
Compositional Foundation Models for Hierarchical Planning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To make effective decisions in novel environments with long-horizon goals, it is
2 crucial to engage in hierarchical reasoning across spatial and temporal scales. This
3 entails planning abstract subgoal sequences, visually reasoning about the under-
4 lying plans, and executing actions in accordance with the devised plan through
5 visual-motor control. We propose *Compositional Foundation Models for Hierarchi-
6 cal Planning (HiP)*, a foundation model which leverages multiple *expert* foundation
7 model trained on language, vision and action data *individually* jointly together to
8 solve long-horizon tasks. We use a large language model to construct symbolic
9 plans that are grounded in the environment through a large video diffusion model.
10 Generated video plans are then grounded to visual-motor control, through an in-
11 verse dynamics model that infers actions from generated videos. To enable effective
12 reasoning within this hierarchy, we enforce consistency between the models via
13 *iterative refinement*. We illustrate the efficacy and adaptability of our approach in
14 three different long-horizon table-top manipulation tasks.

15 1 Introduction

16 Consider the task of making a cup of tea in an unfamiliar house. To successfully execute this task, an
17 effective approach is to reason hierarchically at multiple levels: an abstract level, *e.g.* the high level
18 steps needed to heat up the tea, a concrete geometric level *e.g.*, how we should physically navigate
19 to and in the kitchen, and at a control level *e.g.* how we should actuate our joints to lift a cup. It is
20 further important that reasoning at each level is self-consistent with each other – an abstract plan to
21 look in cabinets for tea kettles must also be physically plausible at the geometric level and executable
22 given the actuations we are capable of. In this paper, we explore how we can create agents capable of
23 solving novel long-horizon tasks which require hierarchical reasoning.

24 Large “foundation models” have become a dominant paradigm in solving tasks in natural language
25 processing [36, 47, 7], computer vision [26], and mathematical reasoning [27]. In line with this
26 paradigm, a question of broad interest is to develop a “foundation model” that can solve novel and
27 long-horizon decision-making tasks. Some prior works [39, 6] collected paired visual, language
28 and action data and trained a monolithic neural network for solving long-horizon tasks. However,
29 collecting paired visual, language and action data is expensive and hard to scale up. Another line
30 of prior works [10, 28] finetune large language models (LLM) on both visual and language inputs
31 using task-specific robot demonstrations. This is problematic because, unlike the abundance of
32 text on the Internet, paired vision and language robotics demonstrations are not readily available
33 and are expensive to collect. Furthermore, finetuning high-performing language models, such as
34 GPT3.5/4 [37, 36] and PaLM [7], is currently impossible, as the model weights are not open-sourced.

35 The key characteristic of the foundation model is that solving a new task or adapting to a new
36 environment is possible with much less data compared to training from scratch for that task or

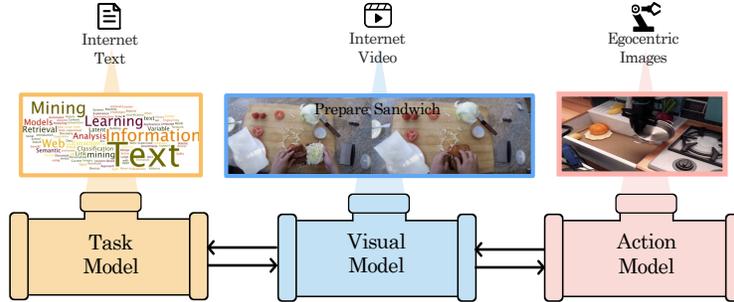


Figure 1: **Compositional Foundation Models for Hierarchical Planning.** HiP uses a task model, represented using a LLM, to create an abstract plan, a visual model, represented using a video model, to generate an image trajectory plan, and an ego-centric action model to infer actions from the image trajectory.

37 domain. Instead of building a foundation model for long-term planning by collecting paired language-
 38 vision-action data, in this work we seek a scalable alternative – can we reduce the need for a costly
 39 and tedious process of collecting paired data across three modalities and yet be relatively efficient
 40 at solving new planning tasks? We propose *Compositional Foundation Models for Hierarchical*
 41 *Planning* (HiP), a foundation model that is a composition of different *expert* models trained on
 42 language, vision, and action data *individually*. Because these models are trained individually, the
 43 data requirements for constructing the foundation models are substantially reduced (Figure 1). Given
 44 an abstract language instruction describing the desired task, HiP uses a large language model to
 45 find a sequence of sub-tasks (i.e., planning). HiP then uses a large video diffusion model to capture
 46 geometric and physical information about the world and generates a more detailed plan in form
 47 of an observation-only trajectory. Finally, HiP uses a large pre-trained inverse model that maps a
 48 sequence of ego-centric images into actions. The compositional design choice for decision-making
 49 allows separate models to reason at different levels of the hierarchy, and jointly make expert decisions
 50 without the need for ever collecting expensive paired decision-making data across modalities.

51 Given three models trained independently, they can produce inconsistent outputs that can lead to
 52 overall planning failure. For instance, a naïve approach for composing models is to take the maximum-
 53 likelihood output at each stage. However, a step of plan which is high likelihood under one model, *i.e.*
 54 looking for a tea kettle in a cabinet may have zero likelihood under a separate model, *i.e.* if there is no
 55 cabinet in the house. It is instead important to sample a plan that jointly maximizes likelihood across
 56 every expert model. To create consistent plans across our disparate models, we propose an *iterative*
 57 *refinement* mechanism to ensure consistency using feedback from the downstream models [28].
 58 At each step of the language model’s generative process, intermediate feedback from a likelihood
 59 estimator conditioned on an image of the current state is incorporated into the output distribution.
 60 Similarly, at each step of the video model generation, intermediate feedback from the action model
 61 refines video generation. This *iterative refinement procedure* promotes consensus among the different
 62 models and thereby enables hierarchically consistent plans that are both responsive to the goal
 63 and executable given the current state and agent. Our proposed iterative refinement approach is
 64 computationally efficient to train, as it does not require any large model finetuning. Furthermore, we
 65 do not require access to the model’s weights and our approach works with any models that offer only
 66 input and output API access.

67 In summary, we propose a compositional foundation model for hierarchical planning that leverages a
 68 composition of foundation models, learned separately on different modalities of Internet and ego-
 69 centric robotics data, to construct long-horizon plans. We demonstrate promising results on three
 70 long-horizon tabletop manipulation environments.

71 2 Compositional Foundation Models for Hierarchical Planning

72 We propose HiP, a foundation model that decomposes the problem of generating action trajectories
 73 for long-horizon tasks specified by a language goal g into three levels of hierarchy: **(1)** Task planning
 74 – inferring a language subgoal w_i conditioned on observation $x_{i,1}$ and language goal g ; **(2)** Visual
 75 planning – generating a physically plausible plan as a sequence of image trajectories $\tau_x^i = \{x_{i,1:T}\}$,
 76 one for each given language subgoal w_i and observation at first timestep $x_{i,1}$; **(3)** Action planning –

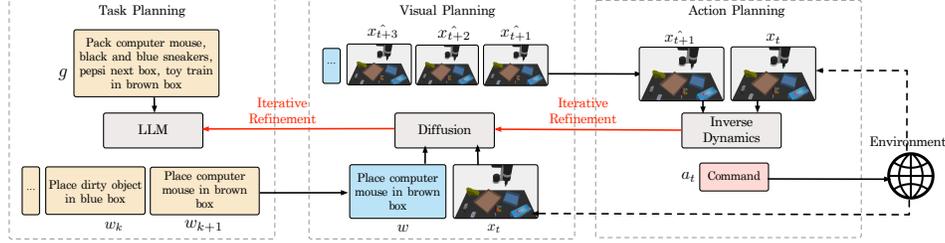


Figure 2: **Planning with HiP.** Given a language goal g and current observation x_t , LLM generates next subgoal w with feedback from a visual plausibility model. Then, Diffusion uses observation x_t and subgoal w to generate observation trajectory τ_x with feedback from an action feasibility model. Finally, action planning uses inverse dynamics to generate action a_t from current x_t and generated observation x_{t+1} (action planning).

77 inferring a sequence of action trajectories $\tau_a^i = \{a_{i,1:T-1}\}$ from the image trajectories τ_x^i executing
 78 the plan. Figure 2 illustrates the model architecture and a pseudocode is provided in Algorithm 1.

79 Let p_Θ model this hierarchical decision-making process. Given our three levels of hierarchy, p_Θ can
 80 be factorized into the following: task distribution p_θ , visual distribution p_ϕ , and action distribution
 81 p_ψ . The distribution over plans, conditioned on the goal and an image of the initial state, can be
 82 written under the Markov assumption as:

$$p_\Theta(W, \{\tau_x^i\}, \{\tau_a^i\} | g, x_{1,1}) = \underbrace{\left(\prod_{i=1}^N p_\theta(w_i | g) \right)}_{\text{task planning}} \underbrace{\left(\prod_{i=1}^N p_\phi(\tau_x^i | w_i, x_{i,1}) \right)}_{\text{visual planning}} \underbrace{\left(\prod_{i=1}^N \prod_{t=1}^{T-1} p_\psi(a_{i,t} | x_{i,t}, x_{i,t+1}) \right)}_{\text{action planning}} \quad (1)$$

83 We seek to find action trajectories τ_a^i , image trajectories τ_x^i and subgoals $W = \{w_i\}$ which maximize
 84 the above likelihood. Please see Appendix B for a derivation of this factorization. In the following
 85 sub-sections, we describe the form of each of these components, how they are trained, and how they
 86 are used to infer a final plan for completing the long-horizon task.

87 2.1 Task Planning via Large Language Models

88 Given a task specified in language g and the current observation $x_{i,1}$, we use a pretrained LLM as the
 89 task planner, which decomposes the goal into a sequence of subgoals. The LLM aims to infer the next
 90 subgoal w_i given a goal g and models the distribution $p_{\text{LLM}}(w_i | g)$. As the language model is trained
 91 on a vast amount of data on the Internet, it captures powerful semantic priors on what steps should be
 92 taken to accomplish a particular task. To adapt the LLM to obtain a subgoal sequence relevant to our
 93 task, we prompt it with some examples of domain specific data consisting of high-level goals paired
 94 with desirable subgoal sequences.

95 However, directly sampling subgoals using a LLM can lead to samples that are inconsistent with the
 96 overall joint distribution in Eqn (1), as the subgoal w_i not only affects the marginal likelihood of task
 97 planning but also the downstream likelihoods of the visual planning model. Consider the example in
 98 Figure 2 where the agent is tasked with packing computer mouse, black and blue sneakers, pepsi box
 99 and toy train in brown box. Let’s say the computer mouse is already in the brown box. While the
 100 subgoal of placing computer mouse in brown box has high-likelihood under task model $p_\theta(w_i | g)$, the
 101 resulting observation trajectory generated by visual model $p_\phi(\tau_x^i | w_i, x_{i,1})$ will have a low-likelihood
 102 under p_ϕ given the subgoal is already completed. Next, we describe how we use iterative refinement
 103 to capture this dependency between language decoding and visual planning to properly sample from
 104 Eqn (1).

105 **Consistency with Visual Planning.** To ensure that we sample subgoal w_i that maximizes the joint
 106 distribution in Eqn (1), we should sample a subgoal that maximizes the following joint likelihood

$$w_i^* = \arg \max_{w_i} p_{\text{LLM}}(w_i | g) p_\phi(\tau_x^i | w_i, x_{i,1}), \quad (2)$$

107 i.e. a likelihood that maximizes both conditional subgoal generation likelihood from a LLM and the
 108 likelihood of sampled videos τ_x^i given the language instruction and current image $x_{i,1}$. One way to
 109 determine the optimal subgoal w_i^* is to generate multiple w_i from LLM and score them using the

Algorithm 1 Decision Making with HiP

```
1: Models: Large language model  $p_{LLM}$ , Subgoal classifier  $f_\phi$ , Noise model of diffusion  $\epsilon_\phi$ , Observation  
   trajectory classifier  $g_\psi$ , Inverse dynamics  $p_\psi$   
2: Hyperparameters: Guidance scales  $\omega, \omega'$ , Denoising diffusion steps  $K$   
3: Input: Current observation  $x_t$ , Language goal  $g$   
4: # Task Planning  
5: for  $i = 1 \dots M$  do  
6:   Generate subgoal  $w_i \sim p_{LLM}(w_i|g)$   
7: end for  
8: Collect candidate subgoals  $W \leftarrow \{w_i\}_{i=1}^M$   
9: # Iterative Refinement from Visual Planning  
10:  $w \leftarrow \arg \max_w f_\phi(x_t, W, g)$   
11: # Visual Planning  
12: Initialize  $(\tau_x)_K \sim \mathcal{N}(0, I)$   
13: for  $k = K \dots 1$  do  
14:   # Iterative Refinement from Action Planning  
15:    $\hat{\epsilon} \leftarrow \epsilon_\phi((\tau_x)_k, x_t, k) + \omega(\epsilon_\phi((\tau_x)_k, x_t, w, k) - \epsilon_\phi((\tau_x)_k, x_t, k)) - \omega' \nabla_{(\tau_x)_k} \log g_\psi(1|(\tau_x)_k)$   
16:    $(\tau_x)_{k-1} \leftarrow \text{Denoise}((\tau_x)_k, \hat{\epsilon})$   
17: end for  
18: # Action Planning  
19: Extract  $(x_t, x_{t+1})$  from  $(\tau_x)_0$   
20: Execute  $a_t \leftarrow p_\psi(a_t|x_t, x_{t+1})$ 
```

110 likelihood of videos sampled from our video model $p_\phi(\tau_x^i|w_i, x_{i,1})$. However, the video generation
111 process is computationally expensive, so we take a different approach.

112 The likelihood of video generation $p_\phi(\tau_x^i|w_i, x_{i,1})$ primarily corresponds to the feasibility of a
113 language subgoal w_i with respect to the initial image $x_{i,1}$. Thus an approximation of Eqn (2) is to
114 directly optimize the conditional density

$$w_i^* = \arg \max_{w_i} p(w_i|g, x_{i,1}). \quad (3)$$

115 We can rewrite Eqn (3) as

$$w_i^* = \arg \max_{w_i} \log p_{LLM}(w_i|g) + \log \left(\frac{p(x_{i,1}|w_i, g)}{p(x_{i,1}|g)} \right)$$

116 We estimate the density ratio $\frac{p(x_{i,1}|w_i, g)}{p(x_{i,1}|g)}$ with a multi-class classifier $f_\phi(x_{i,1}, \{w_j\}_{j=1}^M, g)$ that
117 chooses the appropriate subgoal w_i^* from candidate subgoals $\{w_j\}_{j=1}^M$ generated by the LLM. The
118 classifier implicitly estimates the relative log likelihood estimate of $p(x_{i,1}|w_i, g)$ and use these logits
119 to estimate the log density ratio with respect to each of the M subgoals and find w_i^* that maximizes the
120 estimate [45]. We use a dataset $\mathcal{D}_{\text{classify}} := \{x_{i,1}, g, \{w_j\}_{j=1}^M, i\}$ consisting of observation $x_{i,1}$, goal
121 g , candidate subgoals $\{w_j\}_{j=1}^M$ and the correct subgoal label i to train f_ϕ . For further architectural
122 details, please refer to Appendix C.1.

123 2.2 Visual Planning with Video Generation

124 Upon obtaining a language subgoal w_i from task planning, our visual planner generates a plausible
125 observation trajectory τ_x^i conditioned on current observation $x_{i,1}$ and subgoal w_i . We use a video
126 diffusion model for visual planning given its success in generating text-conditioned videos [20, 53].
127 To provide our video diffusion model with a rich prior for physically plausible motions, we pretrain
128 it $p_\phi(\tau_x^i|w_i, x_{i,1})$ on a large-scale text-to-video dataset Ego4D [13]. We then finetune it on the
129 task-specific video dataset $\mathcal{D}_{\text{video}} := \{\tau_x^i, w_i\}$ consisting of observation trajectories τ_x^i satisfying
130 subgoal w_i . For further architectural details, please refer to Appendix C.2.

131 However, analogous to the consistent sampling problem in task planning, directly sampling observa-
132 tion trajectories with video diffusion can lead to samples that are inconsistent with the overall joint
133 distribution in Eqn (1). The observation trajectory τ_x^i not only affects the marginal likelihood of
134 visual planning, but also the downstream likelihood of the action planning model.

135 **Consistency with Action Planning.** To ensure observation trajectories τ_x^i that correctly maximize
 136 the joint distribution in Eqn (1), we optimize an observation trajectory that maximizes the following
 137 joint likelihood

$$(\tau_x^i)^* = \arg \max_{\tau_x^i} p_\phi(\tau_x^i | w_i, x_{i,1}) \prod_{t=1}^{T-1} p_\psi(a_{i,t} | x_{i,t}, x_{i,t+1}), \quad (4)$$

138 i.e. an image sequence that maximizes both conditional observation trajectory likelihood from video
 139 diffusion and the likelihood of sampled actions τ_a^i given the observation trajectory τ_x^i .

140 To sample such an observation trajectory, we could iteratively bias the denoising of video diffusion
 141 using the log-likelihood of the sampled actions $\prod_{t=1}^{T-1} \log p_\psi(a_{i,t} | x_{i,t}, x_{i,t+1})$. While this solution is
 142 principled, it is slow as it requires sampling of entire action trajectories and calculating the corre-
 143 sponding likelihoods during every step of the denoising process. Thus, we approximate the sampling
 144 and the likelihood calculation of action trajectory $\prod_{t=1}^{T-1} p_\psi(a_{i,t} | x_{i,t}, x_{i,t+1})$ with a binary classifier
 145 $g_\psi(\tau_x^i)$ that models if the observation trajectory τ_x^i leads to a high-likelihood action trajectory.

146 We learn a binary classifier g_ψ to assign high likelihood to feasible trajectories sampled from our video
 147 dataset $\tau_x^i \sim \mathcal{D}_{\text{video}}$ and low likelihood to infeasible trajectories generated by randomly shuffling
 148 the order of consecutive frames in feasible trajectories. Once trained, we can use the likelihood
 149 $\log g_\psi(1 | \tau_x^i)$ to bias the denoising of the video diffusion and maximize the likelihood of the ensuing
 150 action trajectory. For further details on binary classifier, please refer to Appendix D.2.

151 2.3 Action Planning with Inverse Dynamics

152 After generating an observation trajectory τ_x^i from visual planning, our action planner generates
 153 an action trajectory τ_a^i from the observation trajectory. We leverage egocentric internet images for
 154 providing our action planner with useful visual priors. Our action planner is parameterized as an
 155 inverse dynamics model [1, 38] that infers the action $a_{i,t}$ given the observation pair $(x_{i,t}, x_{i,t+1})$:

$$a_{i,t} \sim p_\psi(a_{i,t} | x_{i,t}, x_{i,t+1})$$

156 **Training.** To imbue the inverse dynamics p_ψ with useful visual priors, we initialize it with VC-
 157 1 [32] weights, pretrained on ego-centric images and ImageNet. We then finetune it on dataset
 158 $\mathcal{D}_{\text{inv}} := \{\tau_x^i, \tau_a^i\}$ consisting of paired observation and action trajectories by optimizing:

$$\max_{\psi} \mathbb{E}_{\tau \in \mathcal{D}_{\text{inv}}} [\log p_\psi(a_{i,t} | x_{i,t}, x_{i,t+1})]$$

159 For further architectural details, please refer to Appendix C.3.

160 3 Experimental Evaluations

161 We evaluate the ability of HiP to solve long-horizon planning tasks that are drawn from distributions
 162 with substantial variation, including the number and types of objects and their arrangements. We
 163 then study the effects of iterative refinement and of pretraining on overall performance of HiP. We
 164 also compare against an alternative strategy of visually grounding the LLM without any task-specific
 165 data. In addition, we study how granularity of subgoals affects HiP’s performance, ablate over
 166 choice of visual planning model and analyze sensitivity of iterative refinement to hyperparameters
 167 (Appendix F).

168 3.1 Evaluation Environments

169 We evaluate HiP on three environments, `paint-block`, `object-arrange`, and `kitchen-tasks`
 170 which are inspired by combinatorial planning tasks in Mao et al. [33], Shridhar et al. [43] and Xing
 171 et al. [50] respectively.

- 172 • `paint-block`: A robot has to manipulate blocks in the environment to satisfy language goal
 173 instructions, such as *stack pink block on yellow block and place green block right of them*. However,
 174 objects of correct colors may not be present in the environment, in which case, the robot needs to
 175 first pick up white blocks and put them in the appropriate color bowls to paint them. After that, it

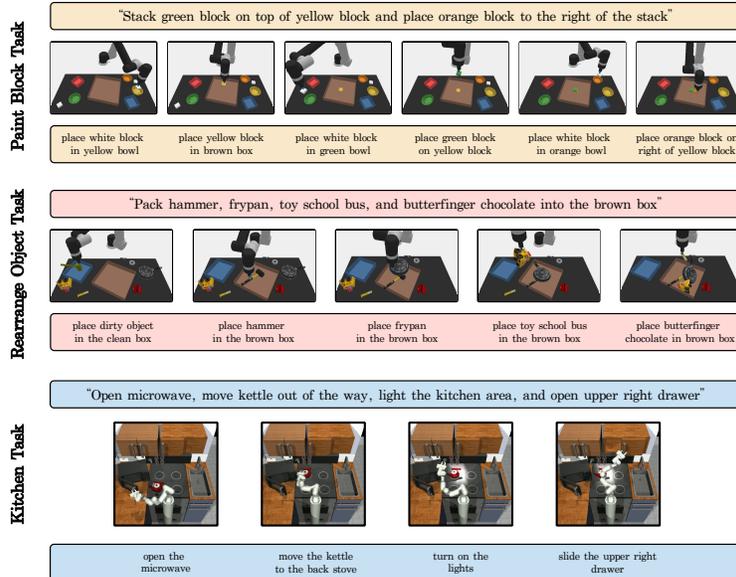


Figure 3: **Example Executions.** Example long-horizon generated plans on tasks in paint-block, object-arrange, and kitchen-tasks domains.

176 should perform appropriate pick-and-place operations to stack a pink block on the yellow block
 177 and place the green block right of them. A new task \mathcal{T} is generated by randomly selecting 3 final
 178 colors (out of 10 possible colors) for the blocks and then sampling a relation (out of 3 possible
 179 relations) for each pair of blocks. The precise locations of individual blocks, bowls, and boxes are
 180 fully randomized across different tasks. Tasks have 4 \sim 6 subgoals.

- 181 • **object-arrange:** A robot has to place appropriate objects in the brown box to satisfy language
 182 goal instructions such as *place shoe, tablet, alarm clock, and scissor in brown box*. However, the
 183 environment may have distractor objects that the robot must ignore. Furthermore, some objects
 184 can be dirty, indicated by a lack of texture and yellow color. For these objects, the robot must first
 185 place them in a blue cleaning box and only afterwards place those objects in the brown box. A new
 186 task \mathcal{T} is generated by randomly selecting 7 objects (out of 55 possible objects), out of which 3 are
 187 distractors, and then randomly making one non-distractor object dirty. The precise locations of
 188 individual objects and boxes are fully randomized across tasks. Tasks usually have 3 \sim 5 subgoals.
- 189 • **kitchen-tasks:** A robot has to complete kitchen subtasks to satisfy language goal instructions
 190 such as *open microwave, move kettle out of the way, light the kitchen area, and open upper right*
 191 *drawer*. However, the environment may have objects irrelevant to the subtasks that the robot
 192 must ignore. Furthermore, some kitchen subtasks specified in the language goal might already be
 193 completed, and the robot should ignore those tasks. There are 7 possible kitchen subtasks: opening
 194 the microwave, moving the kettle, switching on lights, turning on the bottom knob, turning on
 195 the top knob, opening the left drawer, and opening the right drawer. A new task \mathcal{T} is generated
 196 by randomly selecting 4 out of 7 possible kitchen subtasks, randomly selecting an instance of
 197 microwave out of 3 possible instances, randomly selecting an instance of kettle out of 4 possible
 198 instances, randomly and independently selecting texture of counter, floor and drawer out of 3
 199 possible textures and randomizing initial pose of kettle and microwave. With 50% probability, one
 200 of 4 selected kitchen subtask is completed before the start of the task. Hence, tasks usually have
 201 3 \sim 4 subtasks (i.e. subgoals).

202 **Train and Test Tasks.** For all environments, we sample two sets of tasks $\mathcal{T}_{train}, \mathcal{T}_{test} \sim p(\mathcal{T})$. We
 203 use the train set of tasks \mathcal{T}_{train} to create datasets $\mathcal{D}_{classify}, \mathcal{D}_{video}, \mathcal{D}_{inv}$ and other datasets required
 204 for training baselines. We ensure the test set of tasks \mathcal{T}_{test} contains novel combinations of object
 205 colors in paint-block, novel combinations of object categories in object-arrange, and novel
 206 combinations of kitchen subtasks in kitchen-tasks.

207 **Evaluation Metrics.** We quantitatively evaluate a model by measuring its task completion rate for
 208 paint-block and object-arrange, and subtask completion rate for kitchen tasks. We use the
 209 simulator to determine if the goal, corresponding to a task, has been achieved. We evaluate a model

210 on \mathcal{T}_{train} (seen) to test its ability to solve long-horizon tasks and on \mathcal{T}_{test} (unseen) to test its ability to
 211 generalize to long-horizon tasks consisting of novel combinations of object colors in `paint-block`,
 212 object categories in `object-arrange`, and kitchen subtasks in `kitchen-tasks`. We sample 1000
 213 tasks from \mathcal{T}_{train} and \mathcal{T}_{test} respectively, and obtain average task completion rate on `paint-block`
 214 and `object-arrange` domains and average subtask completion rate on `kitchen-tasks` domain.
 215 We report the mean and the standard error over those seeds in Table 1.

216 3.2 Baselines

217 There are several existing strategies for constructing robot manipulation policies conditioned on
 218 language goals, which we use as baselines in our experiments:

- 219 • **Goal-Conditioned Policy** A goal-conditioned transformer model $a_{i,t} \sim p(a_{i,t}|x_{i,t}, w_i)$ that out-
 220 puts action $a_{i,t}$ given a language subgoal w_i and current observation $x_{i,t}$ (Transformer BC) [6]. We
 221 provide the model with oracle subgoals and encode these subgoals with a pretrained language en-
 222 coder (Flan-T5-Base). We also compare against goal-conditioned policy with Gato [39] transformer
 223 architecture.
- 224 • **Video Planner** A video diffusion model (UniPi) [12] $\{\tau_x^i\} \sim p(\{\tau_x^i\}|g, x_{i,1}), a_{i,t} \sim$
 225 $p(a_{i,t}|x_{i,t}, x_{i,t+1})$ that bypasses task planning, generates video plans for the entire task $\{\tau_x^i\}$,
 226 and infers actions $a_{i,t}$ using an inverse model.
- 227 • **Action Planners** Transformer models (Trajectory Transformer) [23] and diffusion models (Dif-
 228 fuser) [24, 4] $\{a_{i,t:T-1}\} \sim p(\{a_{i,t:T-1}\}|x_{i,t}, w_i)$ that produce an action sequence $\{a_{i,t:T-1}\}$
 229 given a language subgoal w_i and current visual observation $x_{i,t}$. We again provide the agents with
 230 oracle subgoals and encode these subgoals with a pretrained language encoder (Flan-T5-Base).
- 231 • **LLM as Skill Manager** A hierarchical system (SayCan) [2] with LLM as high level policy that
 232 sequences skills sampled from a repertoire of skills to accomplish a long-horizon task. We use
 233 CLIPort [43] policies as skills and the unnormalized logits over the pixel space it produces as
 234 affordances. These affordances grounds the LLM to current observation for producing next subgoal.

235 3.3 Results

Model	Paint-block		Object-arrange		Kitchen-tasks	
	Seen	Unseen	Seen	Unseen	Seen	Unseen
Transformer BC (oracle subgoals)	8.3 ± 1.9	5.1 ± 1.6	10.2 ± 2.9	7.3 ± 1.7	48.4 ± 21.6	32.1 ± 24.2
Gato (oracle subgoals)	31.2 ± 2.4	28.6 ± 2.9	37.9 ± 3.3	36.5 ± 3.2	70.2 ± 10.8	66.8 ± 12.2
Trajectory Transformer (oracle subgoals)	22.1 ± 2.1	22.3 ± 2.5	30.5 ± 2.3	29.8 ± 2.9	66.4 ± 20.7	52.1 ± 22.3
Action Diffuser (oracle subgoals)	21.6 ± 2.6	18.2 ± 2.3	29.2 ± 2.4	27.6 ± 2.1	65.9 ± 23.2	55.1 ± 22.8
HiP (Ours, oracle subgoals)	81.2 ± 1.8	79.6 ± 1.9	91.8 ± 2.9	92.3 ± 2.3	92.8 ± 7.1	89.8 ± 7.6
UniPi	37.2 ± 3.8	35.3 ± 3.2	44.1 ± 3.1	44.2 ± 2.9	74.6 ± 14.8	73.4 ± 11.2
SayCan	67.2 ± 3.3	62.8 ± 3.7	70.3 ± 2.6	66.9 ± 2.8	-	-
HiP (Ours)	74.3 ± 1.9	72.8 ± 1.7	75 ± 2.8	75.4 ± 2.6	85.8 ± 9.4	83.5 ± 10.2

Table 1: **Performance on Long-Horizon tasks.** HiP not only outperforms the baselines in solving seen long-horizon tasks but its performance remains intact when solving unseen long-horizon tasks containing novel combination of objects colors in `paint-block`, novel combination of objects categories in `object-rearrange` and novel combination of subtasks in `kitchen-tasks`.

236 We begin by comparing the performance of HiP and baselines to solve long-horizon tasks in
 237 `paint-block`, `object-arrange`, `kitchen-tasks` environments. Table 1 shows that HiP signifi-
 238 cantly outperforms the baselines, although the baselines have an advantage and have access to oracle
 239 subgoals. HiP’s superior performance shows the importance of (i) hierarchy given it outperforms goal-
 240 conditioned policy (Transformer BC and Gato), (ii) task planning since it outperforms video planners
 241 (UniPi), and (iii) visual planning given it outperforms action planners (Trajectory Transformer, Action
 242 Diffuser). It also shows the importance of representing skills with video-based planners which can
 243 be pre-trained on Internet videos and can be applied to tasks (such as `kitchen-tasks`). SayCan, in
 244 contrast, requires tasks to be decomposed into primitives paired with an affordance function, which
 245 can be difficult to define for many tasks like the kitchen task. Thus, we couldn’t run SayCan on
 246 `kitchen-tasks` environment. Finally, to quantitatively show how the errors in $f_\theta(x_{i,1}, w_i, g)$ affect
 247 the performance of HiP, we compare it to HiP with oracle subgoals. For further details on the training
 248 and evaluation of HiP, please refer to Appendix D. For implementation details on Gato and SayCan,
 249 please refer to Appendix E. For runtime analysis of HiP, please refer to Appendix G.

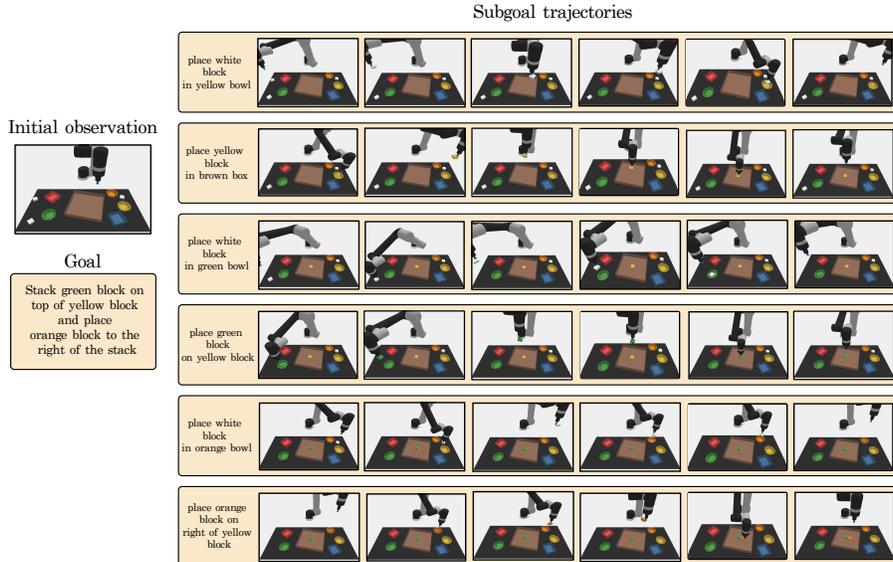


Figure 4: Execution trajectory of HiP on a novel long-horizon task in paint-block environment.

250 **Combinatorial Generalization to Unseen Long-horizon Tasks.** We also quantitatively test the
 251 ability of HiP to generalize to unseen long-horizon tasks, consisting of novel combinations of object
 252 colors in paint-block, object categories in object-arrange, and subtasks in kitchen-tasks.
 253 Table 1 shows that HiP’s performance remains intact when solving unseen long-horizon tasks, and
 254 still significantly outperforms the baselines. Figure 4 visualizes the execution of HiP in unseen
 255 long-horizon tasks in paint-block.

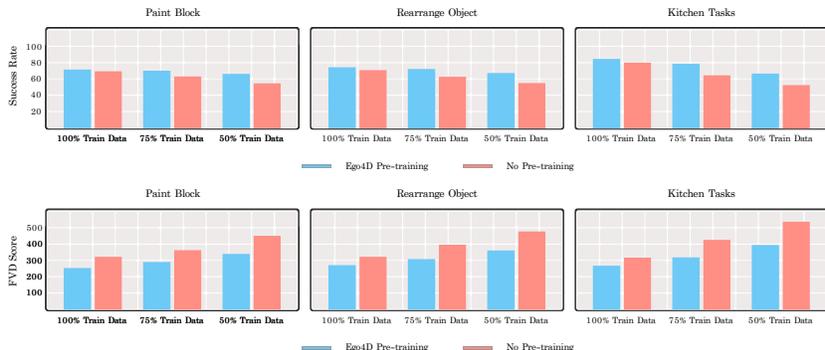


Figure 5: Pretraining video diffusion model with the Ego4D dataset consistently yields higher success rate and lower FVD scores (lower is better), even with reduced training dataset sizes. With pretraining, the model’s FVD score escalates less gradually and its success rate falls less steeply as the dataset size shrinks.

256 **Pre-training Video Diffusion Model.** We investigate how much our video diffusion model benefits
 257 from pre-training on the Internet-scale data. We report both the success rate of HiP and Fréchet
 258 Video Distance (FVD) score that quantifies the similarity between generated videos and ground truth
 259 videos, where lower scores indicate greater similarity in Figure 5. We see that pretraining video
 260 diffusion leads to a higher success rate and lower FVD score. If we reduce the training dataset to 75%
 261 and 50% of the original dataset, the FVD score for video diffusion models (both, with and without
 262 Ego4D dataset pretraining) increases and their success rate falls. However, the video diffusion model
 263 with Ego4D dataset pretraining consistently gets higher success rate and lower FVD scores across
 264 different dataset sizes. As we decrease the domain-specific training data, it is evident that the gap in
 265 performance between the model with and without the Ego4D pre-training widens. For details on how
 266 we process the Ego4D dataset, please refer to Appendix D.2.

267 **Pre-training Inverse Dynamics Model.** We also analyze the benefit of pre-training our inverse
 268 dynamics model and report the mean squared error between the predicted and ground truth actions

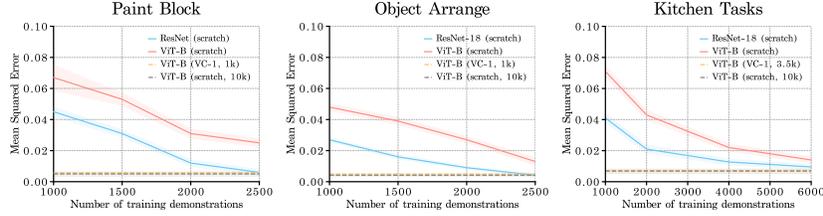


Figure 6: **Pretraining inverse dynamics model.** In paint-block and object-arrange (kitchen-tasks), when initialized with VC-1 weights, inverse dynamics model matches the performance of a randomly initialized model trained on 10K trajectories with just 1K (3.5k) trajectories. Even a smaller ResNet-18 model requires 2.5K (6k) trajectories to approach the same performance.

269 in Figure 6. The pre-training comes in the form of initializing the inverse dynamics model with
 270 weights from VC-1 [32], a vision-transformer (ViT-B) [9] trained on ego-centric images with masked-
 271 autoencoding objective [16]. In paint-block and object-arrange, we see that the inverse
 272 dynamics, when initialized with weights from VC-1, only requires 1K labeled robotic trajectories
 273 to achieve the same performance as the inverse dynamics model trained on 10K labeled robotic
 274 trajectories but without VC-1 initialization. We also compare against an inverse dynamics model
 275 parameterized with a smaller network (ResNet-18). However, the resulting inverse dynamics model
 276 still requires 2.5K robotic trajectories to get close to the performance of the inverse dynamics
 277 model with VC-1 initialization in paint-block and object-arrange. In kitchen-tasks, inverse
 278 dynamics, when initialized with weights from VC-1, only requires 3.5k labeled robotic trajectories
 279 to achieve the same performance as the inverse dynamics model trained on 10K labeled robotic
 280 trajectories but without VC-1 initialization. When parameterized with ResNet-18, the inverse
 281 dynamics model still requires 6k robotic trajectories to get close to the performance of the inverse
 282 dynamics model with VC-1 initialization.

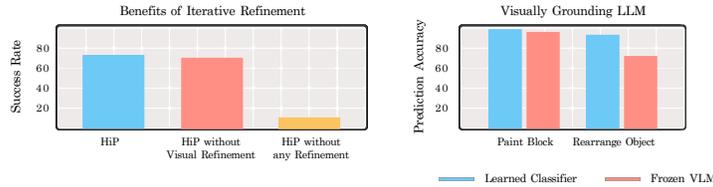


Figure 7: **Ablation Studies.** (Left) While task plan refinement is critical to HiP’s performance, visual plan refinement improves HiP’s performance by a smaller margin in paint block environment. (Right) While frozen pretrained VLM (MiniGPT4) matches the performance of a learned classifier in paint-block environment, its performance deteriorates in a more visually complex object-arrange environment.

283 **Importance of Task Plan and Visual Plan Refinements.** We study the importance of refinement
 284 in task and visual planning in Figure 7. We compare to HiP without visual plan refinement and
 285 HiP without visual and task plan refinement in paint block environment. We see that task plan
 286 refinement for visual grounding of LLM is critical to the performance of HiP. Without it, the task
 287 plan is agnostic to the robot’s observation and predicts subgoals that lead to erroneous visual and
 288 action planning. Furthermore, visual plan refinement improves the performance of HiP as well, albeit
 289 by a small margin. For more details on the hyperparameters used, please refer to Appendix D.4.

290 **Exploring Alternate Strategies for Visually Grounding LLM.** We use a learned classifier
 291 $f_{\theta}(x_{i,1}, w_i, g)$ to visually ground the LLM. We explore if we can use a frozen pretrained Vision-
 292 Language Model (MiniGPT-4 [58]) as a classifier in place of the learned classifier. Although we
 293 didn’t use any training data, we found the prompt engineering using the domain knowledge of
 294 the task to be essential in using the Vision-Language Model (VLM) as a classifier (see Appendix
 295 D.1 for details). We use subgoal prediction accuracy to quantify the performance of the learned
 296 classifier and the frozen VLM. Figure 7 illustrates that while both our learned multi-class classifier
 297 and frozen VLM perform comparably in the paint-block environment, the classifier significantly
 298 outperforms the VLM in the more visually complex object-arrange environment. We detail the
 299 two common failure modes of the VLM approach in object-arrange environment in Appendix
 300 D.1. As VLMs continue to improve, it is possible that their future versions match the performance of
 301 learned classifiers and thus replace them in visually complex domains as well. For further details on
 302 the VLM parameterization, please refer to Appendix D.1.

References

- 303
- 304 [1] P. Agrawal, A. V. Nair, P. Abbeel, J. Malik, and S. Levine. Learning to poke by poking:
305 Experiential learning of intuitive physics. *Advances in neural information processing systems*,
306 29, 2016.
- 307 [2] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, K. Gopalakrishnan,
308 K. Hausman, A. Herzog, et al. Do as i can, not as i say: Grounding language in robotic
309 affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- 310 [3] A. Ajay, A. Kumar, P. Agrawal, S. Levine, and O. Nachum. Opal: Offline primitive discovery
311 for accelerating offline reinforcement learning. *arXiv preprint arXiv:2010.13611*, 2020.
- 312 [4] A. Ajay, Y. Du, A. Gupta, J. Tenenbaum, T. Jaakkola, and P. Agrawal. Is conditional generative
313 modeling all you need for decision-making? *arXiv preprint arXiv:2211.15657*, 2022.
- 314 [5] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican,
315 M. Reynolds, et al. Flamingo: A visual language model for few-shot learning. *NeurIPS*, 2022.
316 URL <https://arxiv.org/abs/2204.14198>.
- 317 [6] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Haus-
318 man, A. Herzog, J. Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv*
319 *preprint arXiv:2212.06817*, 2022.
- 320 [7] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W.
321 Chung, C. Sutton, S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv*
322 *preprint arXiv:2204.02311*, 2022.
- 323 [8] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, E. Li, X. Wang, M. De-
324 hghani, S. Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint*
325 *arXiv:2210.11416*, 2022.
- 326 [9] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
327 M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for
328 image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 329 [10] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson,
330 Q. Vuong, T. Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint*
331 *arXiv:2303.03378*, 2023.
- 332 [11] Y. Du, S. Li, and I. Mordatch. Compositional visual generation with energy based models. In
333 *Advances in Neural Information Processing Systems*, 2020.
- 334 [12] Y. Du, M. Yang, B. Dai, H. Dai, O. Nachum, J. B. Tenenbaum, D. Schuurmans, and P. Abbeel.
335 Learning universal policies via text-guided video generation. *arXiv e-prints*, pages arXiv-2302,
336 2023.
- 337 [13] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang,
338 M. Liu, X. Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In
339 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
340 18995–19012, 2022.
- 341 [14] M. U. Gutmann and A. Hyvärinen. Noise-contrastive estimation of unnormalized statistical
342 models, with applications to natural image statistics. *J. Mach. Learn. Res.*, 13:307–361, 2012.
- 343 [15] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In
344 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–
345 778, 2016.
- 346 [16] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable
347 vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
348 *Recognition*, pages 16000–16009, 2022.
- 349 [17] J. Hermans, V. Begy, and G. Louppe. Likelihood-free mcmc with amortized approximate ratio
350 estimators. In *International conference on machine learning*, pages 4239–4248. PMLR, 2020.

- 351 [18] J. Ho and T. Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
352 2022.
- 353 [19] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural*
354 *Information Processing Systems*, 2020.
- 355 [20] J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi,
356 D. J. Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv*
357 *preprint arXiv:2210.02303*, 2022.
- 358 [21] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners:
359 Extracting actionable knowledge for embodied agents. In *International Conference on Machine*
360 *Learning*, pages 9118–9147. PMLR, 2022.
- 361 [22] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch,
362 Y. Chebotar, et al. Inner monologue: Embodied reasoning through planning with language
363 models. *arXiv preprint arXiv:2207.05608*, 2022.
- 364 [23] M. Janner, Q. Li, and S. Levine. Offline reinforcement learning as one big sequence modeling
365 problem. In *Advances in Neural Information Processing Systems*, 2021.
- 366 [24] M. Janner, Y. Du, J. Tenenbaum, and S. Levine. Planning with diffusion for flexible behavior
367 synthesis. In *International Conference on Machine Learning*, 2022.
- 368 [25] M. Janner, Y. Du, J. Tenenbaum, and S. Levine. Planning with diffusion for flexible behavior
369 synthesis. In *International Conference on Machine Learning*, 2022.
- 370 [26] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C.
371 Berg, W.-Y. Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- 372 [27] A. Lewkowycz, A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. Ramasesh, A. Slone,
373 C. Anil, I. Schlag, T. Gutman-Solo, et al. Solving quantitative reasoning problems with language
374 models. *arXiv preprint arXiv:2206.14858*, 2022.
- 375 [28] S. Li, X. Puig, C. Paxton, Y. Du, C. Wang, L. Fan, T. Chen, D.-A. Huang, E. Akyürek,
376 A. Anandkumar, et al. Pre-trained language models for interactive decision-making. *Advances*
377 *in Neural Information Processing Systems*, 35:31199–31212, 2022.
- 378 [29] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as
379 policies: Language model programs for embodied control. *arXiv preprint arXiv:2209.07753*,
380 2022.
- 381 [30] K. Lin, C. Agia, T. Migimatsu, M. Pavone, and J. Bohg. Text2motion: From natural language
382 instructions to feasible plans. *arXiv preprint arXiv:2303.12153*, 2023.
- 383 [31] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint*
384 *arXiv:1711.05101*, 2017.
- 385 [32] A. Majumdar, K. Yadav, S. Arnaud, Y. J. Ma, C. Chen, S. Silwal, A. Jain, V.-P. Berges,
386 P. Abbeel, J. Malik, et al. Where are we in the search for an artificial visual cortex for embodied
387 intelligence? *arXiv preprint arXiv:2303.18240*, 2023.
- 388 [33] J. Mao, T. Lozano-Perez, J. B. Tenenbaum, and L. P. Kaelbling. PDSketch: Integrated Domain
389 Programming, Learning, and Planning. In *Advances in Neural Information Processing Systems*
390 *(NeurIPS)*, 2022.
- 391 [34] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual representation
392 for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- 393 [35] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and
394 M. Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion
395 models. *arXiv preprint arXiv:2112.10741*, 2021.
- 396 [36] OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023.

- 397 [37] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,
398 K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback.
399 *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- 400 [38] D. Pathak, P. Mahmoudieh, G. Luo, P. Agrawal, D. Chen, Y. Shentu, E. Shelhamer, J. Malik,
401 A. A. Efros, and T. Darrell. Zero-shot visual imitation. In *Proceedings of the IEEE conference*
402 *on computer vision and pattern recognition workshops*, pages 2050–2053, 2018.
- 403 [39] S. Reed, K. Zolna, E. Parisotto, S. G. Colmenarejo, A. Novikov, G. Barth-Maron, M. Gimenez,
404 Y. Sulsky, J. Kay, J. T. Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*,
405 2022.
- 406 [40] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis
407 with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
408 *and Pattern Recognition*, pages 10684–10695, 2022.
- 409 [41] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. Denton, S. K. S. Ghasemipour, B. K. Ayan,
410 S. S. Mahdavi, R. G. Lopes, et al. Photorealistic text-to-image diffusion models with deep
411 language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- 412 [42] T. Salimans and J. Ho. Progressive distillation for fast sampling of diffusion models. *arXiv*
413 *preprint arXiv:2202.00512*, 2022.
- 414 [43] M. Shridhar, L. Manuelli, and D. Fox. Cliport: What and where pathways for robotic manipula-
415 tion. In *Conference on Robot Learning*, pages 894–906. PMLR, 2022.
- 416 [44] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning
417 using nonequilibrium thermodynamics. In *International Conference on Machine Learning*,
418 2015.
- 419 [45] A. Srivastava, S. Han, K. Xu, B. Rhodes, and M. U. Gutmann. Estimating the density ratio
420 between distributions with high discrepancy using multinomial logistic regression. *ArXiv*,
421 abs/2305.00869, 2023.
- 422 [46] M. Sugiyama, T. Suzuki, and T. Kanamori. *Density ratio estimation in machine learning*.
423 Cambridge University Press, 2012.
- 424 [47] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal,
425 E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv*
426 *preprint arXiv:2302.13971*, 2023.
- 427 [48] J. Urain, N. Funk, G. Chalvatzaki, and J. Peters. Se (3)-diffusionfields: Learning cost functions
428 for joint grasp and motion optimization through diffusion. *arXiv preprint arXiv:2209.03855*,
429 2022.
- 430 [49] Z. Wang, J. J. Hunt, and M. Zhou. Diffusion policies as an expressive policy class for offline
431 reinforcement learning. *arXiv preprint arXiv:2208.06193*, 2022.
- 432 [50] E. Xing, A. Gupta, S. Powers, and V. Dean. Kitchenshift: Evaluating zero-shot general-
433 ization of imitation-based policy learning under domain shifts. In *NeurIPS 2021 Work-*
434 *shop on Distribution Shifts: Connecting Methods and Applications*, 2021. URL <https://openreview.net/forum?id=DdglKo8hBq0>.
435
- 436 [51] M. Yang, Y. Du, B. Dai, D. Schuurmans, J. B. Tenenbaum, and P. Abbeel. Probabilistic
437 adaptation of text-to-video models. *arXiv preprint arXiv:2306.01872*, 2023.
- 438 [52] S. Yang, O. Nachum, Y. Du, J. Wei, P. Abbeel, and D. Schuurmans. Foundation models for
439 decision making: Problems, methods, and opportunities. *arXiv preprint arXiv:2303.04129*,
440 2023.
- 441 [53] S. Yu, K. Sohn, S. Kim, and J. Shin. Video probabilistic diffusion models in projected latent
442 space. *arXiv preprint arXiv:2302.07685*, 2023.

- 443 [54] A. Zeng, A. Wong, S. Welker, K. Choromanski, F. Tombari, A. Purohit, M. Ryoo, V. Sindhwani,
444 J. Lee, V. Vanhoucke, et al. Socratic models: Composing zero-shot multimodal reasoning with
445 language. *arXiv preprint arXiv:2204.00598*, 2022.
- 446 [55] E. Zhang, Y. Lu, W. Wang, and A. Zhang. Lad: Language augmented diffusion for reinforcement
447 learning. *arXiv preprint arXiv:2210.15629*, 2022.
- 448 [56] H. Zheng, W. Nie, A. Vahdat, K. Azizzadenesheli, and A. Anandkumar. Fast sampling of
449 diffusion models via operator learning. In *International Conference on Machine Learning*,
450 pages 42390–42402. PMLR, 2023.
- 451 [57] Z. Zhong, D. Rempe, D. Xu, Y. Chen, S. Veer, T. Che, B. Ray, and M. Pavone. Guided
452 conditional diffusion for controllable traffic simulation. *arXiv preprint arXiv:2210.17366*, 2022.
- 453 [58] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny. Minigpt-4: Enhancing vision-language
454 understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

455 Appendix

456 In this Appendix, we first discuss related work in Section A and then describe how we factorize the
 457 hierarchical decision-making process in Section B. In Section C, we then detail the background
 458 and architecture for visually grounded task planning, visual planning with video diffusion, and
 459 action planning with inverse dynamics model. In Section D, we discuss the training and evaluation
 460 details for the different levels of planning in HiP and the corresponding training hyperparameters.
 461 In Section E, we discuss implementation details for Gato and SayCan. In Section F, we showcase
 462 additional ablation studies comparing different approaches to enforce consistency across the levels of
 463 hierarchy, analyzing effect of granularity of subgoals on performance of HiP, ablating on the choice
 464 of video planning model and analyzing sensitivity of iterative refinement to hyperparameters. In
 465 Section G, we analyze runtime of different components of HiP. Finally, we discuss limitation and
 466 conclusion of our work in Section H.

467 A Related Work

468 The field of foundation models for decision-making [52] has seen significant progress in recent years.
 469 A large body of work explored using large language models as zero-shot planners [21, 22, 29, 30, 2],
 470 but it is often difficult to directly ground the language model on vision. To address this problem
 471 of visually grounding the language model, other works have proposed to directly fine-tune large
 472 language models for embodied tasks [28, 10]. However, such an approach requires large paired
 473 vision and language datasets that are difficult to acquire. Most similar to our work, SayCan [2]
 474 uses an LLM to hierarchically execute different tasks by breaking language goals into a sequence of
 475 instructions, which are then inputted to skill-based value functions. While SayCan assumes this fixed
 476 set of skill-based value functions, our skills are represented as video-based planners [12], enabling
 477 generalization to new skills.

478 Another set of work has explored how to construct continuous space planners with diffusion mod-
 479 els [25, 4, 49, 55, 48, 57, 12]. Existing works typically assume task-specific datasets from which the
 480 continuous-space planner is derived [24, 3, 55]. Most similar to our work, UniPi [12] proposes to use
 481 videos to plan in image space and similarly relies on internet videos to train image space planners.
 482 We build on top of UniPi to construct our foundation model for hierarchical planning, and illustrate
 483 how UniPi may be combined with LLMs to construct longer horizon continuous video plans.

484 Moreover, some works [28, 54, 5] explored how different foundation models may be integrated with
 485 each other. In Flamingo [5], models are combined through joint finetuning with paired datasets, which
 486 are difficult to collect. In contrast both Zeng et al. [54] and Li et al. [28] combine different models
 487 zero-shot using either language or iterative consensus. Our work proposes to combine language,
 488 video, and ego-centric action models together by taking the product of their learned distributions
 489 [11]. We use a similar iterative consensus procedure as in Li et al. [28] to sample from the entire joint
 490 distribution and use this combined distribution to construct a hierarchical planning system.

491 B Factorizing Hierarchical Decision-Making Process

492 We model the hierarchical decision-making process described in Section 2 with p_{Θ} which can be
 493 factorized into the task distribution p_{θ} , visual distribution p_{ϕ} , and action distribution p_{ψ} .

$$p_{\Theta}(W, \{\tau_x^i\}, \{\tau_a^i\} | g, x_{1,1}) = \prod_{i=1}^N p_{\theta}(w_i | g, x_{i,1}, w_{<i}, \tau_x^{<i}, \tau_a^{<i}) \prod_{i=1}^N p_{\phi}(\tau_x^i | w_{\leq i}, x_{i,1}, g, \tau_x^{<i}, \tau_a^{<i}) \\ \prod_{i=1}^N p_{\psi}(\tau_a^i | \tau_x^{\leq i}, w_{\leq i}, x_{i,1}, g, \tau_a^{<i})$$

494 Here, given random variables Y^i , $Y^{<i}$ and $Y^{\leq i}$ represents $\{Y^1, \dots, Y^{i-1}\}$ and $\{Y^1, \dots, Y^i\}$
 495 respectively. Now, we apply Markov assumption: Given current observation $x_{i,1}$, future variables
 496 $(w_i, \tau_x^i, \tau_a^i)$ and past variables $(w_j, \tau_x^j, \tau_a^j \ \forall j < i)$ are conditionally independent.

$$p_{\Theta}(W, \{\tau_x^i\}, \{\tau_a^i\} | g, x_{1,1}) = \prod_{i=1}^N p_{\theta}(w_i | g, x_{i,1}) \prod_{i=1}^N p_{\phi}(\tau_x^i | w_i, x_{i,1}, g) \prod_{i=1}^N p_{\psi}(\tau_a^i | \tau_x^i, w_i, x_{i,1}, g)$$

497 We model task distribution p_θ with a large language model (LLM) which is independent of observation
 498 $x_{i,1}$. Since the image trajectory $\tau_x^i = \{x_{i,1:T}\}$ describes a physically plausible plan for achieving
 499 subgoal w_i from observation $x_{i,1}$, it is conditionally independent of goal g given subgoal w_i and
 500 observation $x_{i,1}$. Furthermore, we assume that an action $a_{i,t}$ can be recovered from observation at
 501 the same timestep $x_{i,t}$ and the next timestep $x_{i,t+1}$. Thus, we can write the factorization as

$$p_\Theta(W, \{\tau_x^i\}, \{\tau_a^i\} | g, x_{1,1}) = \left(\prod_{i=1}^N p_\theta(w_i | g) \right) \left(\prod_{i=1}^N p_\phi(\tau_x^i | w_i, x_{i,1}) \right) \left(\prod_{i=1}^N \prod_{t=1}^{T-1} p_\psi(a_{i,t} | x_{i,t}, x_{i,t+1}) \right)$$

502 C Background and Architecture

503 C.1 Task Planning

504 **Background on Density Ratio Estimation.** Let p and q be two densities, such that q is absolutely
 505 continuous with respect to p , denoted as $q \ll p$ i.e. $q(\mathbf{x}) > 0$ wherever $p(\mathbf{x}) > 0$. Then, their ratio
 506 is defined as $r(\mathbf{x}) = p(\mathbf{x})/q(\mathbf{x})$ over the support of p . We can estimate this density ratio $r(\mathbf{x})$ by
 507 training a binary classifier to distinguish between samples from p and q [46, 14, 17]. More recent
 508 work [45] has shown one can introduce auxiliary densities $\{m_i\}_{i=1}^M$ and train a multi-class classifier
 509 to distinguish samples between M classes to learn a better-calibrated and more accurate density ratio
 510 estimator. Once trained, the log density ratio can be estimated by $\log r(\mathbf{x}) = \hat{h}_p(\mathbf{x}) - \hat{h}_q(\mathbf{x})$, where
 511 $\hat{h}_i(\mathbf{x})$ is the unnormalized log probability of the input sample under the i^{th} density, parameterized by
 512 the model.

513 **Learning a Classifier to Visually Ground Task Planning.** We estimate the density ratio $\frac{p(x_{i,1}|w_i,g)}{p(x_{i,1}|g)}$
 514 with a multi-class classifier $f_\phi(x_{i,1}, \{w_j\}, g)$ trained to distinguish samples amongst the conditional
 515 distributions $p(x_{i,1}|w_i, g), \dots, p(x_{i,1}|w_M, g)$ and the marginal distribution $p(x_{i,1}|g)$. Upon conver-
 516 gence, the classifier learns to assign high scores to $(x_{i,1}, w_i, g)$ if w_i is the subgoal corresponding to
 517 the observation $x_{i,1}$ and task g and low scores otherwise.

518 **Architecture.** We parameterize f_ϕ as a 4-layer multi-layer perceptron (MLP) on top of an ImageNet-
 519 pretrained vision encoder (ResNet-18 [15]) and a frozen pretrained language encoder (Flan-T5-
 520 Base [8]). The vision encoder encodes the observation $x_{i,1}$, and the text encoder encodes the subgoals
 521 w_j and the goal g . The encoded observation, the encoded subgoals, and the encoded goal are
 522 concatenated, and passed through a MLP with 3 hidden layers of sizes 512, 256, and 128. The
 523 output dimension for MLP (i.e., number of classes for multi-classification) M is 6 for paint-block
 524 environment, 5 for object-arrange environment and 4 for kitchen-tasks environment.

525 **Choice of Large Language Model.** We use GPT3.5-turbo [37] as our large language model.

526 C.2 Visual Planning

527 **Background.** Diffusion Probabilistic Models [44, 19] learn the data distribution $h(\mathbf{x})$ from a
 528 dataset $\mathcal{D} := \{\mathbf{x}^i\}$. The data-generating procedure involves a predefined forward noising process
 529 $q(\mathbf{x}_{k+1}|\mathbf{x}_k)$ and a trainable reverse process $p_\phi(\mathbf{x}_{k-1}|\mathbf{x}_k)$, both parameterized as conditional Gaussian
 530 distributions. Here, $\mathbf{x}_0 := \mathbf{x}$ is a sample, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{K-1}$ are the latents, and $\mathbf{x}_K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ for a
 531 sufficiently large K . Starting with Gaussian noise, samples are then iteratively generated through a
 532 series of “denoising” steps. Although a tractable variational lower-bound on $\log p_\phi$ can be optimized
 533 to train diffusion models, Ho et al. [19] propose a simplified surrogate loss:

$$\mathcal{L}_{\text{denoise}}(\theta) := \mathbb{E}_{k \sim [1, K], \mathbf{x}_0 \sim h, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\|\epsilon - \epsilon_\phi(\mathbf{x}_k, k)\|^2]$$

534 The predicted noise $\epsilon_\theta(\mathbf{x}_k, k)$, parameterized with a deep neural network, estimates the noise $\epsilon \sim$
 535 $\mathcal{N}(\mathbf{0}, \mathbf{I})$ added to the dataset sample \mathbf{x}_0 to produce noisy \mathbf{x}_k .

536 **Guiding Diffusion Models with Text.** Diffusion models are most notable for synthesizing high-
 537 quality images [41, 35] and videos [20, 53] from text descriptions. Modeling the conditional data
 538 distribution $q(\mathbf{x}|\mathbf{y})$ makes it possible to generate samples satisfying the text description \mathbf{y} . To enable
 539 conditional data generation with diffusion, Ho and Salimans [18] modified the original training setup
 540 to learn both a conditional $\epsilon_\phi(\mathbf{x}_k, \mathbf{y}, k)$ and an unconditional $\epsilon_\phi(\mathbf{x}_k, k)$ model for the noise. The
 541 unconditional noise is represented, in practice, as the conditional noise $\epsilon_\phi(\mathbf{x}_k, \emptyset, k)$, where a dummy

542 value \emptyset takes the place of \mathbf{y} . The perturbed noise $\epsilon_\phi(\mathbf{x}_k, \emptyset, k) + \omega(\epsilon_\phi(\mathbf{x}_k, \mathbf{y}, k) - \epsilon_\phi(\mathbf{x}_k, \emptyset, k))$ (i.e.
 543 *classifier-free guidance*) is used to later generate samples.

544 **Video Diffusion in Latent Space.** As diffusion models generally perform denoising in the input
 545 space [19], the optimization and inference become computationally demanding when dealing with
 546 high-dimensional data, such as videos. Inspired by recent works [40, 53], we first use an autoencoder
 547 v_{enc} to learn a latent space for our video data. It projects an observation trajectory τ_x (i.e., video) into
 548 a 2D tri-plane representation [53] $\tau_z = [\tau_z^T, \tau_z^H, \tau_z^W]$ where $\tau_z^T, \tau_z^H, \tau_z^W$ capture variations in the
 549 video across time, height, and width respectively. We then diffuse over this learned latent space [53].

550 **Latent Space Video Diffusion for Visual Planning.** Our video diffusion model $p_\phi(\tau_x^i | w_i, x_{i,1})$
 551 generates video τ_x^i given a language subgoal w_i and the current observation $x_{i,1}$. It is param-
 552 eterized through its noise model $\epsilon_\phi((\tau_z^i)_k, w_i, x_{i,1}, k) := \epsilon_\phi((\tau_z^i)_k, l_{\text{enc}}(w_i), v_{\text{enc}}(x_{i,1}), k)$ where
 553 $\tau_z^i := v_{\text{enc}}(\tau_x^i)$ is the latent representation of video τ_x^i over which we diffuse. We condition the noise
 554 model ϵ_ϕ on subgoal w_i using a pretrained language encoder l_{enc} and on current observation $x_{i,1}$
 555 using video encoder v_{enc} . To use v_{enc} with a single observation $x_{i,1}$, we first tile the observation along
 556 the temporal dimension to create a video.

557 **Architecture.** We now detail the architectures of different components:

558 • **Video Autoencoder** We borrow our architecture for v_{enc} from PVDM [53] which uses transformers
 559 to project video $\tau_x \in \mathbb{R}^{T \times H \times W}$ to latent codes $\tau_z = [\tau_z^T, \tau_z^H, \tau_z^W]$ where $\tau_z^T \in \mathbb{R}^{C \times H' \times W'}$,
 560 $\tau_z^H \in \mathbb{R}^{C \times T \times W'}$, $\tau_z^W \in \mathbb{R}^{C \times H' \times T}$. Here, $T = 50$ represents the time horizon of a video, $H = 48$
 561 represents video height, $W = 64$ represents video width, $C = 4$ represents latent codebook
 562 dimension, $H' = 12$ represents latent height, and $W' = 8$ represents latent width.

563 • **Language Encoder** We use Flan-T5-Base [8] as the pretrained frozen language encoder l_{enc} .

564 • **Noise Model** We borrow PVDM-L architecture [53] which uses 2D UNet architecture, similar to
 565 the one in Latent Diffusion Model (LDM) [40], to represent $p(\tau_z | \tau_z')$. In our case, $\tau_z = v_{\text{enc}}(\tau_x^i)$
 566 and $\tau_z' = v_{\text{enc}}(x_{i,1})$. To further condition noise model ϵ_ϕ on $l_{\text{enc}}(w_i)$, we augment the 2D UNet
 567 Model with cross-attention mechanism borrowed by LDM [40].

568 For implementing these architectures, we used the codebase <https://github.com/sihyun-yu/PVDM>
 569 which contains the code for PVDM and LDM.

570 **Classifier for Consistency between Visual Planning and Action Planning.** To ensure consistency
 571 between visual planning and action planning, we want to sample observation trajectories that maxi-
 572 mizes both conditional observation trajectory likelihood from diffusion and the likelihood of sampled
 573 actions given the observation trajectory (see equation 4). To approximate likelihood calculation of
 574 action trajectory, we learn a binary classifier g_ψ that models if the observation trajectory leads to
 575 a high likelihood action trajectories. Since diffusion happens in latent space and we use gradients
 576 from g_ψ to bias the denoising of the video diffusion, $g_\psi(\tau_z^i)$ takes the observation trajectory in
 577 latent space. The binary classifier g_ψ is trained to distinguish between observation trajectories in
 578 latent space sampled from video dataset $\tau_z^i = v_{\text{enc}}(\tau_x^i), \tau_x^i \sim \mathcal{D}_{\text{video}}$ (i.e. label of 1) and observa-
 579 tion trajectories in latent space sampled from video dataset whose frames were randomly shuffled
 580 $(\tau_z^i)' = v_{\text{enc}}(\sigma(\tau_x^i)), \tau_x^i \sim \mathcal{D}_{\text{video}}$ (i.e. label of 0). Here, σ denotes the random shuffling of frames.
 581 To randomly shuffle frames in an observation trajectory (of length 50), we first randomly select 5
 582 frames in the observation trajectory. For each of the selected frame, we randomly permute it with its
 583 neighboring frame (i.e. either with the frame before it or with the frame after it). Once g_ψ is trained,
 584 we use it to bias the denoising of the video diffusion

$$\hat{\epsilon} := \epsilon_\phi((\tau_z)_k, v_{\text{enc}}(x_t), k) + \omega(\epsilon_\phi((\tau_z)_k, v_{\text{enc}}(x_t), l_{\text{enc}}(w), k) - \epsilon_\phi((\tau_z)_k, v_{\text{enc}}(x_t), k)) \\ - \omega' \nabla_{(\tau_z)_k} \log g_\psi(1 | (\tau_z)_k)$$

585 Here, $\hat{\epsilon}$ is the noise used in denoising of the video diffusion and ω, ω' are guidance hyperparameters.

586 **Classifier Architecture.** The classifier $g_\psi(\tau_z = [\tau_z^T, \tau_z^H, \tau_z^W])$ has a ResNet-9 encoder that converts
 587 τ_z^T, τ_z^H , and τ_z^W to latent vectors, then concatenate those latent vectors and passes the concatenated
 588 vector through an MLP with 2 hidden layers of sizes 256 and 128 and an output layer of size 1.

589 **C.3 Action Planning**

590 To do action planning, we learn an inverse dynamics model to $p_\psi(a_{i,t}|x_{i,t}, x_{i,t+1})$ predicts 7-
 591 dimensional robot states $s_{i,t} = p_\psi(x_{i,t})$ and $s_{i,t+1} = p_\psi(x_{i,t+1})$. The first 6 dimensions of the robot
 592 state represent joint angles and the last dimension of the robot state represents the gripper state (i.e.,
 593 whether it's open or closed). The first 6 action dimension is represented as joint angle difference
 594 $a_{i,t}[:6] = s_{i,t+1}[:6] - s_{i,t}[:6]$ while the last action dimension is gripper state of next timestep
 595 $a_{i,t}[-1] = s_{i,t+1}[-1]$.

596 **Architecture.** We use ViT-B [9] (VC-1 [32] initialization) along with a linear layer to parameterize
 597 p_ψ . ViT-B projects the observation $x_{i,t} \in \mathbb{R}^{48 \times 64 \times 3}$ to 768 dimensional latent vector from which the
 598 linear layer predicts the 7 dimensional state $s_{i,t}$.

599 **D Training and Evaluation**

600 **D.1 Task Planning**

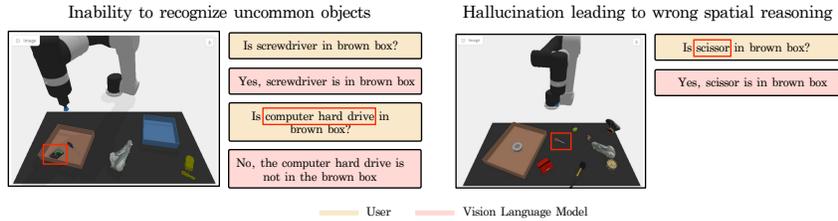


Figure 8: **Failure in VLM.** to recognize uncommon objects like computer hard drives and occasional hallucination of object presence, leading to incorrect visual reasoning.

601 **Training Objective and Dataset for Learned Classifier.** We use a softmax cross-entropy loss
 602 to train the multi-class classifier $f_\phi(x_{i,1}, \{w_j\}_{j=1}^M, g)$ to classify an observation $x_{i,1}$ into one of
 603 the M given subgoal. We train it using the classification dataset $\mathcal{D}_{\text{classify}} := \{x_{i,1}, g, \{w_j\}_{j=1}^M, i\}$
 604 consisting of observation $x_{i,1}$, goal g , candidate subgoals $\{w_j\}_{j=1}^M$ and the correct subgoal label
 605 i . The classification dataset for paint-block, object-arrange, and kitchen-tasks consists of
 606 58k, 82k and 50k datapoints respectively.

607 **Vision-Language Model (VLM) as a Classifier.** We use a frozen pretrained Vision-Language Model
 608 (VLM) (MiniGPT4 [58]) as a classifier. We first sample a list of all possible subgoals $W = \{w_i\}_{i=1}^M$
 609 from the LLM given the language goal g . We then use the VLM to eliminate subgoals from W
 610 that have been completed. For each subgoal, we question the VLM whether that subgoal has been
 611 completed. For example, consider the subgoal "Place white block in yellow bowl". To see if the
 612 subgoal has been completed, we ask the VLM "Is there a block in yellow bowl?". Consider the
 613 subgoal "Place green block in brown box" as another example. To see if the subgoal has been
 614 completed, we ask the VLM "Is there a green block in brown box?". Furthermore, if the VLM
 615 says "yes" and the subgoal has been completed, we also remove other subgoals from W that should
 616 have been completed, such as "Place white block in green bowl". Once we have eliminated the
 617 completed subgoals, we use the domain knowledge to determine which subgoal to execute out of
 618 all the remaining subgoals. As an example, if the goal is to "Stack green block on top of blue block
 619 in brown box" and we have a green block in green bowl and a blue block in blue bowl, we should
 620 execute the subgoal "Place blue block in brown box" before the subgoal "Place green block on blue
 621 block". While this process of asking questions from VLM to determine the remaining subgoals and
 622 then sequencing the remaining subgoals doesn't require any training data, it heavily relies on the
 623 task's domain knowledge.

624 **Failure Modes of VLM.** We observe two common failure modes of the VLM approach in
 625 object-arrange environment and visualize them in Figure 8. First, because the model is not
 626 trained on any in-domain data, it often fails to recognize uncommon objects, such as computer hard
 627 drives, in the observations. Second, it occasionally hallucinates the presence of objects at certain
 628 locations and thus leads to incorrect visual reasoning.

629 **VLM as a Subgoal Predictor.** We also tried to prompt the VLM with 5 examples of goal g
630 and subgoal candidates $\{w_i\}_{i=1}^M$ and then directly use it to generate the next subgoal w_i given the
631 observation $x_{i,1}$ and the goal g . However, it completely failed. We hypothesize that the VLM fails to
632 directly generate the next subgoal due to its inability to perform in-context learning.

633 **Evaluation.** We evaluate the trained classifier f_ϕ and the frozen VLM for subgoal prediction accuracy
634 on 5k unseen datapoints, consisting of observation, goal, candidate subgoals and correct subgoal,
635 generated from test tasks $\mathcal{T}_{\text{test}}$. We average over 4 seeds and show the results in Figure 7.

636 D.2 Visual Planning

637 **Ego4D dataset processing.** We pre-train on *canonical clips* of the Ego4D dataset which are text-
638 annotated short clips made from longer videos. We further divide each canonical clip into 10sec
639 segments from which we derive 50 frames. We resize each frame to 48×64 . We create a pretraining
640 Ego4D dataset of (approximately) 344k short clips, each consisting of 50 frames and a text annotation.
641 We use the loader from R3M [34] codebase (<https://github.com/facebookresearch/r3m>) to load our
642 pretraining Ego4D dataset.

643 **Training Objective and Dataset.** We use pixel-level L1 reconstruction and negative perceptual
644 similarity for training the autoencoder v_{enc} . We borrow this objective from PVDM [53] paper
645 except we don't use adversarial loss. We keep the language encoder frozen. We use denoising loss
646 in video latent space $\mathbb{E}_{k \sim [1, K], \tau_z, w, x \sim \mathcal{D}, \epsilon \sim \mathcal{N}(0, I)} [\|\epsilon - \epsilon_\phi((\tau_z)_k, l_{\text{enc}}(w), v_{\text{enc}}(x), k)\|^2]$ to train the
647 noise model ϵ_ϕ . We replace w with a null token so that ϵ_ϕ learns both a text-conditional model and
648 an unconditional model. We pretrain the autoencoder v_{enc} and the noise model ϵ_ϕ on the processed
649 Ego4D dataset. We then finetune it on our dataset $\mathcal{D}_{\text{video}} := \{\tau_x^i, w_i\}$ consisting of approximately
650 100k observation trajectories of length $T = 50$ and associated text subgoals.

651 **Classifier Training Objective and Dataset.** We use a binary cross-entropy loss to train the binary
652 classifier $g_\psi(\tau_z)$ that predicts if the observation trajectory in latent space $\tau_z = v_{\text{enc}}(\tau_x)$ leads to
653 high-likelihood action trajectory. It is trained using trajectories from video dataset $\tau_x \sim \mathcal{D}_{\text{video}}$.

654 D.3 Action Planning

655 **Training Objective and Dataset.** We train inverse dynamics p_ψ on a dataset \mathcal{D}_{inv} . Since actions are
656 differences between robotic joint states, we train p_ψ to directly predict robotic state $s_{i,t} = p_\psi(x_{i,t})$
657 by minimizing the mean squared error between the predicted robotic state and ground truth robotic
658 state. Hence, $\mathcal{D}_{\text{inv}} := \{\tau_x^i, \tau_s^i\}$ consists of 1k paired observation and robotic state trajectories, each
659 having a length of $T = 50$, in `paint-block` and `object-arrange` domains. In `kitchen-tasks`
660 domain, it consists of 3.5k paired observation and robotic state trajectories, each having a length of
661 $T = 50$.

662 **Evaluation.** We evaluate the trained p_ψ (i.e., VC-1 initialized model and other related models in
663 Figure 6) on 100 unseen paired observation and robotic state trajectories generated from test tasks
664 $\mathcal{T}_{\text{test}}$. We use mean squared error to evaluate our inverse dynamics models. We use 4 seeds to calculate
665 the standard error, represented by the shaded area in Figure 6.

666 D.4 Hyperparameters

667 **Task Planning.** We train f_ϕ for 50 epochs using AdamW optimizer [31], a batch size of 256, a
668 learning rate of $1e - 3$ and a weight decay of $1e - 6$. We used one V100 Nvidia GPU for training the
669 multi-class classifier.

670 **Visual Planning.** We borrow our hyperparameters for training video diffusion from the PVDM
671 paper [53]. We use AdamW optimizer [31], a batch size of 24 and a learning rate of $1e - 4$ for
672 training the autoencoder. We use AdamW optimizer, a batch size of 64, and a learning rate of $1e - 4$
673 for training the noise model. During the pretraining phase with the Ego4D dataset, we train the
674 autoencoder for 5 epochs and then the noise model for 5 epochs. During the finetuning phase with
675 $\mathcal{D}_{\text{video}}$, we train the autoencoder for 10 epochs and then the noise model for 40 epochs. We used two
676 A6000 Nvidia GPUs for training these diffusion models. We train g_ψ for 10 epochs using AdamW
677 optimizer, a batch size of 256 and a learning rate of $1e - 4$. We used one V100 Nvidia GPU for
678 training the binary classifier. During classifier-free guidance, we use $\omega = 4$ and $\omega' = 1$.

679 **Action Planning.** We train VC-1 initialized inverse dynamics model for 20 epochs with AdamW
680 optimizer [31], a batch size of 256 and a learning rate of $3e - 5$. We trained other randomly initialized
681 ViT-B inverse dynamics models and randomly initialized ResNet-18 inverse dynamics models for 20
682 epochs with AdamW optimizer, a batch size of 256, and a learning rate of $1e - 4$. We used one V100
683 Nvidia GPU for training these inverse dynamics models.

684 E Implementation Details for Gato and SayCan

685 In training the visual and action planning for HiP, we use 100k robot videos for visual planner and
686 the inverse dynamics, when trained from scratch, utilizes 10k state-action trajectory pairs. In order to
687 ensure fair comparison, we use 110k datapoints for training Gato [39] and SayCan [2].

688 **Gato.** We borrow the Gato [39] architecture from [Vima Codebase](#) and use it for training a language
689 conditioned policy with imitation learning. We use 110k (language, observation trajectory, action
690 trajectory) datapoints in each of the three domains for training Gato. Furthermore, we provide oracle
691 subgoals to Gato.

692 **SayCan.** We borrow the SayCan [2] algorithm from the [SayCan codebase](#) and adapt it to our
693 settings. Following the recommendations of SayCan codebase, we use CLIPort policies as primitives.
694 CLIPort policies take in top-down RGBD view and outputs pick and place pixel coordinates. Then,
695 an underlying motion planner picks the object from the specified pick-coordinate and places the
696 object at the specified place-coordinate. We train CLIPort policies on 110k (language, observation,
697 action) datapoints in `paint-block` and `object-arrange` domain. The SayCan paper uses value
698 function as an affordance function to select the correct subgoal given current observation and high
699 level goal. However, CLIPort policies don't have a value function. The SayCan codebase uses a
700 hardcoded scoring function which doesn't apply to `object-arrange` domain. To overcome these
701 issues, we use the LLM grounding strategy from [?]. It uses unnormalized logits over the pixel
702 space given by CLIPort policies as affordance and uses it to ground LLM to current observation and
703 thus predict the subgoal. We then compare SayCan with HiP and other baselines on `paint-block`
704 and `object-arrange` domain in Table 1. While SayCan outperforms other baselines, HiP still
705 outperforms it both on seen and unseen tasks of `paint-block` and `object-arrange` domain. We
706 couldn't run SayCan on `kitchen-tasks` domain as there's no clear-cut primitive in that domain.
707 This points to a limitation of SayCan which requires tasks to be expressed in terms of primitives with
708 each primitive paired with an affordance function.

709 F Additional Ablation Studies

710 F.1 Consistency between task planning and visual planning

711 To make task planning consistent with visual planning, we need to select subgoal w_i^* which maximizes
712 the joint likelihood (see equation 2) of LLM $p_{\text{LLM}}(w_i|g)$ and video diffusion $p_{\phi}(\tau_x^i|w_i, x_{i,1})$. While
713 generating videos for different subgoal candidates w_i and calculating the likelihood of the generated
714 video is computationally expensive, we would still like to evaluate its performance in subgoal
715 prediction given it is theoretically grounded. To this end, we first sample M subgoals $W = \{w_j\}_{j=1}^M$
716 from the LLM. Then, we calculate $w_i^* = \arg \max_{w \in W} \log p_{\phi}(\tau_x^i|w, x_{i,1})$ and use w_i^* as our predicted
717 subgoal. Since $\log p_{\phi}(\tau_x^i|w, x_{i,1})$ is intractable, we estimate its variational lower-bound as an
718 approximation. We use this approach for subgoal prediction in `paint-block` environment and
719 compare its performance to that of the learned classifier. It achieves a subgoal prediction accuracy of
720 $54.3 \pm 7.2\%$ whereas the learned classifier achieves a subgoal prediction accuracy of $98.2 \pm 1.5\%$
721 in `paint-block` environment. Both approaches outperform the approach of randomly selecting
722 a subgoal from W (i.e., no task plan refinement), which yields a subgoal prediction accuracy of
723 16.67% given $M = 6$. The poor performance of the described approach could result from the fact
724 that the diffusion model only coarsely approximates the true distribution $p(\tau_x^i|w_i, x_{i,1})$, which results
725 in loose variational lower-bound and thus uncalibrated likelihoods from the diffusion model. A larger
726 diffusion model could better approximate $p(\tau_x^i|w_i, x_{i,1})$, resulting in tighter variational lower-bound
727 and better-calibrated likelihoods.

728 **F.2 Consistency between visual planning and action planning**

729 To make visual planning consistent with action planning, we need to select observation trajectory
 730 $(\tau_x^i)^*$ which maximizes joint likelihood (see equation 4) of conditional video diffusion $p_\phi(\tau_x^i|w_i, x_{i,1})$
 731 and inverse model $\prod_{t=1}^{T-1} p_\psi(a_{i,t}|x_{i,t}, x_{i,t+1})$. While sampling action trajectories and calculating
 732 their likelihoods during every step of the denoising process is computationally inefficient, we would
 733 still like to evaluate its effectiveness in visual plan refinements. However, we perform video diffusion
 734 in latent space while our inverse model is in observation space. Hence, for purpose of this experiment,
 735 we learn another inverse model $\bar{p}_\psi(\tau_a^i|\tau_z^i)$ that uses a sequence model (i.e. a transformer) to produce
 736 an action trajectory τ_a^i given an observation trajectory in latent space τ_z^i . We train \bar{p}_ψ for 20
 737 epochs on 10k paired observation and action trajectories, each having a length of $T = 50$. We use
 738 AdamW optimizer, a batch size of 256 and a learning rate of $1e - 4$ during training. To generate an
 739 observation trajectory that maximizes the joint likelihood, we first sample 30 observation trajectories
 740 from video diffusion $p_\phi(\tau_x^i|w_i, x_{i,1})$ conditioned on subgoal w_i and observation $x_{i,1}$. For each
 741 generated observation trajectory τ_x^i , we sample a corresponding action trajectory τ_a^i and calculate
 742 its corresponding log-likelihood $\log \bar{p}_\psi(\tau_a^i|v_{\text{enc}}(\tau_x^i))$. We select the observation trajectory τ_x^i with
 743 highest log-likelihood. Note that we only use \bar{p}_ψ for visual plan refinement and use p_ψ for action
 744 execution to ensure fair comparison. If we use this approach for visual plan refinement with HiP, we
 745 obtain a success rate of 72.5 ± 1.9 on unseen tasks in paint-block environment. This is comparable
 746 to the performance of HiP with visual plan refinements from learned classifier g_ψ which obtains a
 747 success rate of 72.8 ± 1.7 on unseen tasks in paint-block environment. In contrast, HiP without
 748 any visual plan refinement obtains a success rate of 71.1 ± 1.3 on unseen tasks in paint-block
 749 environment. These results show that g_ψ serves as a good approximation for estimating whether an
 750 observation trajectory leads to a high-likelihood action trajectory, while still being computationally
 751 efficient.

752 **Architecture for \bar{p}_ψ .** We use a transformer model to represent $\bar{p}_\psi(\tau_a|\tau_z = [\tau_z^T, \tau_z^H, \tau_z^W])$. We first
 753 use a ResNet-9 encoder to convert τ_z^T , τ_z^H , and τ_z^W to latent vectors. We then concatenate those
 754 latent vectors and project the resulting vector to a hidden space of 64 dimension using a linear layer.
 755 We then pass the 64 dimensional vector to a trajectory transformer model [23] which generates an
 756 action trajectory τ_a of length 50. The trajectory transformer uses a transformer architecture with 4
 757 layers and 4 self-attention heads.

758 **F.3 How granularity of subgoals affects performance of HiP ?**

759 We conduct a study in paint-block environment to analyze how granularity of subgoals affect HiP.
 760 In our current setup, a subgoal in paint-block domain is of form "Place <block color> block
 761 in/on/to <final block location>" and involves a pick and a place operation. We refer to our
 762 current setup as HiP (standard). We introduce two additional level of subgoal granularity:

- 763 • *Only one pick or place operation:* The subgoal will be of form "Pick <block color> block
 764 in/on <initial block location>" or "Place <block color> block in/on/to <final
 765 block location>". It will involve either one pick or one place operation. We refer to the
 766 model trained in this setup as HiP (more granular).
- 767 • *Two pick and place operations:* The subgoal will be of form "Place <1st block color>
 768 block in/on/to <final 1st block location> and Place <2nd block color> block
 769 in/on/to <final 2nd block location>". It will involve two pick and place operations.
 770 We refer to the model trained in this setup as HiP (less granular).

Model	HiP (more granular)	HiP (Standard)	HiP (less granular)	UniPi
Paint-block (Seen)	74.5 ± 1.8	74.3 ± 1.9	61.8 ± 3.1	37.2 ± 3.8
Paint-block (Unseen)	73.1 ± 2.1	72.8 ± 1.7	58.2 ± 3.4	35.3 ± 3.2

Table 2: **Granularity of Subgoals.** Performance of HiP as we vary the granularity of subgoals. Initially, it doesn't get affected but then starts to deteriorate when subgoals become too coarse.

771 Note that UniPi has the least granularity in terms of subgoals as it tries to imagine the entire trajectory
 772 from goal description. Table 2 in the rebuttal document compares HiP (standard), HiP (more granular),

773 HiP (less granular) and UniPi on seen and unseen tasks in paint-block environment. We observe
 774 that HiP (standard) and HiP (more granular) have similar success rates where HiP (less granular)
 775 has a lower success rate. UniPi has the lowest success rate amongst these variants. We hypothesize
 776 that success rate of HiP remains intact when we decrease the subgoal granularity as long as the
 777 performance of visual planner doesn’t degrade. Hence, HiP (standard) and HiP (more granular)
 778 have similar success rates. However, when the performance of visual planner degrades as we further
 779 decrease the subgoal granularity, we see a decline in success rate as well. That’s why HiP (less
 780 granular) sees a decline in success rate and UniPi has the lowest success rate amongst all variants.

781 **F.4 Ablation on Visual Planning Model**

Model	Paint-block		Object-arrange		Kitchen-tasks	
	Seen	Unseen	Seen	Unseen	Seen	Unseen
HiP (RSSM)	70.2 ± 2.4	69.5 ± 1.6	59.6 ± 3.8	59.2 ± 3.9	50.6 ± 16.2	46.8 ± 19.4
HiP	74.3 ± 1.9	72.8 ± 1.7	75 ± 2.8	75.4 ± 2.6	85.8 ± 9.4	83.5 ± 10.2

Table 3: **Ablating Visual Planner.** While performance gap between HiP and HiP (RSSM) is small in Paint-block, it widens in more visually complex domains, such as Object-arrange and Kitchen-tasks, thereby showing the importance of video diffusion model.

782 To show the benefits of video diffusion model, we perform an ablation where we use (text-conditioned)
 783 recurrent state space model (RSSM), taken from DreamerV3 [?], as visual model for HiP. We borrow
 784 the RSSM code from [dreamerv3-torch codebase](#). To adapt RSSM to our setting, we condition RSSM
 785 on subgoal (i.e. subgoal encoded into a latent representation by Flan-T5-Base) instead of actions.
 786 Hence, sequence model of RSSM becomes $h_t = f(h_{t-1}, z_{t-1}, w)$ where w is latent representation of
 787 subgoal. Furthermore, we don’t predict any reward since we aren’t in a reinforcement learning setting
 788 and don’t predict continue vector since we decode for a fixed number of steps. Hence, we remove
 789 reward prediction and continue prediction from the prediction loss. To make the comparisons fair,
 790 we pretrain RSSM with Ego4D data as well. We report the results in Table 3. We see that HiP with
 791 video diffusion model outperforms HiP with RSSM in all the three domains. While the performance
 792 gap between HiP(RSSM) and HiP (i.e. using video diffusion) is small in paint-block environment,
 793 it widens in object-arrange and kitchen-tasks domains as the domains become more visually
 794 complex.

795 **F.5 Analyzing sensitivity of iterative refinement to hyperparameters**

HiP	$\omega' = 0.5$	$\omega' = 0.75$	$\omega' = 1.0$	$\omega' = 1.25$	$\omega' = 1.5$	$\omega' = 1.75$	$\omega' = 2.0$
Paint-block (Seen)	71.8 ± 2.3	72.3 ± 2.0	74.3 ± 1.9	73.9 ± 2.2	72.1 ± 1.7	70.4 ± 2.4	68.2 ± 1.9
Paint-block (Unseen)	71.1 ± 2.5	71.4 ± 1.8	72.8 ± 1.7	73.1 ± 1.5	71.4 ± 1.5	69.3 ± 2.7	66.8 ± 1.4

Table 4: **Sensitivity of visual iterative refinement to guidance scale.** Performance of HiP as we vary the guidance scale ω' . HiP performs best when $\omega' \in \{1, 1.25\}$ but performance degrades for higher values of ω' .

796 The subgoal classifier doesn’t introduce any test time hyperparameters and we use standard hyper-
 797 parameters ($1e - 3$ learning rate, $1e - 6$ weight decay, 256 batch size, 50 epochs, Adam optimizer)
 798 for its training which remains fixed across all domains. We observed that the performance changes
 799 minimally across different hyperparameters, given a learning rate decay over training. However,
 800 the observation trajectory classifier g_ψ introduces an additional test time hyperparameter ω' which
 801 appropriately weights the gradient from observation trajectory classifier. Table 4 in the rebuttal
 802 document varies ω' between 0.5 and 2 in intervals of 0.25 and shows success rate of HiP. We see that
 803 HiP gives the best performance when $\omega' \in \{1, 1.25\}$ but it’s performance degrades for higher values
 804 of ω' .

805 **G Analyzing Runtime of HiP**

806 We provide average runtime of HiP for a single episode in all the three domains in Table 5 of the
 807 rebuttal document. We average across 1000 seen tasks in each domain. We break the average runtime

Domain	Subgoal candidate generation	Subgoal classification	Visual planning per subgoal	Action planning per subgoal	Action execution per subgoal	Episodic runtime
Paint-block	1.85s	0.41s	7.32s	0.91s	6.35s	80.61
Object-arrange	1.9s	0.43s	7.39s	0.89s	9.57s	78.71
Kitchen-tasks	1.81s	0.41s	7.35s	0.98s	1.28s	40.37

Table 5: **Run time of HiP.** Average episodic run-time of HiP, along with average run-time of its different components for `Paint-block`, `Object-arrange` and `Kitchen-tasks` domains. While HiP has similar planning times across different domains, it has different action execution times and episodic runtimes across domains due to differences in simulation properties and average number of subgoals.

808 by different components: task planning (subgoal candidate generation and subgoal classification),
809 visual planning, action planning and action execution. We execute the action plan for a subgoal in
810 open-loop and then get observation from the environment for deciding the next subgoal. From Table 5,
811 we see that majority of the planning time is taken by visual planning. Recent works [24, 42, 56] have
812 proposed techniques to reduce sampling time in diffusion models, which can be incorporated into our
813 framework for improving visual planning speed in the future.

814 H Limitations and Conclusion

815 **Limitations.** Our approach has several limitations. As high-quality foundation models for visual
816 sequence prediction and robot action generation do not exist yet, our approach relies on smaller-scale
817 models that we directly train. Once high-quality video foundation models are available, we can use
818 them to guide our smaller-scale video models [51] which would reduce the data requirements of our
819 smaller-scale video models. Furthermore, our method uses approximations to sample from the joint
820 distribution between all the model. An interesting avenue for future work is to explore more efficient
821 and accurate methods to ensure consistent samples from the joint distribution.

822 **Conclusion.** In this paper, we have presented an approach to combine many different foundation
823 models into a consistent hierarchical system for solving long-horizon robotics problems. Currently,
824 large pretrained models are readily available in the language domain only. Ideally, one would train a
825 foundation model for videos and ego-centric actions, which we believe will be available in the near
826 future. However, our paper focuses on leveraging separate foundation models trained on different
827 modalities of internet data, instead of training a single big foundation model for decision making.
828 Hence, for the purposes of this paper, given our computational resource limitations, we demonstrate
829 our general strategy with smaller-scale video and ego-centric action models trained in simulation,
830 which serve as proxies for larger pretrained models. We show the potential of this approach in
831 solving three long-horizon robot manipulation problem domains. Across environments with novel
832 compositions of states and goals, our method significantly outperforms the state-of-the-art approaches
833 towards solving these tasks.

834 In addition to building larger, more general-purposed visual sequence and robot control models, our
835 work suggests the possibility of further using other pretrained models in other modalities, such as
836 touch and sound, which may be jointly combined and used by our sampling approach. Overall, our
837 work paints a direction towards decision making by leveraging many different powerful pretrained
838 models in combination with a tiny bit of training data. We believe that such a system will be
839 substantially cheaper to train and will ultimately result in more capable and general-purpose decision
840 making systems.