions. Below, we detail the process and structure used to generate the text actions for each action sequence in the dataset.

The generation of text actions for VH-1.5M involves a systematic and automated process. This process ensures consistency and variety in the text actions, which are essential for robust training and evaluation in vision-and-language tasks. The key steps in this process are as follows:

**Verb Selection:** A list of verbs related to various actions (e.g., "walk through," "close," "drink") is predefined. For each identified action sequence directory, a verb is randomly selected from the relevant list. This selection ensures a diverse representation of actions.

**Object Name Extraction:** Each directory represents the object acted upon, which signifies the object affected by the action. However, if the action does not involve an object, such as "walk through" or "turn left," no extraction is necessary.

**Phrase Construction:** Two types of phrases are constructed for each action sequence:

**Next Timestep Phrase:** Describes the immediate next action in the sequence. For example, "next timestep: redeposit the plate".

**Goal State Phrase:** Describes the intended final action or goal of the sequence. For example, "the goal state: redeposit plate".

**Prompt File Creation:** The constructed phrases are saved in a prompt json file within the respective action sequence directory. This JSON file contains two keys: "next" and "goal," corresponding to the next timestep phrase and goal state phrase, respectively.

### D.1 MORE EXAMPLES OF THE SAMPLES

We give some samples in the sequence of the task, which are shown in Figure 6. Note that samples in one sequence are arranged in chronological order, with the timestep increasing from top to bottom.

## E MORE EXAMPLES OF GENERATING IMAGES

More examples of generated images from EgoPlan can be seen in Figure 8. Each line represents a task, and the task prompts are, in order: "capture the chicken", "grasp juice", "grasp the hairproduct", "open the cabinet", "open the microwave", "go left", "make a left", "make a left-hand turn", "make a right", "turn right", "turn to the right", "walk straight ahead".

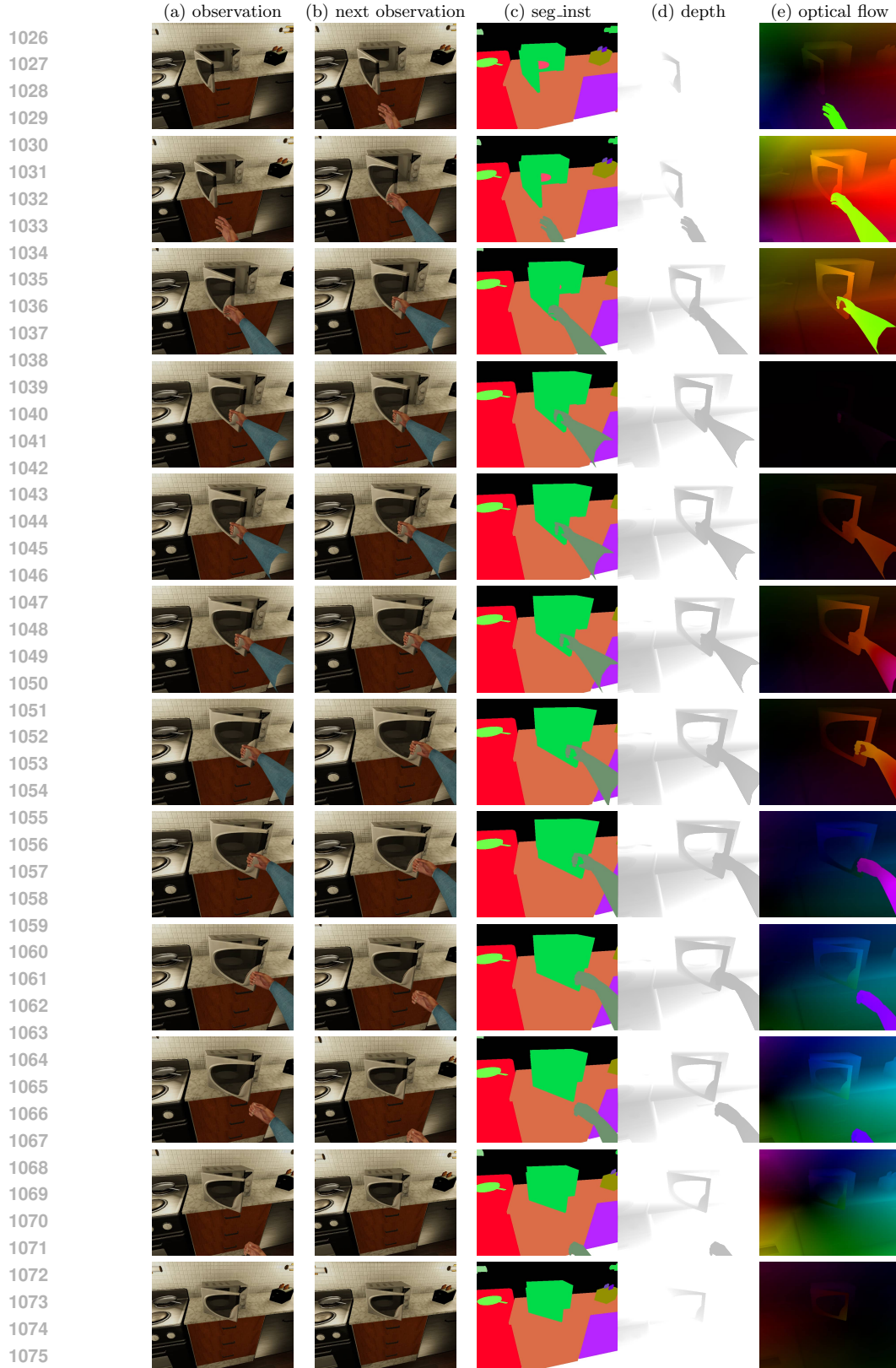


Figure 6: Samples in the sequence of closing the microwave.

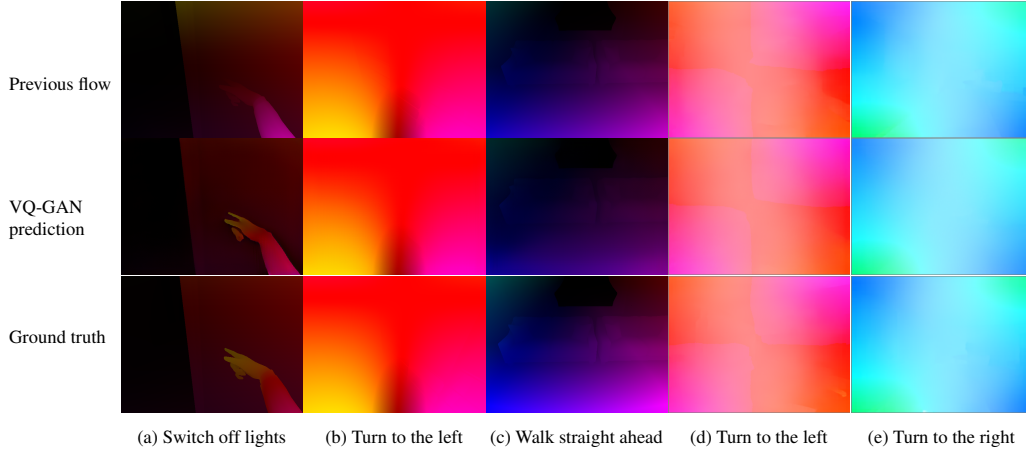


Figure 7: Examples of optical flow prediction by VQ-GAN. The first 3 columns are optical flow from the VirtualHome environment. The last 2 columns are optical flow from the Habitat2.0 environment.

As illustrated in Figure 7a and 7c, the quality of optical flow prediction for details is promising. Furthermore, as demonstrated in Figure 7d and 7e, the VQ-GAN trained on the VH-1.5M dataset can easily generalize to other environments.

Table 6: User study for the subgoal generation. The user score is the percentage of images that users consider to meet the criteria out of the total 1000 images.

	Close	Drink	Grab	Open	Put back	Put in
Mean user score(%)	66.5	71.75	55	66.375	62.125	64.625
	Sit	Stand up	Switch off	Switch on	Walk through	
Mean user score(%)	79.875	78.75	73.375	77.875	79	

## F USER STUDY OF SUBGOAL IMAGE GENERATION

We also conduct a user study on the image generation of the subgoal. A total of 8 users evaluated whether the generated image met the criteria of the subgoal described in the text. Each user evaluates 100 generated images for each action, and the evaluation results are shown in Table 6. The results show that most of the generated subgoal images can represent the meaning of the text subgoals. More examples of generating figures can be seen in Figure 9. It is worth mentioning that after our dataset and training, the subgoal prediction model exhibits certain scene understanding ability. For example, the "power up the lightswitch" subgoal illuminates objects (like walls) in a room in the view scene, which is interesting compared to previous work where image generation of subgoals helps decision making.

## G PROMPT OF TASK-DECOMPOSITION AND LOW-LEVEL ACTION SELECTION

We conducted experiments with detailed query prompt for each task as follows:

Listing 2: query for action selection.

```
Start working. The picture of what you can see has been given above, the
picture is what you see from the first person perspective as the
person in the room. Analyze the scene and all the items in the
picture to make a task plan to complete the instruction.
The instruction is as follows:
"""
{"instruction": [INSTRUCTION]}
"""
The history is as follows:
```



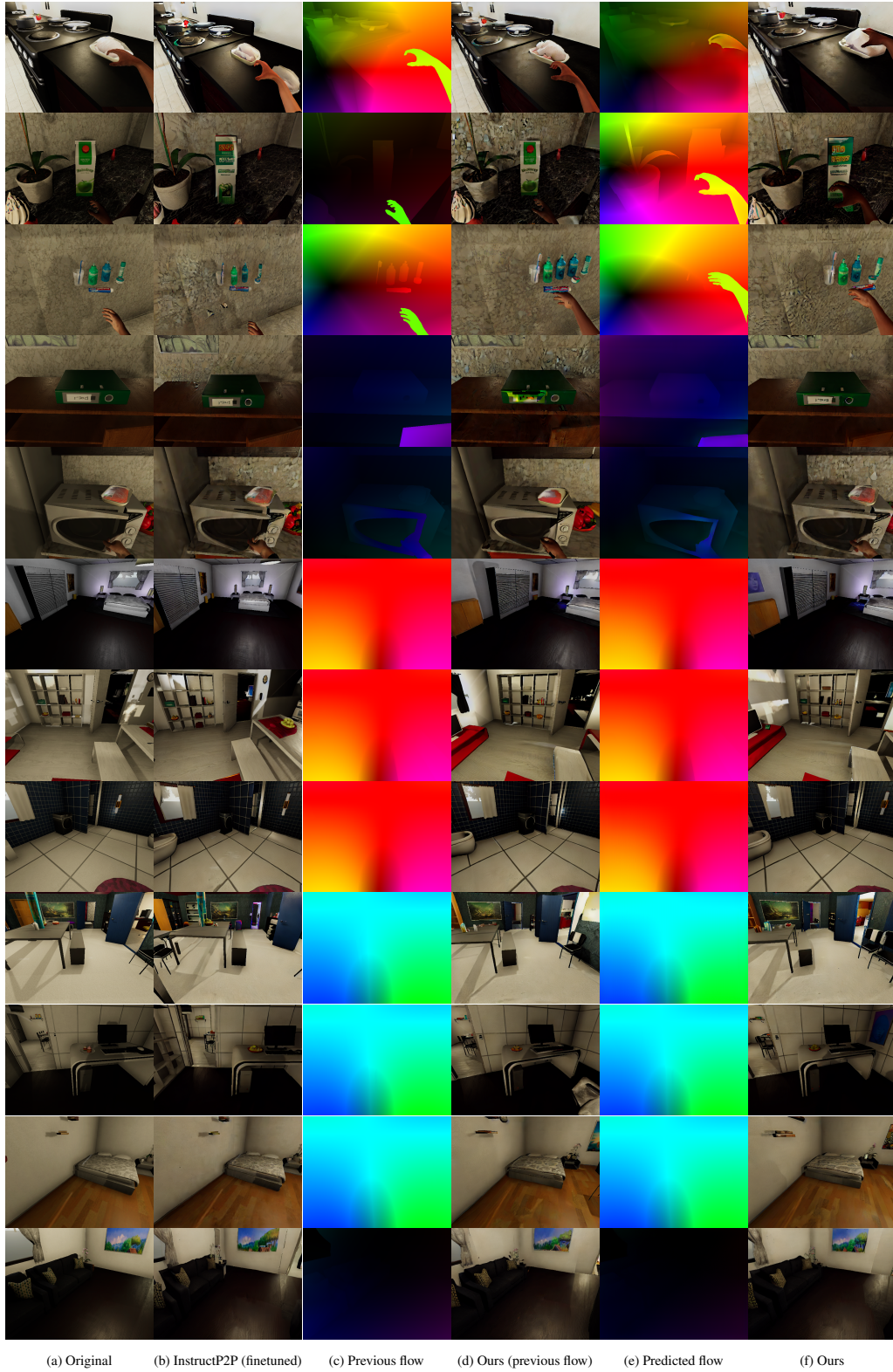


Figure 8: Examples of the generated image of the EgoPlan in VirtualHome. We can find that in some hand reconstruction and direction understanding scenes, the model without introducing optical flow prior information often performs poorly.

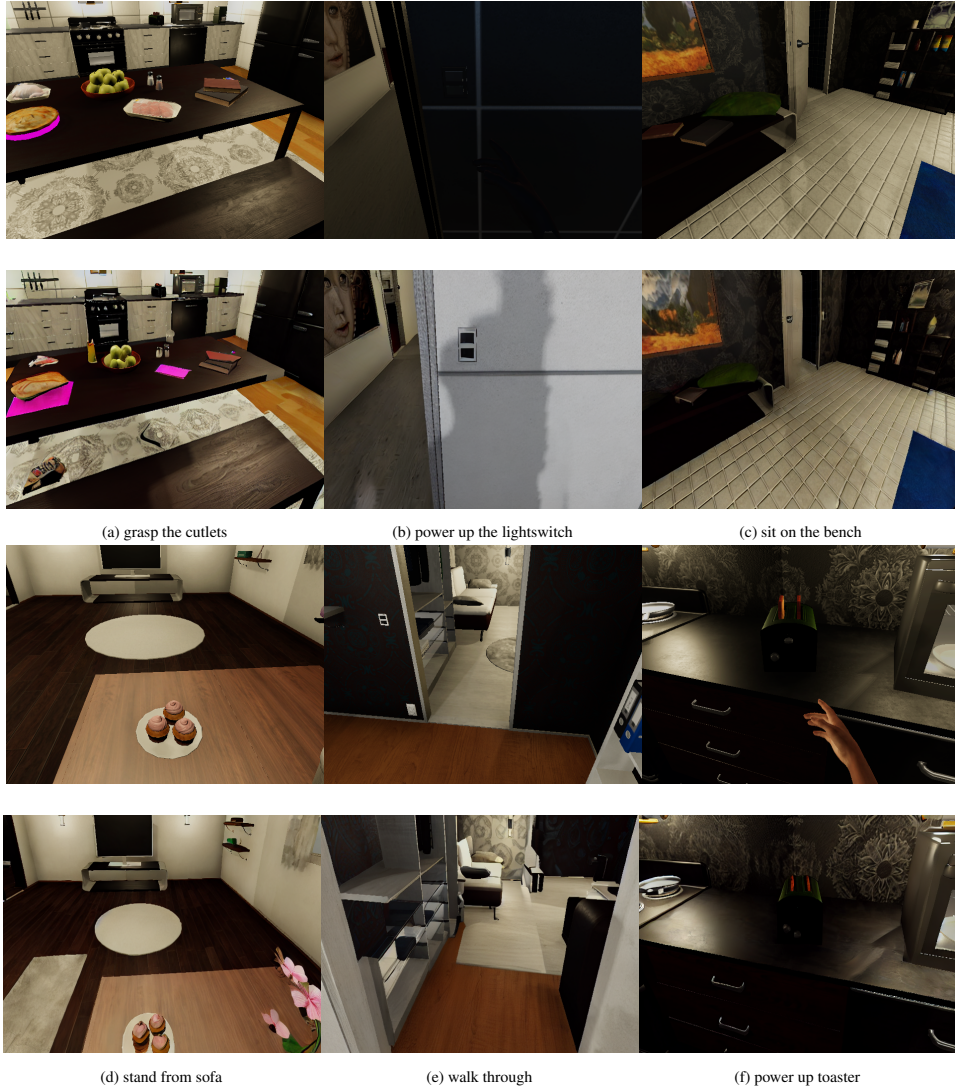


Figure 9: Examples of the generated image subgoals. The first and third rows is the original image, and the second and forth rows is the image subgoal generated based on the text subgoal.

```

"""
{"history": [HISTORY]}
"""
You return should follow these rules:
1. Make sure you provide 4 lines of output each time, the first line is
the ["Preoperation"] and the secondline is the ["Postoperation"] of
the action to be taken in the current task plan, and the third line
is the action to be taken in the plan, which is the ["task_sequence
"]. The fourth line is the natural language expression of the action
taken, namely ["step_instructions"]. When output the answer, do not
attach "step_instructions", "task_sequence", etc.
2. In addition to these, other problem such as input images is too dark
and historical actions is empty, please DO NOT output.
3. Make sure that element of the ["step_instructions"] explains
corresponding element of the ["task_sequence"]. That is, the fourth
line explains the third line.
4. DO NOT USE undefined verbs. USE ONLY verbs in "HUMAN ACTION LIST".
5. The first line and the second line are detailed explanation of the
forth line. For the task in the forth line, it must be explained in
two parts: ["Preoperation"] and ["Postoperation"] in the first and

```

second line, separately represents the action state of the agent and item before and after the execution of the task.

6. Look carefully at the output examples provided. DO NOT use any strings or spaces at the end of sentences. Never left ',' at the end of the sentences. STRICTLY ENSURE that the output is always four lines long, with no blank lines.
7. The environment given is a picture that you see from the first person perspective as the person in the room. Analyze the scene and all the items in the picture to make a task plan. If you see a picture that is all black, this means there has been no task planning or execution before, please give a general task plan, but BE SURE to stick to the output format shown earlier.
8. When selecting each action for task planning, carefully think about the function of the action in terms of the two parts ["Preconditions"] and ["Postconditions"] after the action, where ["Preconditions"] represents the state of the environment before the action is executed, and ["Postconditions"] represents the state of the environment after the execution, after which the planning is carried out.
9. All sentences you output should NOT be double-quoted.
10. Please strictly correspond to the actions and items in the instructions, please strictly keep the spelling of the items, for multi-word items, please do not add connection symbols between words, for items composed of single-word, please do not split the word.
11. The history is a string that records the actions performed in the past few steps, separated by " ". Please plan what action to perform at this step based on the historical actions, instructions and the current picture.
12. Make sure that you output a consistent manipulation as a human. For example, grasping an object should not occur in successive steps. Consider whether the current action is similar to the last action in the history. DO NOT output same two actions in row.
13. Every time you do task planning, you should consider whether the historical action in history and the current action have completed the instruction, and if so, output "Stop()" in time.

Adhere to the output format I defined above. Follow the nine rules. Think step by step.

We conducted experiments with detailed environment, role of LMM, action function, few-shot output example prompt for each task as follows:

Listing 3: prompt for environment.

```
[user]
Information about environments and objects are given as a picture that
can be seen from the first person perspective. The picture will be
given in the example latter.
-----
The texts above are part of the overall instruction. Do not start working
yet:
[assistant]
Understood. I will wait for further instructions before starting to work.
```

Listing 4: prompt for role of LMM.

```
[user]
You are an excellent interpreter of human instructions for household
tasks. Given an instruction and information about the working
environment, you break it down into a sequence of human actions.
Please do not begin working until I say "Start working." Instead, simply
output the message "Waiting for next input." Understood?
[assistant]
Waiting for next input.
```

Listing 5: prompt for explanation of action function.

```

1296
1297 [user]
1298 Necessary and sufficient human actions are defined as follows:
1299 """
1300 "HUMAN ACTION LIST"
1301
1302 Walk(arg1): Walks some distance towards a room or object.
1303 Preconditions: If the environment represented by picture doesn't have the
1304 obj1 for the task decomposition you did to perform the action, add a
1305 subtask of Walk(obj1) before the task.
1306
1307 Grab(arg1): Grabs an object.
1308 Preconditions: The object1 property is grabbable (except water). The
1309 character is close to obj1. obj1 is reachable (not inside a closed
1310 container). The character has at least one free hand.
1311 Postconditions: Adds a directed edge: character holds_rh or hold_lh, obj1
1312 . obj1 is no longer on a surface or inside a container.
1313
1314 Open(arg1): Opens an object.
1315 Preconditions: The obj1 property is IS_OPENABLE and the state is closed.
1316 The character is close to obj1. obj1 is reachable (not inside a
1317 closed container). The character has at least one free hand.
1318 Postconditions: The obj1 state is open.
1319
1320 Close(arg1): Closes an object.
1321 Preconditions: The obj1 property is IS_OPENABLE and the state is open.
1322 The character is close to obj1. obj1 is reachable (not inside a
1323 closed container). The character has at least one free hand.
1324 Postconditions: The obj1 state is closed.
1325
1326 Put(arg1, arg2): Puts an object on another object.
1327 Preconditions: The character holds_lh obj1 or character holds_rh obj1.
1328 The character is close to obj2.
1329 Postconditions: Removes directed edges: character holds_lh obj1 or
1330 character holds_rh obj1. Adds directed edges: obj1 on obj2.
1331
1332 PutIn(arg1, arg2): Puts an object inside another object that is OPENABLE,
1333 such as stove and microwave.
1334 Preconditions: The character holds_lh obj1 or character holds_rh obj1.
1335 The character is close to obj2. obj2 is not closed. If obj2 is closed
1336 , The character should open obj2 first and put obj1 in obj2.
1337 Postconditions: Removes directed edges: character holds_lh obj1 or
1338 character holds_rh obj1. Adds directed edges: obj1 inside obj2.
1339
1340 SwitchOn(arg1): Turns an object on.
1341 Preconditions: The obj1 has the property "switch." The obj1 state is off.
1342 The character is close to obj1.
1343 Postconditions: The obj1 state is on.
1344
1345 SwitchOff(arg1): Turns an object off.
1346 Preconditions: The obj1 has the property "switch." The obj1 state is on.
1347 The character is close to obj1.
1348 Postconditions: The obj1 state is off.
1349
1350 Drink(arg1): Drinks from an object.
1351 Preconditions: The obj1 property is drinkable or recipient. The character
1352 is close to obj1.
1353
1354 Sit(arg1): Sit down on an object.
1355 Preconditions: The obj1 property is sittable. The character is close to
1356 obj1.
1357
1358 Stop(): The instruction can end the task sequence after the completion of
1359 the task by the planned instruction.

```

```

Preconditions: After the instruction is decomposed into a series of tasks
, these tasks fulfill all the requirements of the instruction to be
executed in order, that is, the instruction is completed in the
history.
"""
-----
The texts above are part of the overall instruction. Do not start working
yet:
[assistant]
Waiting for next input.

```

Listing 6: prompt for output example.

```

[user]
I will give you some examples of the input and the output you will
generate.
Example 1:
"""
- Input:
The picture of what you can see has been given above.
"instruction": "take the salmon on top of the microwave and put it in the
fridge"
"history": ""
- Output:
The microwave where the salmon is located appears to be distant or out of
reach, and I need to approach it to interact with it.
I am now close enough to the microwave to interact with it, specifically
to reach the salmon.
Walk(<microwave>)
Walk towards the microwave to reach the salmon on top.
"""
-----
Example 2:
"""
- Input:
The picture of what you can see has been given above.
"instruction": "take the salmon on top of the microwave and put it in the
fridge"
"history": "Walk(<microwave>)"
- Output:
The salmon is on top of the microwave and within reach. I have at least
one free hand to grab it.
I am now holding the salmon, which is no longer on the microwave.
Grab(<salmon>)
Grab the salmon from the top of the microwave
"""
-----
Example 3:
"""
- Input:
The picture of what you can see has been given above.
"instruction": "take the salmon on top of the microwave and put it in the
fridge"
"history": "Walk(<microwave>)""Grab(<salmon>)"
- Output:
The fridge appears to be distant or out of reach, and I need to approach
it to interact with it.
I am now close enough to the fridge to put the salmon inside.
Walk(<fridge>)
Walk to the fridge with the salmon
"""
-----
Example 4:
"""
- Input:

```



```

1404 The picture of what you can see has been given above.
1405 "instruction": "take the salmon on top of the microwave and put it in the
1406 fridge"
1407 "history": "Walk(<microwave>)" "Grab(<salmon>)" "Walk(<fridge>)"
1408 - Output:
1409 Before I can put the salmon inside, the fridge must be open.
1410 The fridge is now open, and I can place items inside.
1411 Open(<fridge>)
1412 Open the fridge
1413 ""
1414 -----
1415 Example 5:
1416 ""
1417 - Input:
1418 The picture of what you can see has been given above.
1419 "instruction": "take the salmon on top of the microwave and put it in the
1420 fridge"
1421 "history": "Walk(<microwave>)" "Grab(<salmon>)" "Walk(<fridge>)" "Open(<
1422 fridge>)"
1423 - Output:
1424 I hold the salmon. I am close to the fridge which is now open.
1425 The salmon is now inside the fridge, and my hands are free.
1426 PutIn(<salmon>, <fridge>)
1427 Put the salmon in the fridge
1428 ""
1429 -----
1430 Example 6:
1431 ""
1432 - Input:
1433 The picture of what you can see has been given above.
1434 "instruction": "take the salmon on top of the microwave and put it in the
1435 fridge"
1436 "history": "Walk(<microwave>)" "Grab(<salmon>)" "Walk(<fridge>)" "Open(<
1437 fridge>)" "PutIn(<salmon>, <fridge>)"
1438 - Output:
1439 After placing the salmon inside, the fridge remains open.
1440 The fridge is now closed, securing the salmon inside.
1441 Close(<fridge>)
1442 Close the fridge door
1443 ""
1444 -----
1445 Example 7:
1446 ""
1447 - Input:
1448 The picture of what you can see has been given above.
1449 "instruction": "take the salmon on top of the microwave and put it in the
1450 fridge"
1451 "history": "Grab(<salmon>)" "Walk(<fridge>)" "Open(<fridge>)" "PutIn(<salmon
1452 >, <fridge>)" "Close(<fridge>)"
1453 - Output:
1454 I take the salmon on top of the microwave and put it in the fridge.
1455 The instruction has been finished.
1456 Stop()
1457 Complete the instruction and stop the task planning
1458 ""
1459 -----
1460 The texts above are part of the overall instruction. Do not start working
1461 yet:
1462 [assistant]
1463 Waiting for next input.

```

Listing 7: prompt for output format.

```
[user]
```

```

1458 You divide the actions given in the text into detailed robot actions and
1459 put them together as a python dictionary.
1460 The dictionary has three keys.
1461 """
1462 - dictionary["task_cohesion"]: A dictionary containing information about
1463 the robot's actions that have been split up.
1464 - dictionary["instruction_summary"]: contains a brief summary of the
1465 given sentence.
1466 """
1467 Two keys exist in dictionary["task_cohesion"].
1468 """
1469 - dictionary["task_cohesion"]["task_sequence"]: A dictionary containing
1470 information about the human's actions that have been split up.
1471 - dictionary["task_cohesion"]["step_instructions"]: contains a brief text
1472 explaining why this step is necessary.
1473 -----
1474 The texts above are part of the overall instruction. Do not start working
1475 yet:
1476 [assistant]
1477 Waiting for next input.

```

## H TRAJECTORIES OF SELF-REFLECTION IN NAVIGATION TASKS

When executing navigation tasks, the subgoal is "walk to (<somewhere>)" while the underlying actions include "walk forward", "turn left", and "turn right". Accomplishing the navigation task with such a subgoal constitutes a long-horizon composite task. In addressing these types of tasks, we employ a **React+Reflexion** mechanism that leverages previous actions to perform operations such as obstacle avoidance and target searching. For example, we can illustrate trajectories based on these actions. The trajectories under the self-reflection mechanism are shown as follows:

Listing 8: Trajectory of self-reflection.

```

1487 Trial #1
1488 Environment: <observation image>. You are in the middle of a room.
1489 Looking quickly around you, you see a bed 1, a desk 2, a desk 1, a
1490 computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer
1491 1, a garbagecan 1, a laundryhamper 1, a cabinet 1, a shelf 6, a shelf
1492 5, a shelf 4, a shelf 3, a shelf 2, and a shelf 1.
1493 Task: walk to computer 1.
1494 > think: To solve the task, I need to find a desk, then find a computer.
1495 > think: First I need to find a computer. From the obsersation image, a
1496 computer is on desk 1, from the dynamic model output, you should turn
1497 left to get direction.
1498 > turn left
1499 > success
1500 Environment: <observation image>. You are in the middle of a room.
1501 Looking quickly around you, you see a bed 1, a desk
1502 2, a desk 1, a computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer
1503 2, a drawer 1, a garbagecan 1, a
1504 laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a shelf 3,
1505 a shelf 2, and a shelf 1.
1506 > think: First I need to find a computer. From the obsersation image, a
1507 computer is on desk 1, from the dynamic model output, you should turn
1508 left to get direction.
1509 > turn left
1510 > success
1511 Environment: <observation image>. You are in the middle of a room.
1512 Looking quickly around you, you see a bed 1, a desk
1513 2, a desk 1, a computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer
1514 2, a drawer 1, a garbagecan 1, a
1515 laundryhamper 1, a cabinet 1, a shelf 5, a shelf 4, a shelf 3, a shelf 2,
1516 and a shelf 1.

```

```

1512 > think: You are get right direction of the desk, the desk is in the
1513 middle of the egocentric view.
1514 > walk forward
1515 > Nothing happens.
1516 Environment: <observation image>. You are in the middle of a room.
1517 Looking quickly around you, you see a bed 1, a desk
1518 2, a desk 1, a computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer
1519 2, a drawer 1, a garbagecan 1, a
1520 laundryhamper 1, a cabinet 1, a shelf 5, a shelf 4, a shelf 3, a shelf 2,
1521 and a shelf 1.
1522 > think: You are get right direction of the desk, the desk is in the
1523 middle of the egocentric view.
1524 > walk forward
1525 > Nothing happens.
1526 Environment: <observation image>. You are in the middle of a room.
1527 Looking quickly around you, you see a bed 1, a desk
1528 2, a desk 1, a computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer
1529 2, a drawer 1, a garbagecan 1, a
1530 laundryhamper 1, a cabinet 1, a shelf 5, a shelf 4, a shelf 3, a shelf 2,
1531 and a shelf 1.
1532 > think: You are get right direction of the desk, the desk is in the
1533 middle of the egocentric view.
1534 > walk forward
1535 > Nothing happens.
1536 Status: Fail
1537
1538 Reflection: When we locate the computer and identify the correct
1539 direction, we find that the action to approach the computer is not
1540 responding. Observing first-person perspective images reveals that
1541 there is a bed in front, acting as an obstacle that prevents forward
1542 movement. To get closer to the computer, it is necessary to navigate
1543 around the bed first.
1544
1545 Trial #2
1546 Environment: <observation image>. You are in the middle of a room.
1547 Looking quickly around you, you see a bed 1, a desk 2, a desk 1, a
1548 computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer
1549 1, a garbagecan 1, a laundryhamper 1, a cabinet 1, a shelf 6, a shelf
1550 5, a shelf 4, a shelf 3, a shelf 2, and a shelf 1.
1551 Task: walk to computer 1.
1552 > think: To solve the task, I need to find a desk, then find a computer.
1553 > think: First I need to find a computer. From the obsersation image, a
1554 computer is on desk 1, from the dynamic model output, you should turn
1555 right to get around the bed.
1556 > turn right
1557 > success
1558 Environment: <observation image>. You are in the middle of a room.
1559 Looking quickly around you, you see a bed 1, a desk 2, a desk 1, a
1560 computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer
1561 1, a laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a
1562 shelf 3, a shelf 2, and a shelf 1, a laundryhamper 2.
1563 > think: First I need to find a computer. From the obsersation image, a
1564 computer is on desk 1, from the dynamic model output, you should turn
1565 right to get around the bed.
1566 > turn right

```

```

1566 > success
1567 Environment: <observation image>. You are in the middle of a room.
1568 Looking quickly around you, you see a bed 1, a desk 2, a desk 1, a
1569 computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer
1570 1, a laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a
1571 shelf 3, a shelf 2, and a shelf 1, a laundryhamper 2.
1572 > think: First I need to find a computer. From the observation image, a
1573 computer is on desk 1, currently, there are no obstacles blocking
1574 your path in front of you. From the dynamic model output, you should
1575 walk forward to get around the bed.
1576 > walk forward
1577 > success
1578 Environment: <observation image>. You are in the middle of a room.
1579 Looking quickly around you, you see a bed 1, a desk 2, a desk 1, a
1580 computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer
1581 1, a laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a
1582 shelf 3, a shelf 2, and a shelf 1, a laundryhamper 2.
1583 > think: First I need to find a computer. From the observation image, a
1584 computer is on desk 1, currently, there are no obstacles blocking
1585 your path in front of you. From the dynamic model output, you should
1586 walk forward to get around the bed.
1587 > walk forward
1588 > success
1589 Environment: <observation image>. You are in the middle of a room.
1590 Looking quickly around you, you see a bed 1, a desk 2, a desk 1, a
1591 computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer
1592 1, a laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a
1593 shelf 3, a shelf 2, and a shelf 1, a laundryhamper 2.
1594 > think: First I need to find a computer. From the observation image, a
1595 computer is on desk 1, currently, there are no obstacles blocking
1596 your path in front of you. From the dynamic model output, you should
1597 walk forward to get around the bed.
1598 > walk forward
1599 > success
1600 Environment: <observation image>. You are in the middle of a room.
1601 Looking quickly around you, you see a bed 1, a desk
1602 2, a desk 1, a computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer
1603 2, a drawer 1, a garbagecan 1, a
1604 laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a shelf 3,
1605 a shelf 2, and a shelf 1.
1606 > think: First I need to find a computer. From the observation image, a
1607 computer is on desk 1, from the dynamic model output, you should turn
1608 left to get direction.
1609 > turn left
1610 > success
1611 Environment: <observation image>. You are in the middle of a room.
1612 Looking quickly around you, you see a bed 1, a desk
1613 2, a desk 1, a computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer
1614 2, a drawer 1, a garbagecan 1, a
1615 laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a shelf 3,
1616 a shelf 2, and a shelf 1.
1617 > think: First I need to find a computer. From the observation image, a
1618 computer is on desk 1, from the dynamic model output, you should turn
1619 left to get direction.

```

```

1620 > turn left
1621 > success
1622 Environment: <observation image>. You are in the middle of a room.
1623     Looking quickly around you, you see a bed 1, a desk
1624     2, a desk 1, a computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer
1625     2, a drawer 1, a garbagecan 1, a
1626     laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a shelf 3,
1627     a shelf 2, and a shelf 1.
1628 > think: First I need to find a computer. From the observation image, a
1629     computer is on desk 1, from the dynamic model output, you should turn
1630     left to get direction.
1631 > turn left
1632 > success
1633 Environment: <observation image>. You are in the middle of a room.
1634     Looking quickly around you, you see a bed 1, a desk
1635     2, a desk 1, a computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer
1636     2, a drawer 1, a garbagecan 1, a
1637     laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a shelf 3,
1638     a shelf 2, and a shelf 1.
1639 > think: First I need to find a computer. From the observation image, a
1640     computer is on desk 1, from the dynamic model output, you should turn
1641     left to get direction.
1642 > turn left
1643 > success
1644 Environment: <observation image>. You are in the middle of a room.
1645     Looking quickly around you, you see a bed 1, a desk 2, a desk 1, a
1646     computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer
1647     1, a laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a
1648     shelf 3, a shelf 2, and a shelf 1, a laundryhamper 2.
1649 > think: First I need to find a computer. From the observation image, a
1650     computer is on desk 1. From the dynamic model output, you should walk
1651     forward to get around the bed.
1652 > walk forward
1653 > success
1654 Environment: <observation image>. You are in the middle of a room.
1655     Looking quickly around you, you see a bed 1, a desk 2, a desk 1, a
1656     computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer
1657     1, a laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a
1658     shelf 3, a shelf 2, and a shelf 1, a laundryhamper 2.
1659 > think: First I need to find a computer. From the observation image, a
1660     computer is on desk 1. From the dynamic model output, you should walk
1661     forward to get around the bed.
1662 > walk forward
1663 > success
1664 Environment: <observation image>. You are in the middle of a room.
1665     Looking quickly around you, you see a bed 1, a desk 2, a desk 1, a
1666     computer 1, a drawer 5, a drawer 4, a drawer 3, a drawer 2, a drawer
1667     1, a laundryhamper 1, a cabinet 1, a shelf 6, a shelf 5, a shelf 4, a
1668     shelf 3, a shelf 2, and a shelf 1, a laundryhamper 2.
1669 > think: First I need to find a computer. From the observation image, a
1670     computer is on desk 1. From the dynamic model output, you should walk
1671     forward to get around the bed.
1672 > walk forward
1673 > success

```

Status: Success

## I TRAINING DETAILS

**VQ-GAN** Our model was trained on the VH-1.5M dataset. The training set consists of motion trajectories from the first 49 rooms, while the last room was used as the validation set. Each room contains approximately 30,000 frames of images. The images were normalized and augmented using random cropping and horizontal flipping. We trained the model from scratch, where the input for each frame was the optical flow of the previous frame. The model was tasked with predicting the optical flow of the next frame based on this input. We used a batch size of 12 and trained the model for 50 epochs with an initial learning rate of  $3.5 \cdot 10^{-5}$ . The training process was conducted on eight NVIDIA A100 GPUs, each with 40GB of memory, and the total training time was approximately four days.

**Instructpix2pix** Our model was trained on the VH-1.5M dataset. The training set consists of motion trajectories from the first 49 rooms, while the last room was used as the validation set. Each room contains approximately 30,000 frames of images. The images were normalized and augmented using random cropping and horizontal flipping. We trained the model from pretraining, where the input for each frame was the previous frame. The model was tasked with predicting the next frame based on this input. We used a batch size of 32 and trained the model for 50000 epochs. We use cosine annealing to drop the learning rate from  $10^{-4}$  to  $10^{-5}$  for the first 20,000 training rounds. The training process was conducted on eight NVIDIA A100 GPUs, each with 40GB of memory, and the total training time was approximately two days.

**ControlNet** Our model was trained on the VH-1.5M dataset. The training set consists of motion trajectories from the first 49 rooms, while the last room was used as the validation set. Each room contains approximately 30,000 frames of images. The images were normalized and augmented using random cropping and horizontal flipping. We initialize the model weights to 0 then train the model from scratch, where the input for each frame was the optical flow of the previous frame. The model was tasked with predicting the next frame based on this input. We used a batch size of 24 and trained the model for 80000 epochs. We use cosine annealing to drop the learning rate from  $10^{-4}$  to  $10^{-5}$  for the first 40,000 training rounds. The training process was conducted on eight NVIDIA A100 GPUs, each with 40GB of memory, and the total training time was approximately four days.

## J COMPLEXITY ANALYSIS

**In the section, we will discuss why does our pipeline employ one-step planning? This is actually based on striking a balance between decision accuracy and complexity.** For our problem settings, we use Partially Observable Markov Decision Process (POMDP) (Smallwood & Sondik, 1973) to define the decision making process due to egocentric view is the partial observation for Egoplan agent. When we use our dynamics model to do multi-step prediction, due to the number of possible future states goes up very fast, can we guarantee a significant improvement in decision making (task completion success rate) without an explosive increase in the number of decisions in GPT4V? We calculated the relationship between action decision accuracy (compared to a skillfull human expert) and the number of GPT4V decisions for different dynamic model prediction steps in different tasks (the first six virtualhome tasks), and constructed a statistic  $\frac{\text{accuracy}}{\text{complexity}}$  (the larger statistic indicates the more "effective" and "skillfull" of agent's decision-making assisted by this dynamic model), the results are as shown in the Table 7. These results point out that in some long-horizon tasks with huge changes in perspective, when an agent with egocentric view is performing model predictive control, using some AI agent decision technology, multi-step prediction often brings a lot of computational complexity.

## K SUCCESS RATE (COMPLETE RATE) OF VIRTUALHOME TASKS

The final success of the long-range tasks on virtualhome tasks is shown in Figure 11. The final completion rate reflects the probability that the agent will reach the end point in the long-term

	1-step	2-step	3-step	4-step
take and place	7.12	1.21	0.17	0.07
take and put1	7.01	1.34	0.25	0.03
take and put2	6.98	1.14	0.29	0.06
take and drink	6.74	1.06	0.31	0.09
turn on sit	7.32	1.31	0.35	0.02
put apple	7.22	0.95	0.39	0.07

Table 7: The indicators of decision accuracy and decision numbers  $\frac{accuracy}{complexity}$  for dynamic model autoregressive prediction with different number of steps.

task trajectory, and thus reflects the stability of the pipeline. See Figure 10 for a more intuitive representation of Table 4.

## L REAL WORLD EXPERIMENTS

We conducted experiments in real scenarios. We placed the necessary experimental items in different rooms. In the subsequent version of our work, we will present the experimental environment in more detail. Now, we are presenting the success rate of the task in the real scenario. Compared to the baseline completion rate, Egoplan achieved a higher task completion rate in real scenarios, For our methods, we adopt the diffusion policy (Chi et al., 2023) method as our low-level policy. In real world tasks, we collected approximately 10 trajectories with about 100 frames to fine-tune our world model.

Since the Qwen-2.5-VL-32B model is much smaller than GPT4V and may not have been pre-trained on specialized embodied reasoning datasets, its performance is much worse compared to GPT4V. We also replaced our world model with Cosmos-Predict1-7B-Video2World, which supports both text and video input. We unified our input as the baseline of the current egocentric view (a single frame image) as Cosmos(frame). We found that if we input historical videos into Cosmos as Cosmos(video), then Cosmos would undergo significant improvements and approach the performance of our method (but our method only inputs current observations). We have attached the results in Tables 8 and 9.

Table 8: The number of success on 12 tasks for all the methods. Tasks 1-6 occur inside one room, while tasks 7-12 take place in two rooms. Each task was executed 100 times.

Task	Qwen-2.5+React	Qwen-2.5-VL	React	Reflexion	VL+P2P	Qwen-2.5-VL+OF
take and place	0	0	2	2	4	4
take and put1	0	0	2	3	5	4
take and put2	0	0	1	2	3	4
take and drink	0	0	1	2	4	5
turn on sit	0	0	1	2	4	4
put apple	0	0	1	1	4	4
take and place2	0	0	0	1	3	4
take and place3	0	0	1	1	3	4
take and put3	0	0	1	2	3	4
take open and put	0	0	1	2	4	4
take put and open	0	0	1	2	4	4
take and put4	0	0	0	1	2	3
Task	Cosmos(frame)	GR-SUSIE	GR-MG	Cosmos(video)	Ours(text goal)	Ours(image goal)
take and place	3	4	4	8	6	8
take and put1	4	4	3	8	6	9
take and put2	3	4	5	10	7	9
take and drink	4	4	3	6	6	8
turn on sit	4	4	3	6	5	6
put apple	3	3	4	7	6	9
take and place2	3	3	3	4	3	5
take and place3	3	3	3	6	4	8
take and put3	3	3	6	4	4	6
take open and put	3	3	4	6	5	5
take put and open	3	4	5	5	5	6
take and put4	3	3	4	5	5	4

We trained several end-to-end RL/Implicit Learning methods on our dataset. For these methods, we tried our best to uniformly select appropriate model architectures. Here are the various models we chose and their respective effects in Table 10.

Table 9: The number of success on 12 tasks for all the methods. Tasks 1-6 occur inside one room, while tasks 7-12 take place in two rooms. Each task was executed 100 times.

Task	GPT4+React	GPT4V	React	Reflexion	GPT4V+P2P
take and place	0	2	8	10	12
take and put1	0	2	7	8	10
take and put2	0	3	5	6	8
take and drink	0	1	3	5	8
turn on sit	0	2	5	6	8
put apple	0	3	3	5	10
take and place2	0	1	2	3	7
take and place3	0	1	3	6	9
take and put3	0	0	4	4	9
take open and put	0	1	3	4	7
take put and open	0	2	2	8	12
take and put4	0	1	4	6	9
Task	GPT4V+OF	SUSIE	GR-MG	Ours(text goal)	Ours(image goal)
take and place	14	11	14	16	21
take and put1	13	11	12	17	20
take and put2	10	9	10	14	18
take and drink	12	9	11	13	15
turn on sit	13	13	11	12	15
put apple	11	11	11	13	17
take and place2	9	5	6	9	10
take and place3	12	9	11	13	15
take and put3	11	9	9	11	13
take open and put	10	12	14	14	15
take put and open	12	12	11	14	14
take and put4	12	10	12	14	15

- **LCBC (Language-Conditioned Behavior Cloning)** For the LCBC baseline, we use the same architecture and hyperparameters as the low-level policy in SUSIE. We encode the language instruction using MUSE (Yinfei Yang et al. Multilingual Universal Sentence Encoder for Semantic Retrieval) and feed it into the ResNet-50 image encoder using FiLM conditioning.
- **PPO** We use ResNet-50 as image emcoder, 3 256-unit MLP layers are used as backbone. When PPO agent accomplish each sub-goal we give reward = 1. When PPO agent accomplish goal we give reward = 10.
- **GCBC (Goal Conditional Behavior Cloning)** For the GCBC baseline, we need to emphasize that this is actually the method of SUSIE. SUSIE uses `Instructpix2pix` to generate sugoa! and then applies the GCBC method in the downstream model.
- **GCIL (Goal Conditional Imitation Learning)** Observation and goal image are passed into ResNet-50 image encoder. 3 256-unit MLP layers are used as backbone.

Among these methods, we found that three types of information mainly guide the strategy learning: language instructions (LCBC), sub-goals (GCBG, GCIL), and rewards (PPO). Based on the results, the preliminary conclusion we can draw is that sub-goals are the most useful for guiding the learning of strategies.

Table 10: The number of success on 12 tasks for all the methods. Tasks 1-6 occur inside one room, while tasks 7-12 take place in two rooms. Each task was executed 100 times.

	LCBC	GCBC(SUSIE)	GCIL	PPO	Ours(image goal)
take and place	10	11	10	8	21
take and put1	8	11	12	9	20
take and put2	6	9	13	7	18
take and drink	5	9	8	10	15
turn on sit	5	13	12	10	15
put apple	8	11	10	10	17
take and place2	3	5	7	7	10
take and place3	4	9	7	8	15
take and put3	5	9	7	4	13
take open and put	5	12	10	9	15
take put and open	4	12	12	9	14
take and put4	4	10	14	10	15



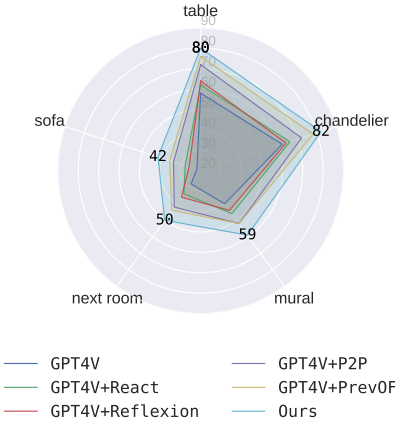
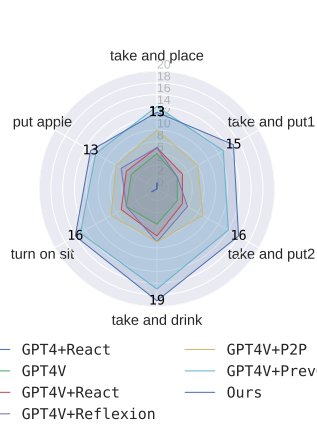
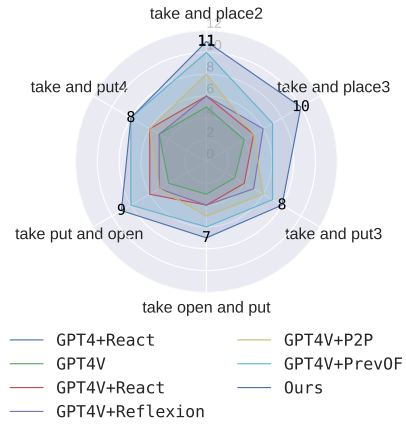


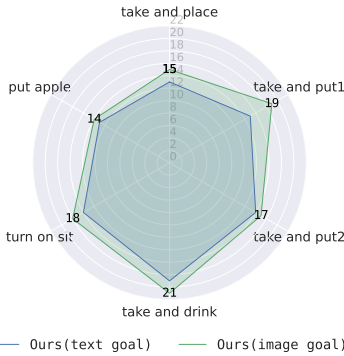
Figure 10: The success rate on 5 navigation tasks for all the methods in Habitat 2.0. GPT4+React is omitted due to its poor performance.



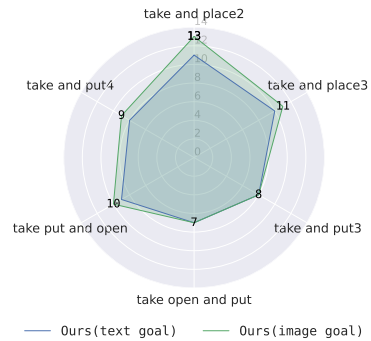
(a) Ours (text subgoal) and other baselines on task 1-6



(b) Ours (text subgoal) and other baselines on task 7-12



(c) Image subgoal and text subgoal on task 1-6, compared with SuSIE



(d) Image subgoal and text subgoal on task 7-12, compared with SuSIE

Figure 11: The success rate on 12 tasks for all the methods. Note that tasks 1-6 occur inside one room, while tasks 7-12 take place in two rooms.