VIDEO2POLICY: SCALING UP MANIPULATION TASKS IN SIMULATION THROUGH INTERNET VIDEOS

Anonymous authors

Paper under double-blind review

ABSTRACT

Simulation offers a promising approach for cheaply scaling training data for generalist policies. To scalably generate data from diverse and realistic tasks, existing algorithms either rely on large language models (LLMs) that may hallucinate tasks not interesting for robotics; or digital twins, which require careful real-tosim alignment and are hard to scale. To address these challenges, we introduce Video2Policy, a novel framework that leverages large amounts of internet RGB videos to reconstruct tasks based on everyday human behavior. Our approach comprises two phases: (1) task generation through object mesh reconstruction and 6D position tracking; and (2) reinforcement learning utilizing LLM-generated reward functions and iterative in-context reward reflection for the task. We demonstrate the efficacy of Video2Policy by reconstructing over 100 videos from the Something-Something-v2 (SSv2) dataset, which depicts diverse and complex human behaviors on 9 different tasks. Our method can successfully train RL policies on such tasks, including complex and challenging tasks such as throwing. Furthermore, we show that a generalist policy trained on the collected sim data generalizes effectively to new tasks and outperforms prior approaches. Finally, we show the performance of our policies improves by simply including more internet videos. We believe that the proposed Video2Policy framework is a step towards generalist policies that can execute practical robotic tasks based on everyday human behavior.

032

004

010 011

012

013

014

015

016

017

018

019

020

021

022

024

025

026

027

1 INTRODUCTION

Training generalist policies requires collecting large quantities of diverse robotic expert data. However, the approach of collecting data through teleoperation is constrained by high operation costs, while collecting data from autonomous policies can be unsafe, or result in low-quality data. Simulation offers an appealing alternative to real-world data that does not suffer from these challenges, and can be used to train general and robust policies (Hwangbo et al., 2019; Andrychowicz et al., 2020). Recent work has explored automatically generating diverse and relevant tasks in simulation as a way to create a scalable pipeline for generating robotics data (Deitke et al., 2022; Wang et al., 2023b;c; Makatura et al., 2023).

041 However, existing methods primarily rely on text-only task specification using Large Language 042 Models (LLMs), which do not have grounded robotics knowledge. They often produce tasks that 043 are not diverse, have uninteresting behavior, or use uninteresting object assets and are thus less 044 useful for training generalist policies. To better capture the real-world distributions of task behaviors and objects, we propose leveraging RGB videos from the internet to create corresponding tasks. Unlike Real2Sim approaches that construct digital twins (Hsu et al., 2023; Torne et al., 2024) for 046 a single scene, we want to train a generalist policy for multiple scenes and therefore we do not 047 require perfect reconstructions. Instead, we leverage large amounts of internet videos to capture 048 task-relevant information such as object assets and scene layouts. We then generate simulated tasks using a Vision-Language Model (VLM) that can take the video, video captions, object meshes, sizes, and 6D poses, and produce corresponding **task codes**, which can be executed to generate scenes. 051

Beyond task proposals, we require an efficient and automatic way to solve tasks. Naively applying
 reinforcement or imitation learning is challenging as it requires manual human effort for each task
 to create demonstrations or reward functions. Inspired by the recent advancements of LLMs in

069

070 071



Figure 1: The **Video2Policy** framework can leverage internet videos to generate simulation tasks and learn policies for them automatically, which can be considered a data engine for generalist policies.

code generation for various tasks (Achiam et al., 2023; Roziere et al., 2023), some researchers have
proposed automating policy learning or deployment through using an LLM to produce policy code
directly (Huang et al., 2023b; Liang et al., 2023; Wang et al., 2023b), or to produce reward function
code (Ma et al., 2023; Wang et al., 2023c). Gensim (Wang et al., 2023b) leverages this idea for
unsupervised task generation. However, predicting goals restricts it to simple tasks as it does not
account for dynamic tasks or tasks that involve complex object interactions.

In contrast, reinforcement learning (RL) is effective at solving complex tasks (Schulman et al., 2017;
Ye et al., 2021; Hafner et al., 2023; Wang et al., 2024; Springenberg et al., 2024). RoboGen (Wang et al., 2023c) leverages LLM-generated reward functions for RL. However, it is hard to scale as it requires manual success functions. However, since we leverage both text prompts and explicit visual prior knowledge from human videos, we can leverage that information for better success functions.

083 We propose Video2Policy, a framework to leverage internet RGB videos in an automated pipeline 084 to create simulated data to train generalist police, as shown in Fig. 1. It produces code for task 085 generation in simulation using VLMs and learns policies for those tasks via reinforcement learning with reward reflection. This autonomous approach allow us to easily scale up data generation by 087 using internet videos to produce a large amount of visually grounded tasks. Our framework leverages 880 both behavioral and object diversity from the videos, enabling generalization both at the object level as well as task level. To further improve generalization instead of simply cloning the trajectory from 089 the video, we apply position randomization to the initial states. Specifically, our framework consists 090 of two phases: (1) We reconstruct the object meshes involved in the tasks from videos as task assets, 091 and extract the 6D pose of each object (2) We leverage the VLM to write the task code based on 092 visual information as well as text prompts, and iteratively generate the reward function using reward 093 reflection based on training logs and evaluation results. Ultimately, we obtain the corresponding 094 refined task code along with the learned policy model that demonstrates behavior similar to the 095 input video, enabling it to generate high-quality data relevant to daily human tasks. 096

For experimental evaluation, we focus on table-top manipulation tasks with a single robot arm in IsaacGym simulator(Makoviychuk et al., 2021). To validate the effectiveness of our approach, we 098 conduct experiments on the Something-Something V2 (SSv2) dataset (Goyal et al., 2017), a human video dataset with diverse scenes and language instructions of human daily behaviors. We recon-100 struct over 100 videos on 9 distinct tasks from the dataset in total, including involving multiple 101 objects. To further evaluate more complex behaviors and object relationships, we also utilize three 102 self-recorded videos. Significantly, the results indicate the policies learned by our method signifi-103 cantly outperform those learned from previous LLM-driven methods. We achieve 82% success rates 104 for the 9 tasks on average. To demonstrate generalization capabilities, we train a general policy by 105 imitation learning from the simulation data collected from the learned policies from 50 videos of the same behavior. We find that it can achieve 75% success rates on the tasks from the 10 unseen 106 videos with the same behavior. More importantly, as the number of training tasks increases, the 107 performance of the general policy improves, which shows the scalability of our framework.

108 2 RELATED WORK

110

Real2Sim Scene Generation for Robotics Generating realistic and diverse scenes for robotics has 111 recently emerged as a significant challenge aimed at addressing data issues through simulation. 112 Some researchers have developed Real2Sim pipelines for image-to-scene or video-to-scene gener-113 ation. Certain studies (Hsu et al., 2023; Torne et al., 2024) focus on constructing digital twins that 114 facilitate the transition from the real world to simulation; however, these approaches often depend on 115 specific real-world scans or extensive human assistance. Dai et al. introduces the concept of digital 116 cousins, while Chen et al. (2024) employs inverse graphics to enhance data diversity. Neverthe-117 less, their approach to diversity primarily involves replacing various assets. The lack of task-level 118 diversity hinders the ability to capture the distribution of real-world tasks, and the constraints of 119 specific data formats complicate the scalability of robotics data. Although Real2Code (Mandi et al., 2024) aims to build simulation scenes from RGB images, it focuses on articulation parts and requires 120 in-domain code data for fine-tuning. 121

122 Scaling up Simulation Tasks Previously, researchers aimed to build simulation benchmarks to fa-123 cilitate scalable skill learning and standardized workflows (Li et al., 2023; Gu et al., 2023; Srivastava 124 et al., 2022; Nasiriany et al., 2024). Most of these benchmarks were constructed manually, making 125 them difficult to scale. Recently, some researchers have focused on text-to-scene generation, emphasizing the creation of diverse scenes. Works such as Deitke et al. (2022); Makatura et al. (2023); 126 Liang et al.; Chen et al. (2023) utilize procedural asset generation or domain randomization, while 127 Jun & Nichol (2023); Yu et al. (2023a); Poole et al. (2022) engage in text-to-3D asset generation. 128 Although these approaches can achieve asset-level or scene-level diversity in robotic tasks, they fall 129 short in delivering task-level diversity. Gensim (Wang et al., 2023b) attempts to generate rich sim-130 ulation environments and expert demonstrations using large language models (LLMs) to achieve 131 task-level diversity. However, text-based task generation tends to be arbitrary regarding object se-132 lection and their relationships, limiting its ability to represent the true distribution of tasks in the real 133 world. In contrast, our work leverages real-world RGB videos to create corresponding simulation 134 tasks that better reflect the real-world distributions of tasks and objects, facilitating easier scalability 135 due to the abundance and accessibility of internet video data.

136 **Policy Learning via LLMs** To enable automatic policy learning with high quality, researchers are 137 increasingly turning to large language models (LLMs) for assistance. Some studies (Liang et al., 138 2023; Huang et al., 2023a; Lin et al., 2023; Wang et al., 2023a) propose generating structured code 139 outputs for decision-making problems, most of which rely on predefined primitives. Other works 140 (Yu et al., 2023b; Ma et al., 2023; Wang et al., 2023c) generate reward functions using LLMs for 141 reinforcement learning. Nevertheless, Eureka (Ma et al., 2023) requires predefined success functions 142 for iterative updates of reward functions, while RoboGen (Wang et al., 2023c) selects the highest reward as the initial state for the next sub-task, which introduces noise due to the variability in 143 the generated reward functions. In contrast, our work generates success functions by leveraging 144 visual prior knowledge from the provided videos and updates the reward functions iteratively using 145 a Chain-of-Thought (CoT) approach (Wei et al., 2022). 146

147 148

3 BACKGROUND

149 150

To present our work systematically, we formulate it as a two-level hierarchical solution of Markov
Decision Processes (MDPs), corresponding to the two-phase pipeline. In contrast to some hierarchical methods targeting action levels (McGovern et al., 1998; Hauskrecht et al., 2013), we focus on
the level of reward functions.

Low-level MDP of Controlling Problem Specifically, the scene reconstruction phase is to construct the MDP $\mathcal{G} = \langle \mathcal{S}, \mathcal{A}, T, \mathcal{R}_{0|1} \rangle$ from videos, where $\mathcal{S} \in \mathbb{R}^m$ represents the states of the environment, \mathcal{A} is the action space of the agent, and T is the transition probability function. $\mathcal{R}_{0|1}$ is the 0-1 reward function that distinguish whether the trajectory is successful. $\mathcal{R}_{0|1}$ can be a scalar evaluation function. To solve this MDP, we will train a policy π through reinforcement learning in the Isaac Gym simulation. To achieve high performance, we leverage the LLM to sample various reward functions to learn policies in a high-level manner, based on the evaluation results from $\mathcal{R}_{0|1}$ and some CoT (Wei et al., 2022) instructions. 162 High-level MDP of Reward Designs The aim of reward design is to create a shaped reward function 163 that simplifies the optimization of a challenging given reward function, such as the sparse 0-1 reward 164 function $\mathcal{R}_{0|1}$. Following the definition of Reward Design Problem (RDP) from previous works 165 (Singh et al., 2009; Ma et al., 2023), we consider a high-level MDP $\hat{\mathcal{G}} = \langle \hat{\mathcal{S}}, \hat{\mathcal{A}}, \hat{T}, F \rangle$. Here, $\hat{\mathcal{A}}$ is 166 the space of reward functions. Each time we choose an action $\hat{R} \in \mathcal{A}$ in the MDP $\hat{\mathcal{G}}$, we will train a 167 policy π by RL for the low-level MDP $\mathcal{G} = \langle \mathcal{S}, \mathcal{A}, T, \hat{\mathcal{R}} + \mathcal{R}_{0|1} \rangle$. The horizon of the MDP $\hat{\mathcal{G}}$ is the 168 iteration number in the second phase. \hat{S} includes the training and evaluation information during RL and the policy model π . \hat{T} is the state transition function, and F is the reward function that produces 170 a scalar evaluation of any policy π . Specifically, F is equal to $\mathcal{R}_{0|1}$. Thus, the high-level MDP's 171 goal is to find a reward function $\hat{\mathcal{R}} \in \hat{A}$ to maximize the success rates of the low-level policies. 172

173 174

175

4 GENERATING SIMULATED TASKS AND POLICIES FROM HUMAN VIDEOS

The proposed framework, **Video2Policy**, steps further for task proposal and policy learning through internet videos, to provide diverse and realistic tasks as well as the corresponding learned policies. It consists of two phases: task scene generation and policy learning. In Sec. 4.1, we introduce the pipeline for reconstructing scenes from RGB videos. Subsequently, in Sec. 4.2, we demonstrate how to generate the code for the task and learn policies to solve it. Finally, we provide an example of training a generalist policy within our framework in Sec. 4.3.

182 183

4.1 Scene Reconstruction from Videos

Since the goal of this work is to learn policies from videos rather than automatically retarget trajectories, our scene reconstruction phase focuses on reconstructing the manipulated objects along with their relative relationships. To master the skill demonstrated in the video, we allow for random positions and orientations of each object in the initial states. Consequently, the pipeline for the scene reconstruction is illustrated in Fig. 2: (1) detect and segment the manipulated objects in the video using text captions; (2) reconstruct the object meshes from the video and estimate the actual sizes of the meshes; (3) perform 6D position tracking for each object in the video. Afterward, we obtain a JSON file that includes the video and object information, from which the task code is generated.

Object Grounding We first use Grounding DINO (Liu et al., 2023) to detect the manipulated objects in the first frame of the video based on their names. Since the SSv2 dataset (Goyal et al., 2017) provides both video captions and object labels, we use these as text prompts for detection. For self-collected videos of more challenging behaviors, we provide the video captions and object names manually. Afterward, we perform video segmentation using the SAM-2 model (Ravi et al., 2024). Specifically, we segment the objects in the first frame using bounding boxes obtained during detection and select five positive pixel points in each object mask for video segmentation. This process yields segmented videos that contain the masks of the manipulated objects.

200 Object Reconstruction With the segmentation masks of each object for each frame, we perform mesh reconstruction based on these images. Since most internet videos are recorded from 201 a single viewpoint, we leverage the InstantMesh model (Xu et al., 2024) to reconstruct the 3D 202 meshes, which supports mesh generation from a single image. Typically, we choose the first 203 frame to reconstruct the meshes; however, for those with objects that are significantly occluded 204 in the first frame, we utilize the last frame instead. Due to the training mechanism in the cur-205 rent mesh generation methods, the meshes are commonly normalized during reconstruction, lead-206 ing to mismatched sizes compared to the corresponding objects in the video. To establish a 207 more realistic size relationship between objects, we propose a simple and efficient size estima-208 tion method. We predict the camera intrinsic matrix K and the depth $d_{i,j}$ of the image $I_{i,j}$ using 209 UniDepth (Piccinelli et al., 2024), where i, j are the pixel coordinates. Given the masked region 210 M of the object, we can calculate the maximum distance for the masked region in reality D_{image} : $D_{\text{image}} = \max_{(i_1, j_1), (i_2, j_2) \in M} \|\mathbf{p}(i_1, j_1) - \mathbf{p}(i_2, j_2)\|, \text{ where } \mathbf{p}(i, j) = \mathbb{K}^{-1} \cdot [x, y, 1]^{\mathsf{T}} \cdot d_{i,j}. \text{ Here } \mathbf{p}(i_1, j_1) + \mathbf{p}(i_2, j_2) \|\mathbf{p}(i_1, j_1) - \mathbf{p}(i_2, j_2)\|, \text{ where } \mathbf{p}(i_1, j_1) = \mathbb{K}^{-1} \cdot [x, y, 1]^{\mathsf{T}} \cdot d_{i,j}.$ 211 is the 3D position of each masked pixel in the camera coordinate system. We can then calculate the 212 213 maximum distance of vertices in the mesh object, denoted as D_{mesh} . The scale ratio ρ for the mesh object is defined as $\rho = D_{\text{image}}/D_{\text{mesh}}$. The absolute sizes may exhibit some noise due to errors 214 in depth estimation, intrinsic estimation, and object occlusion. However, the relative size of each 215 object is mostly accurate, as D_{image} and D_{image} are calculated within the same camera coordinate



Figure 2: **The scene generation phase of the proposed Video2Policy framework**. We generate the tasks from internet videos into the simulation, which we attempt to solve by reinforcement learning.

system. In practice, we can apply domain randomization to change the absolute sizes of each objectwhile maintaining the same relative size among them.

6D Position Tracking of Objects After reconstructing the objects, we predict the 6D position of
each object throughout the video, which will be fed into GPT-4o for code generation. We utilize
FoundationPose (Wen et al., 2024) in model-based setups to estimate the position and orientation of
the objects. This model takes the object mesh, predicted camera intrinsics, and the depth information
from each frame as inputs. Finally, we automatically generate a URDF file for each object based on
the mesh file and the calculated scaled size. Ablations for the tracking information are in App. A.3.

Compared to directly feeding the videos into the task generation process, we explicitly leverage
visual prior knowledge from the videos and provide comprehensive information about the task,
including the 3D meshes, object sizes, and 6D positions of the objects. We consolidate all the
information into a JSON file that describes the task scene. Consequently, the generated tasks can
better align with the distributions of real-world scenes and behaviors.

248 249

250

4.2 TASK CODE GENERATION AND POLICY LEARNING

After extracting the visual information from the video into a task JSON file, we can build the task 251 scene in simulation and learn policies based on GPT-40. This process occurs in two stages. First, 252 we generate the complete task code, which can be executed directly in the Isaac Gym (Makoviy-253 chuk et al., 2021) environment. Second, inspired by recent work on LLM-driven reward function 254 generation, we iteratively fine-tune the generated reward function using in-context reward reflection (Shinn et al., 2023; Ma et al., 2023; Wang et al., 2023a). In contrast to the previous work Eureka 256 (Ma et al., 2023), which is the most similar to ours, we generate the task codes, including the reward 257 functions, from scratch, rather than relying on pre-existing task codes and starting reward reflec-258 tion from manually defined success functions. We formulate this learning process as a two-level 259 hierarchical solution of MDP, defined in Sec. 3.

260 **Task Code Generation** Inspired by previous work (Ma et al., 2023; Wang et al., 2023b;c), we sys-261 tematically introduce the pipeline for general task code generation, which helps to infer codes by 262 prompting in a well-organized way. Notably, the task code consists of six parts: scene informa-263 tion, reset function, success function, observation function, observation space function, and reward 264 function. (1) The scene information refers to the task scene JSON file created from the videos. It 265 contains the task title, video file path, video description, and object information, including sizes, 266 URDF paths, and tracked 6D position lists. (2) The reset function is responsible for positioning the objects according to specific spatial relationships in the beginning. For instance, in the task "uncov-267 ering A from B," both objects A and B are randomly placed on the table by the parent reset function, 268 but object A will be positioned above object B as specified by the generated child reset function. 269 (3) The success function determines the success state. Notably, both the reset and success functions 270 are generated by GPT-40 based on the task description, the provided 6D position list, and Chain-271 of-Thoughts examples (Wei et al., 2022). (4) Furthermore, since we use reinforcement learning to 272 master the skill in simulation, we have access to the states of the objects, and the policy is based on 273 state observations. Thus, in addition to the object states, we also query GPT-40 to determine whether 274 additional observations are necessary. Interestingly, we find that it can include observations such as the distance between the object and the gripper or the normalized velocity toward the target object. 275 (5) Simultaneously, it calculates the correct observation shape for building the neural networks. (6) 276 Regarding the reward function, we follow the instructions from Ma et al. (2023) with CoT examples. We write a template for task code generation, allowing us to query the VLM just once to generate the 278 executable code encompassing all six parts. Most significantly, we generate eight example codes and 279 select one by re-checking the correctness, reasonability and efficiency from GPT-40, which is the 280 base code candidate for the subsequent in-context reward reflection stage. The generated examples 281 are demonstrated in App. A.1 and the robustness evaluation results are in App. A.2. 282

Reinforcement Learning and Reward Function Iteration As mentioned in 3, we propose to op-283 timize the policy through a two-level hierarchical solution of MDP. Given the generated task code, 284 we train policies through reinforcement learning under the reward function \mathcal{R} and success function 285 $\mathcal{R}_{0|1}$. Notably, we assign a high trade-off to the success rewards, formulating the training reward 286 function as $\mathcal{R} + \lambda \mathcal{R}_{0|1}$, $\lambda = 100$. Moreover, following the approach in Eureka (Ma et al., 2023), 287 we apply the in-context reward reflection to design the reward functions iteratively using GPT-40. 288 Each time, we sample N = 8 different reward functions for training policies and collect the train-289 ing and evaluation logs. We then select the best function from the previous iteration and generate 290 new reward functions based on these logs, along with specific instructions and CoT examples. For 291 example, in addition to providing good examples from previous tasks, we prioritize training outputs 292 where the accumulated rewards for successful trajectories exceed those for failed ones. 293

4.3 TRAINING GENERAL POLICIES IN SIMULATION

296 As witnessed by the recent success of policy learning from large-scale datasets with certain formats 297 (Padalkar et al., 2023; Reed et al., 2022; Team et al., 2024; Brohan et al., 2023), we want to inves-298 tigate how to learn a general policy from the internet videos, which directly outputs the executable 299 actions of the robot rather than the un-executable future videos (Du et al., 2024; Qin et al., 2023) or language tokens (Liang et al., 2023; Brohan et al., 2023). We consider our Video2Policy a data 300 engine to generate successful policies from internet videos. Then we can acquire expert trajectories 301 in simulation, which match the video behavior. Notably, those expert trajectories can be any format 302 we want, such as states, 2D images, or 3D images. In this work, we choose RGB image observation 303 as the universal perception format of our general policy. Therefore, we train policies from the videos 304 and collect the successful trajectories from the learned policy. Afterward, we use imitation learning 305 to learn the general policy from the collected dataset by behavior cloning. 306

- 5
- 308 309

307

294

295

EXPERIMENTS

310

In this section, we present detailed evaluations of the proposed Video2Policy framework on internet video datasets about manipulations. Specifically, the experiments are designed to answer the fol-311 lowing questions: (a) How does the generated scene from the video look like, including the objects 312 as well as the visual information, in Sec. 5.1. b) How does the policy trained under our framework 313 perform compared to the videos and the baselines, in Sec. 5.1. (c) What is the performance of the 314 general policy learned from diverse internet videos, and can it generalize to novel scenes, in Sec. 315 5.2. (d) What affects the proposed Video2Policy framework most for policy learning, in Sec. 5.3. 316

Experimental Setup We use the Issac Gym (Makoviychuk et al., 2021) as the simulation engine 317 for all the experiments, which is commonly used in robotics tasks due to the advantages of efficient 318 computing and highly realistic physics. We focus on the table manipulation tasks, and all the objects 319 will randomly reset on the table in the beginning. The horizon of all the tasks is set to 300 and 320 the parallel environments are 8192, equally. For each task, we average the success rates over 10 321 evaluation episodes across 3 runs with different seeds. 322

Video Data Source To reconstruct scenes from internet RGB videos, we choose the Something 323 Something V2 (SSv2) dataset (Goyal et al., 2017), a common and diverse video dataset for the



Figure 3: Some visualization of the tasks generated from SSv2 Video Dataset.

335 robotics community. It includes diverse behaviors concerning manipulating something with some-336 thing by human hands. To further investigate the ability of our framework on more complex objects 337 or behaviors, we record three in-the-wild videos of different behaviors by ourselves. Notably, all the 338 videos we use in the experiment are 3 channels with RGB only, with the highest accessibility. For 339 the video quality, the small motion of the camera is tolerated, and we scale the resolution to 1024.

340 Scene Generation As mentioned in Sec. 4.1, we do 6D position tracking for all the objects. Con-341 sidering that we randomize the initial states of all the objects, we only feed the 6D pos from the 342 first frame and the final frame into the prompts. It is reasonable in most tasks because those two 343 frames are significant to infer the relationship among objects. Even if this simplification will miss 344 the motion information for some behavior, e.g. throwing, we also provide the task description to 345 design the reward function so that the LLM will generate the velocity reward components. More-346 over, we explicitly calculate the difference between the 6D pos and feed the information into the 347 LLM to think about the success function. For most of the SSv2 videos of a single object, there are severe occlusions, making it difficult to reconstruct the mesh asset. We manually choose the first 348 frame or the last frame to reconstruct the mesh and predict the 6D position in the same pipeline. We 349 generate the task code in a curriculum manner (Ma et al., 2023; Wang et al., 2023b) after obtaining 350 the visual information. From the beginning, we provide the example code of reach a block 351 and grasp a block. Then we will add the successfully generated task examples into the task 352 pool for the next one. Finally, it can even learn to use the direction velocity reward for dynamic 353 tasks and resolve them. Some demos for the generated tasks are in Fig. 3. 354

Reinforcement Learning For the policy learning, we choose the PPO (Schulman et al., 2017) algo-355 rithm in a well-tuned implementation codebase (Makoviichuk & Makoviychuk, 2021; Ma et al., 356 2023). We share the same parameters recommended in the codebase across all tasks and all 357 baseliens. As for the evaluation metric, we write the ground-truth success function for each gen-358 erated task. Unlike Eureka (Ma et al., 2023), we do not allow access to the evaluation metric during 359 training, and we manually evaluate the results for the final models. For SSv2 dataset, we make 5 360 iterations and sample 8 reward functions at each iteration during reinforcement learning. In our 361 collected videos, we make 8 iterations and sample 8 reward functions. 362

364

5.1 POLICY LEARNING FROM VIDEOS

365 Policy Learning Baselines Since there is few works transforming the internet RGB videos into 366 policies for the task, we focus on the policy learning part in the experiments. We benchmark our 367 method against the following baselines. (1) Code-as-Policy (CoP) (Liang et al., 2023), which queries 368 the LLM with all the states in the environment to write the executable code for the robot. To ensure better performance of CoP, we use the close-loop control and regenerate code policies every 50 369 steps. (2) RoboGen, which does not require a success function and learns without reward reflection 370 iteration. (3) Eureka, which generates code for both the reward and the success function using an 371 LLM and does not use video information. To make fair comparisons, we use the same object meshes 372 and task codes generated from the videos for all the baselines. 373

374 Performance Analysis of Learned Policies We compare the performance of our method with the 375 above baselines on 9 videos sampled from 9 tasks generated from SSv2 dataset, as shown in Tab. 1. We find that the proposed Video2Policy method outperforms the baseline in most tasks, with 376 smaller variance across seeds. RoboGen (Wang et al., 2023c) and Eureka (Ma et al., 2023) achieve 377 comparable results to ours in the videos of a single object. However, for multiple objects, the

Table 1: Results of Learned Policies for Videos in SSv2 dataset (3 seeds). The mean \pm std of the success rates are shown in the table. Our method outperforms the other baselines to a degree and achieves smaller variance in general.

Task (Succ.)	Code-as-Policy	RoboGen	Eureka	Video2Policy
single object				
Push sth. left	0.17 ± 0.13	0.93 ± 0.05	$\textbf{1.00} \pm 0.00$	$\textbf{1.00} \pm 0.00$
Push sth. right	0.75 ± 0.12	$\textbf{1.00} \pm 0.00$	$\textbf{1.00} \pm 0.00$	$\textbf{1.00}\pm0.00$
Lift up sth.	0.33 ± 0.21	0.28 ± 0.09	0.83 ± 0.13	$\textbf{0.93} \pm 0.05$
Tip sth. over	$\textbf{1.00}\pm0.00$	0.97 ± 0.05	0.67 ± 0.47	$\textbf{1.00}\pm0.00$
multiple objects				
Cover sth. with sth.	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	$\textbf{0.07} \pm 0.05$
Uncover sth. from sth.	0.10 ± 0.08	0.67 ± 0.26	0.63 ± 0.38	$\textbf{0.97} \pm 0.05$
Push sth. with sth.	0.03 ± 0.05	0.03 ± 0.05	$\textbf{0.47} \pm 0.38$	0.43 ± 0.40
Push sth. next to sth.	0.80 ± 0.16	0.23 ± 0.05	0.83 ± 0.12	$\textbf{1.00}\pm0.00$
Drop sth. in front of sth.	0.53 ± 0.31	0.37 ± 0.09	0.87 ± 0.17	$\textbf{0.93} \pm 0.05$
Average	0.41	0.50	0.70	0.82

Table 2: Results of Learned Policies for Videos of self-collected videos (3 seeds). The mean \pm std of the success rates are shown in the table. For the videos concerning more complex objects or behaviors, such as throwing, we achieved significantly better performance compared to the baselines.

Task (Succ.)	Code-as-Policy	RoboGen	Eureka	Video2Policy
Insert Fork into Container Sliding Remoter to Mouse Throw Carlic into Bowl	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 0.27 \pm 0.38 \\ 0.53 \pm 0.12 \\ 0.03 \pm 0.05 \end{array}$	$\begin{array}{c} 0.57 \pm 0.40 \\ 0.87 \pm 0.05 \\ 0.37 \pm 0.29 \end{array}$	$0.93 \pm 0.05 \\ 0.97 \pm 0.05 \\ 0.70 \pm 0.36$
Average	0.00 ± 0.00	0.03 ± 0.05	0.60	0.70 ± 0.30

performance of RoboGen drops a lot, while the Eureka has much larger variances during training.It indicates that the success function is significant for more complex tasks. Moreover, the visual information of the object relationship is important for proposing the success function.

To further demonstrate the ability of Video2Policy in policy learning, we propose three harder tasks from self-collected videos, including the non-convex object manipulation and dynamic tasks. The results are illustrated on Tab. 2, which shows that our method can still resolve the harder tasks with a success rate of 87% on average. Instead, the three baselines perform much worse, especially for the CoP and RoboGen. CoP fails in all cases because the script policy receives feedback from the environment less frequently and cannot control the speed of the object. For example, the task throw garlic into bowl will reset the bowl into some place where the agent cannot reach. However, CoP can only pick and place into a precise position, which fails in the throwing motion. For RoboGen, the reward function has a higher probability of making mistakes without iterative reflection. We also conduct the ablations for the accuracy of the success functions in Sec. 5.3.

5.2 POLICY GENERALIZATION ANALYSIS FROM DIVERSE VIDEOS

To further demonstrate the scalability of our framework, we target training a general policy from diverse videos. As mentioned in Sec. 4.3, we regard the Video2Policy as a data engine to generate expert trajectories in simulation. To validate the generalization ability across unseen videos, we focus on one single behavior, lifting up. Specifically, we sample 100 videos concerning the lifting behavior from SSv2 dataset, generate the scenes, and train policies for the corresponding task. Af-terward, we collect 100 successful trajectories from each policy model, including the 256×256 size of RGB image observation and the 7-dim actions. Then we train an image-based policy model by imitation learning from the collected trajectories. Finally, we sample another 10 different lifting videos and evaluate the performance on those tasks with novel objects.

Training Details We choose imitation learning algorithms to learn the general policy and use the Be havior Cloning (BC) implemented by ourselves. For the model architecture, we apply the pre-trained
 Resnet18 (He et al., 2016) as the backbone to extract features and stack 2 frames as observations.



Figure 4: **Performance of the trained general policy on 10 unseen task instances.** BC-V2P outperforms BC-CoP on **9** of 10 significantly and demonstrates stronger generalization ability.

Table 3: Ablation Results toward the visual information, success function picking, reward function iteration, and reward function sampling. Removing any component results in a performance drop.

Task (Succ.)	Lifting up	Uncover	Throw	Avg.
Video2Policy (V2P)	0.93 ± 0.05	$\textbf{0.97} \pm 0.05$	$\textbf{0.70} \pm 0.36$	0.87
V2P w.o. visual information	0.90 ± 0.08	0.77 ± 0.13	0.57 ± 0.17	0.75
V2P w.o. success picking	0.97 ± 0.05	0.43 ± 0.40	0.30 ± 0.42	0.57
V2P w.o. iterative reward designs	0.53 ± 0.29	0.73 ± 0.24	0.17 ± 0.17	0.48
V2P w.o. multiple reward samples	0.73 ± 0.12	0.53 ± 0.24	0.27 ± 0.38	0.51

Then a policy head is built in a 3-layer MLP with hidden states of 512 dim. Additionally, the policies are trained on the collected trajectories for 30 epochs with a batch size of 1024. The learning rate of the policy head is 3e-4, while set to 3e-5 for the Resnet backbone. For evaluation, we evaluate the final checkpoint on the 10 lifting tasks for 3 seeds, and we average the results on 10 trajectories for each seed. The evaluation tasks include 5 objects with novel shapes as well as novel textures, and 5 objects with unseen categories.

463 Baseline of General Policy For gen-

445

446

447

464 eral policies, we consider the pro-465 posed Video2Policy as a novel data engine for expert trajectory collec-466 tion. Therefore, other data engines 467 are the baseline, which can generate 468 successful trajectories in our recon-469 structed scenes. We compare the fol-470 lowing models: (1) Code-as-Policy 471 (CoP), which can be applied to the 472 novel tasks directly with state input 473 instead of image observations; (2) 474 BC-CoP, which trains the BC gen-475 eral policy from the data collected 476 by CoP; (3) BC-V2P, same as BC-477 CoP but using the data collected by V2P. The CoP baseline demonstrates 478 how well the state-based general pol-479 icy can be, while the latter BC-CoP 480



Figure 5: Scalability of the general policy model (BC-V2P) towards the number of training tasks on lifting behavior (3 seeds). Learning from more videos can enhance the generalization ability of the BC-V2P model.

and BC-V2P results illustrate the performance of the general policy under different data engines.

Generalization to Unseen Videos The performances of the models are shown in Fig. 4. After training on 100 task instances, our BC-V2P model works on all the novel tasks and achieves remarkable
 generalization performance, marked in red. Specifically, the average success rate reaches 75%,
 while CoP is 32% and the BC-CoP is 26%. It indicates that our proposed Video2Policy framework can obtain more informative policy priors from the videos. Notably, the CoP results are based on

486 states and the variances from CoP policies are larger than those from V2P policies. And the model 487 performs worse on the objects with unseen categories compared to the objects with novel shapes. 488

Scalability Analysis Furthermore, it is significant to investigate the scalability of our framework 489 towards the number of generated tasks from videos. We analyze it by training the general policy on 490 $N \in \{10, 20, 30, \dots, 100\}$ tasks and evaluate the same 10 evaluation tasks, shown in Fig. 5. Overall, 491 increasing the number of videos continues to improve the performance of the BC-V2P model. Note 492 that current real2sim methods only use a few scenes such as 6 scenes (Torne et al., 2024). In contrast, 493 our method can leverage large amounts of scenes and learn policies, which proves the scalability. 494

5.3 ABLATION STUDY 496

497 Compared to the previous works on LLM-driven RL (Ma et al., 2023; Wang et al., 2023c), we apply 498 the visual information to propose success functions by LLM and pick the best one by LLM.

499 Success Function Evaluation Moreover, we attempt to 500 calculate the correlation between the generated success 501 functions and the manually designed ones, as shown in 502 Fig. 6. We choose the success rates calculated under different success functions as the data points. Compared 504 to Eureka, the success function generated by more vi-505 sual information can be more reasonable and closer to the ground-truth ones, with a higher correlation coefficient 506 qual to 0.83. In addition, the generated success functions 507 with visual information can be over-confident, as the data 508 points lie below the y=x curve. Instead, it can be noise 509 with caption only when generating the success functions. 510



Figure 6: Correlation analysis of generated success functions and the manual ground-truth functions.

Ablation of Code Generation Components Here, we ablate the results of removing each part in 511 the code generation as follows: (1) V2P w.o. visual information, where only the video caption 512 is provided to generate the codes; (2) V2P w.o. success picking, where we do not pick the best 513 success function by sampling 8 different success functions; (3) V2P w.o. iterative reward designs, 514 where we do not apply the iterative reward reflection to fine-tuning the reward functions; (4) V2P 515 w.o. multiple reward samples, where we only train 1 example of the reward function instead of 8 516 samples. The results are on the Tab. 3. It will encounter a performance drop without any part. The 517 most significant one can be the iterative reward reflection component, which finetunes the reward 518 function based on the previous training results. Additionally, w.o. the multiple reward samples and 519 w.o. success picking have larger variances than the others. With those techniques, we achieve better 520 performance than the previous LLM-driven methods (Ma et al., 2023; Wang et al., 2023c).

521 522 523

495

6 DISCUSSION

524 We have proposed Video2Policy, a pipeline for generating simulated tasks from human videos. We 525 show that our design enables us to effectively learn from human videos and generate high quality 526 data. And We show that our data can be used to train a general visuomotor policy that generalizes 527 to unseen tasks. This generalist policy scales favorably with the number of videos used for task 528 generation. We believe this is a step towards generalist robotic policies that can perform a wide 529 range of tasks similar to the wide range of everyday human behavior.

530 **Limitations and Future work** Our approach is based on existing pre-trained foundation models. 531 As such, it is bottlenecked by the quality of these models, particularly mesh reconstruction and 532 reward code generation. However, as these foundation models continue to improve we expect the 533 performance of our method to improve as well. Another limitation is that it does not fully leverage 534 the dynamic information contained in the video. As for the **future work**, while this paper validates the idea of a pipeline for generating data for generalist robotic policies in simulation, we would 536 like such a pipeline to be useful for real-world policies. Two avenues of future work towards this 537 direction are (i) increasing domain randomization (Andrychowicz et al., 2020) such as by LLMgenerated domain randomization schedules (Ma et al., 2024) which will lead to improved sim-to-538 real transfer, or (ii) improved scene reconstruction as more powerful 3D pretrained models become available (Leroy et al., 2024) which will lead to better visual quality of the produced data.

5407REPRODUCIBILITYSTATEMENT5417

The main implementations of our proposed method are in Sec. 4. The settings of the experiments, the details of training policies, and the hyper-parameters are in 5.

References

542

543

544

546 547

548

549

558

559

560

561

562

563

582

583

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020.
- Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho,
 Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. Do as i can, not as i say: Grounding
 language in robotic affordances. In *Conference on robot learning*, pp. 287–318. PMLR, 2023.
 - Zoey Chen, Sho Kiami, Abhishek Gupta, and Vikash Kumar. Genaug: Retargeting behaviors to unseen situations via generative augmentation. *arXiv preprint arXiv:2302.06671*, 2023.
 - Zoey Chen, Aaron Walsman, Marius Memmel, Kaichun Mo, Alex Fang, Karthikeya Vemuri, Alan Wu, Dieter Fox, and Abhishek Gupta. Urdformer: A pipeline for constructing articulated simulation environments from real-world images. arXiv preprint arXiv:2405.11656, 2024.
- Tianyuan Dai, Josiah Wong, Yunfan Jiang, Chen Wang, Cem Gokmen, Ruohan Zhang, Jiajun Wu,
 and Li Fei-Fei. Acdc: Automated creation of digital cousins for robust policy learning. In 8th
 Annual Conference on Robot Learning.
- Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Kiana Ehsani, Jordi Salvador, Winson Han, Eric Kolve, Aniruddha Kembhavi, and Roozbeh Mottaghi. Procthor: Large-scale embodied ai using procedural generation. *Advances in Neural Information Processing Systems*, 35:5982–5994, 2022.
- 571
 572
 573
 574
 Yilun Du, Sherry Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *Advances in Neural Information Processing Systems*, 36, 2024.
- Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al.
 The" something something" video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pp. 5842–5850, 2017.
- Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang, Stone Tao, Xinyue Wei, Yunchao Yao, et al. Maniskill2: A unified benchmark for generalizable manipulation skills. *arXiv preprint arXiv:2302.04659*, 2023.
 - Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 585 Milos Hauskrecht, Nicolas Meuleau, Leslie Pack Kaelbling, Thomas L Dean, and Craig Boutilier.
 586 Hierarchical solution of markov decision processes using macro-actions. *arXiv preprint* 587 *arXiv:1301.7381*, 2013.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- 592 Cheng-Chun Hsu, Zhenyu Jiang, and Yuke Zhu. Ditto in the house: Building articulation models of
 593 indoor scenes through interactive perception. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pp. 3933–3939. IEEE, 2023.

615

616

617

621

632

633

634

594	Sizuan Huang Zhangkai Jiang Hao Dong Vu Oigo Pang Gao and Hongshang Li InstructOact
	Siyuan muang, Zhengkai mang, mao Dong, mu Qiao, reng Oao, and mongsheng Li. mstruct2act.
595	Mapping multi-modality instructions to robotic actions with large language model. arXiv preprint
596	<i>arXiv:2305.11176</i> , 2023a.

- Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer:
 Composable 3d value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023b.
- Jemin Hwangbo, Joonho Lee, Alexey Dosovitskiy, Dario Bellicoso, Vassilios Tsounis, Vladlen Koltun, and Marco Hutter. Learning agile and dynamic motor skills for legged robots. *Science Robotics*, 4(26):eaau5872, 2019.
- Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*, 2023.
- Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r.
 arXiv preprint arXiv:2406.09756, 2024.
- ⁶¹⁰ Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín⁶¹¹ Martín, Chen Wang, Gabrael Levine, Michael Lingelbach, Jiankai Sun, et al. Behavior-1k: A
 ⁶¹² benchmark for embodied ai with 1,000 everyday activities and realistic simulation. In *Conference*⁶¹³ *on Robot Learning*, pp. 80–93. PMLR, 2023.
 - 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.
- William Liang, Sam Wang, Hung-Ju Wang, Yecheng Jason Ma, Osbert Bastani, and Dinesh Jayara man. Environment curriculum generation via large language models. In 8th Annual Conference on Robot Learning.
- Kevin Lin, Christopher Agia, Toki Migimatsu, Marco Pavone, and Jeannette Bohg. Text2motion:
 From natural language instructions to feasible plans. *Autonomous Robots*, 47(8):1345–1365, 2023.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023.
- Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayara man, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via
 coding large language models. *arXiv preprint arXiv:2310.12931*, 2023.
 - Yecheng Jason Ma, William Liang, Hung-Ju Wang, Sam Wang, Yuke Zhu, Linxi Fan, Osbert Bastani, and Dinesh Jayaraman. Dreureka: Language model guided sim-to-real transfer. *arXiv preprint arXiv:2406.01967*, 2024.
- Liane Makatura, Michael Foshey, Bohan Wang, Felix HähnLein, Pingchuan Ma, Bolei Deng, Megan Tjandrasuwita, Andrew Spielberg, Crystal Elaine Owens, Peter Yichen Chen, et al. How can large language models help humans in design and manufacturing? *arXiv preprint arXiv:2307.14377*, 2023.
- Denys Makoviichuk and Viktor Makoviychuk. rl-games: A high-performance framework for reinforcement learning. https://github.com/Denys88/rl_games, May 2021.
- Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021.
- 647 Zhao Mandi, Yijia Weng, Dominik Bauer, and Shuran Song. Real2code: Reconstruct articulated objects via code generation. *arXiv preprint arXiv:2406.08474*, 2024.

648 649 650	Amy McGovern, Doina Precup, Balaraman Ravindran, Satinder Singh, and Richard S Sutton. Hi- erarchical optimal control of mdps. In <i>Proceedings of the Tenth Yale Workshop on Adaptive and</i> <i>Learning Systems</i> , pp. 186–191, 1998.
652 653 654	Soroush Nasiriany, Abhiram Maddukuri, Lance Zhang, Adeet Parikh, Aaron Lo, Abhishek Joshi, Ajay Mandlekar, and Yuke Zhu. Robocasa: Large-scale simulation of everyday tasks for generalist robots. <i>arXiv preprint arXiv:2406.02523</i> , 2024.
655 656 657 658	Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. <i>arXiv preprint arXiv:2310.08864</i> , 2023.
659 660 661	Luigi Piccinelli, Yung-Hsu Yang, Christos Sakaridis, Mattia Segu, Siyuan Li, Luc Van Gool, and Fisher Yu. Unidepth: Universal monocular metric depth estimation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 10106–10116, 2024.
662 663	Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. <i>arXiv preprint arXiv:2209.14988</i> , 2022.
665 666 667	Can Qin, Shu Zhang, Ning Yu, Yihao Feng, Xinyi Yang, Yingbo Zhou, Huan Wang, Juan Carlos Niebles, Caiming Xiong, Silvio Savarese, et al. Unicontrol: A unified diffusion model for controllable visual generation in the wild. <i>arXiv preprint arXiv:2305.11147</i> , 2023.
668 669 670	Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. <i>arXiv preprint arXiv:2408.00714</i> , 2024.
672 673 674	Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. <i>arXiv preprint arXiv:2205.06175</i> , 2022.
675 676 677	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. <i>arXiv preprint arXiv:2308.12950</i> , 2023.
679 680	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> , 2017.
681 682 683	Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection. <i>arXiv preprint arXiv:2303.11366</i> , 2(5):9, 2023.
684 685 686	Satinder Singh, Richard L Lewis, and Andrew G Barto. Where do rewards come from. In <i>Proceed</i> - ings of the annual conference of the cognitive science society, pp. 2601–2606. Cognitive Science Society, 2009.
687 688 689 690 691	Jost Tobias Springenberg, Abbas Abdolmaleki, Jingwei Zhang, Oliver Groth, Michael Bloesch, Thomas Lampe, Philemon Brakel, Sarah Bechtle, Steven Kapturowski, Roland Hafner, et al. Of- fline actor-critic reinforcement learning scales to large models. <i>arXiv preprint arXiv:2402.05546</i> , 2024.
692 693 694 695	Sanjana Srivastava, Chengshu Li, Michael Lingelbach, Roberto Martín-Martín, Fei Xia, Kent Elliott Vainio, Zheng Lian, Cem Gokmen, Shyamal Buch, Karen Liu, et al. Behavior: Benchmark for everyday household activities in virtual, interactive, and ecological environments. In <i>Conference on robot learning</i> , pp. 477–490. PMLR, 2022.
696 697 698 699	Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. <i>arXiv preprint arXiv:2405.12213</i> , 2024.
700	Marcel Torne, Anthony Simeonov, Zechu Li, April Chan, Tao Chen, Abhishek Gupta, and Pulkit Agrawal. Reconciling reality through simulation: A real-to-sim-to-real approach for robust ma-

 Marcel Torne, Anthony Simeonov, Zechu Li, April Chan, Tao Chen, Abinshek Gupta, and Pukit
 Agrawal. Reconciling reality through simulation: A real-to-sim-to-real approach for robust manipulation. arXiv preprint arXiv:2403.03949, 2024.

702 703 704	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. <i>arXiv preprint arXiv:2305.16291</i> , 2023a.
705	
706	Lini Ware Viene Line Zhachare Vien Makit Shridhar Char Day Vieha Oir Dailin Ware
707 708	Huazhe Xu, and Xiaolong Wang. Gensim: Generating robotic simulation tasks via large language models. <i>arXiv pranriat arXiv:2310.01361</i> , 2023b
709	models. <i>urxiv preprint urxiv.2510.01501</i> , 20250.
710	
711 712	discrete and continuous control with limited data. <i>arXiv preprint arXiv:2403.00564</i> , 2024.
713	
714 715 716	Yufei Wang, Zhou Xian, Feng Chen, Tsun-Hsuan Wang, Yian Wang, Katerina Fragkiadaki, Za- ckory Erickson, David Held, and Chuang Gan. Robogen: Towards unleashing infinite data for automated robot learning via generative simulation. <i>arXiv preprint arXiv:2311.01455</i> , 2023c.
717	
718 719 720	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> 35:24824–24837, 2022
721	
722	
723	Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield. Foundationpose: Unified 6d pose estimation
724	and tracking of novel objects. In Proceedings of the IEEE/CVF Conference on Computer vision and Pattern Recognition pp. 17868–17870–2024
725	<i>una Fanern Recognition</i> , pp. 17606–17879, 2024.
726	
727 728	Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. Instantmesh: Efficient 3d mesh generation from a single image with sparse-view large reconstruction models.
729	<i>arxiv preprim arxiv.2404.07191</i> , 2024.
730	
732 733	weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games with limited data. <i>Advances in neural information processing systems</i> , 34:25476–25488, 2021.
734 735 736 737	Tianhe Yu, Ted Xiao, Austin Stone, Jonathan Tompson, Anthony Brohan, Su Wang, Jaspiar Singh, Clayton Tan, Jodilyn Peralta, Brian Ichter, et al. Scaling robot learning with semantically imagined experience. <i>arXiv preprint arXiv:2302.11550</i> , 2023a.
738	
739 740	Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montse Gonzalez Are- nas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, et al. Language to
741	rewards for robotic skill synthesis. arXiv preprint arXiv:2306.08647, 2023b.
742	
743	
744 745	A APPENDIX
746 747	A.1 IMPLEMENTATION DETAILS OF CODE GENERATION
748	Examples of Generated Codes To further illustrate the generated codes described in Sec. 4.2.
749	we take the task insert fork into storage as an example, shown in Fig. 7. We provide
750	detailed prompts and examples to inform the GPT-40 for code generation. In the beginning, we gen-
751	erate 8 code samples and then leverage GPT-40 to select the most reasonable one as the base code.
752 753	Afterward, we perform iterations for the reward function part, during which only the reward function code is required to generate. For each iteration, we learn RL policies under the N different reward

functions and select the one with the highest success rates as the base code for the next iteration.
 Noticed that the VLM can add some states such as dist_to_fork and change the observation space automatically.

from fac	tory tacks task import Task
from fac	tory tasks task import Task
 scene_in	Jory. Casks. Cask Import lask
scene_in	
	<pre>fo = {'task_class_name':, 'objects':}</pre>
lass In """	<pre>sertForkIntoStorageBox(Task):</pre>
inse	rt fork into storage box
def	
	<pre>super()init(envs, urdf_assets_root, scene_info)</pre>
def	<pre>self.task_description = "insert fork into storage box" reset objects states(self, env ids):</pre>
	<pre>return super().reset_objects_states(env_ids)</pre>
dei	<pre>compute_observation_key(sell, states): obs_keys = super().compute_observation_key(states)</pre>
	# Add distance observation between the end-effector and the fork
	<pre>states["dist_to_fork"] = dist_to_fork</pre>
	<pre>>bs_keys.append("dist_to_fork") # add distance the fact and the stores have</pre>
	<pre># Add distance observation between the fork and the storage box dist_fork_to_box = torch.norm(states["fork_pos"] - states["storage_box_pos"], dim=-1)</pre>
	<pre>states["dist_fork_to_box"] = dist_fork_to_box</pre>
	<pre>pus_keys.append("dist_fork_to_box") return obs_keys</pre>
def	<pre>yet_observation_space(self):</pre>
	<pre># Increase the observation space by the number of new observations added</pre>
dof	return obs_space + 2
der	# the fork in the storage
	# xy position of the fork is within the storage box # z position of the fork is within the storage box
	<pre>xy_distance = torch.norm(states["fork_pos"][:, :2] - states["storage_box_pos"][:, :2], dim=-1)</pre>
	 # fork should be inside the storage boy's boundary in 7-avis
	<pre>z_condition = (lower_z_of_fork < upper_z_of_box) & (upper_z_of_fork > lower_z_of_box)</pre>
	<pre># fork should be inside the storage box's boundary in XY-axis xy condition = xy distance < 0.01</pre>
	success = z_condition & xy_condition
def	<pre>return success compute reward(self, actions, states):</pre>
	return reward_function(actions, states)
# REWARD	FUNCTION
torch.j	it.script
	Figure 7: An example of the generated code for the task Insert.
	Figure 7: An example of the generated code for the task Insert.
	Figure 7: An example of the generated code for the task Insert.
	Figure 7: An example of the generated code for the task Insert.
	Figure 7: An example of the generated code for the task Insert.
¥## Anal	Figure 7: An example of the generated code for the task Insert.
 ¥ <i>## Anal</i> **Dis	Figure 7: An example of the generated code for the task Insert.
*## Anal '. **Dis - **	Figure 7: An example of the generated code for the task Insert.
 ∉## Anal L. **Dis - ** - **	Figure 7: An example of the generated code for the task Insert.
### Anal . **Dis - ** - ** - ** - **	Figure 7: An example of the generated code for the task Insert.
### Anal L. **Dis - ** → - ** 2. **Dis - **	Figure 7: An example of the generated code for the task Insert.
<pre>### Anal 1. **Dis - ** - ** - ** 2. **Dis - ** - ** - ** - ** - ** - ** - ** - *</pre>	<pre>Figure 7: An example of the generated code for the task Insert. ysis of Policy Feedback tance from Gripper to Garlic (dist_grasp_reward)** Deservation**: The values fluctuate around 0.1, with a minimum of 0.0018 and a maximum of 0.13. The last value is 0.1272. Suggestion**: The fluctuation indicates variability, but the impact may be insufficient for su Adjust the scale or temperature slightly for better sensitivity. tance from Gerlic to Bowl (dist_graft_c_bowl)** Deservation**: The values are quite low (0.0000 - 0.0029), indicating this component is not significantly contributing. Suggestion**: for better sensitivity</pre>
### Anal L. **Dis - ** - ** - ** - ** - ** - ** - ** - *	<pre>Figure 7: An example of the generated code for the task Insert. ysis of Policy Feedback tance from Gripper to Garlic (dist_grasp_reward)** Deservation**: The values fluctuate around 0.1, with a minimum of 0.0018 and a maximum of 0.13. The last value is 0.1272. Suggestion**: The fluctuation indicates variability, but the impact may be insufficient for su Adjust the scale or temperature slightly for better sensitivity. tance from Gerlic to Bowl (dist_grain_ic_to_bowl)** Deservation**: The values are quite low (0.0000 - 0.0029), indicating this component is not significantly contributing. Suggestion**: Increase the temperature parameter for better gradient. boity Reward (vel_reward)**</pre>
### Anal 1. **Dis - ** - **	Figure 7: An example of the generated code for the task Insert.
**************************************	Figure 7: An example of the generated code for the task Insert.
**************************************	Figure 7: An example of the generated code for the task Insert.
 ### Anal 1. **Dis - ** - ** - ** - ** - ** - ** - ** - *	Figure 7: An example of the generated code for the task Insert. <pre>ysis of Policy Feedback tance from Gripper to Garlic (dist_grasp_reward)** Dbservation**: The values fluctuate around 0.1, with a minimum of 0.0018 and a maximum of 0.13 The last value is 0.1272. Suggestion**: The fluctuation indicates variability, but the impact may be insufficient for su Adjust the scale or temperature slightly for better sensitivity. There for Gerlic to Bowl (dist_grain_tc_to_bowl)** Dbservation**: The values are quite low (0.0000 - 0.0029), indicating this component is not significantly contributing. Suggestion**: Increase the temperature parameter for better gradient. socity Reward (vel_reward)** Dbservation**: The values show variability and seem to be contributing to the learning process last value is 0.1746. Suggestion**: It is effective but might need a slight adjustment in temperature and scaling. al Position Reward (final_position_reward)** Dbservation**: This component consistently shows zeros, indicating it's not achieved in any of spisodes.</pre>
<pre>### Anal 1. **Dis - ** - ** - ** - ** 3. **Vel - ** - ** - ** 1. **Fin - ** - ** - **</pre>	Figure 7: An example of the generated code for the task Insert. <pre>ysis of Policy Feedback tance from Gripper to Garlic (dist_grasp_reward)** Deservation**: The values fluctuate around 0.1, with a minimum of 0.0018 and a maximum of 0.13. The last value is 0.1272. Suggestion**: The fluctuation indicates variability, but the impact may be insufficient for su Adjust the scale or temperature slightly for better sensitivity. There for Gerlic to Bool (dist_garlic_to_bowl)** Deservation**: The values are quite low (0.0000 - 0.0029), indicating this component is not significantly contributing. Suggestion**: Increase the temperature parameter for better gradient. boity Reward (vel_reward)** Deservation**: The values show variability and seem to be contributing to the learning process last value is 0.1746. Suggestion**: It is effective but might need a slight adjustment in temperature and scaling. al Position Reward (final_position_reward)** Deservation**: This component consistently shows zeros, indicating it's not achieved in any of spisodes. Suggestion**: Rewriting or significantly adjusting this component is necessary. </pre>
 ### Anal 1. **Dis - ** - ***	<pre>Figure 7: An example of the generated code for the task Insert. ysis of Policy Feedback tance from Gripper to Garlic (dist_grasp_reward)** Deservation**: The values fluctuate around 0.1, with a minimum of 0.0018 and a maximum of 0.13 The last value is 0.1272. Suggestion**: The fluctuation indicates variability, but the impact may be insufficient for su Adjust the scale or temperature slightly for better sensitivity. tance from Gerlic to Bowl (dist_graic_tc_bowl)** Deservation**: The values are quite low (0.0000 - 0.0029), indicating this component is not significantly contributing. Suggestion**: Increase the temperature parameter for better gradient. ocity Reward (vel_reward)** Deservation**: The values show variability and seem to be contributing to the learning process last value is 0.1746. Suggestion**: It is effective but might need a slight adjustment in temperature and scaling. al Position Reward (final_position_reward)** Deservation**: Thes component consistently shows zeros, indicating it's not achieved in any of spisodes. Suggestion**: Rewriting or significantly adjusting this component is necessary. al Reward (total_reward)** Deservation**: The total reward has shown improvement over time but doesn't lead to success. I </pre>
<pre>### Anal 1. **Dis - ** - ** - ** - ** - ** - ** - ** - *</pre>	Figure 7: An example of the generated code for the task Insert. <pre>ysis of Policy Feedback tance from Gripper to Garlic (dist_grasp_reward)** Deservation**: The values fluctuate around 0.1, with a minimum of 0.0018 and a maximum of 0.13 The last value is 0.1272. Suggestion**: The fluctuation indicates variability, but the impact may be insufficient for su Adjust the scale or temperature slightly for better sensitivity. There for Gerlic to Bowl (dist_grain_tc_to_bowl)** Deservation**: The values are quite low (0.0000 - 0.0029), indicating this component is not significantly contributing. Suggestion**: Increase the temperature parameter for better gradient. socity Reward (vel_reward)** Deservation*:: The values show variability and seem to be contributing to the learning process last value is 0.1746. Suggestion**:: This effective but might need a slight adjustment in temperature and scaling. al Position Reward (final_position_reward)** Deservation**: The scale or significantly adjusting this component is nected in any of spisodes. Suggestion**: Rewriting or significantly adjusting this component is necessary. al Reward (total_reward)** Deservation**: The total reward has shown improvement over time but doesn't lead to success. If fluctuates with no episodes achieving success. </pre>
 ### Anal 1. **Dis - ** - **	Figure 7: An example of the generated code for the task Insert. prise of Policy Feedback tance from Gripper to Garlic (dist_grasp_reward)** Deservation**: The values fluctuate around 0.1, with a minimum of 0.0018 and a maximum of 0.13 The last value is 0.1272. Suggestion**: The fluctuation indicates variability, but the impact may be insufficient for su Adjust the scale or temperature slightly for better sensitivity. The come Garlic to Bowl (dist_grain(c,to_bowl)** Deservation**: The values are quite low (0.0000 - 0.0029), indicating this component is not significantly contributing. Suggestion**: The values show variability and seem to be contributing to the learning process last value is 0.1746. Suggestion**: The values show variability and seem to be contributing to the learning process last value is 0.1746. Suggestion**: This effective but might need a slight adjustment in temperature and scaling. al Position Reward (final_position_reward)** Deservation**: The total reward has shown improvement over time but doesn't lead to success. If Suggestion**: The total reward has shown improvement over time but doesn't lead to success. If Succeases. Suggestion**: Improve the component scales to better balance the total reward for successful spisodes.
<pre>### Anal 1. **Dis - *** - ** - - ** - - ** -</pre>	Figure 7: An example of the generated code for the task Insert. psis of Policy Feedback tance from Gripper to Garlic (dist_grasp_reward)** Observation**: The values fluctuate around 0.1, with a minimum of 0.0018 and a maximum of 0.133 The last value is 0.1272. Suggestion**: The fluctuation indicates variability, but the impact may be insufficient for suc Adjust the scale or temperature slightly for better sensitivity. tance from Garlic to Bowl (dist_garlic_to_bowl)** Observation**: The values are quite low (0.0000 - 0.0029), indicating this component is not significantly contributing. Suggestion**: Increase the temperature parameter for better gradient. ocity Reward (vel_reward)** Observation**: The values show variability and seem to be contributing to the learning process last value is 0.1746. Suggestion**: This effective but might need a slight adjustment in temperature and scaling. al Position Reward (final_position_reward)** Observation**: This component consistently shows zeros, indicating it's not achieved in any of spisodes. Suggestion**: Rewriting or significantly adjusting this component is necessary. al Reward (total_reward)** Observation**: The total reward has shown improvement over time but doesn't lead to success. In fluctuates with no episodes achieving success. Suggestion**: Improve the component scales to better balance the total reward for successful spisodes. Suggestion**: Improve the component scales to better balance the total reward for successful spisodes.
<pre>### Anal 1. **Dis - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - *** - ** -</pre>	Figure 7: An example of the generated code for the task Insert. psis of Policy Feedback tance from Gripper to Garlic (dist_grasp_reward)** Observation**: The values fluctuate around 0.1, with a minimum of 0.0018 and a maximum of 0.132 The last value is 0.1272. Suggestion**: The fluctuation indicates variability, but the impact may be insufficient for suc Adjust the scale or temperature slightly for better sensitivity. tance from Garlic to Bowl (dist_garlic_to_bowl)** Observation**: The values are quite low (0.0000 - 0.0029), indicating this component is not significantly contributing. Suggestion**: Increase the temperature parameter for better gradient. ocity Reward (vel_reward)** Observation**: The values show variability and seem to be contributing to the learning process last value is 0.1746. Suggestion**: This effective but might need a slight adjustment in temperature and scaling. al Position Reward (final_position_reward)** Observation**: This component consistently shows zeros, indicating it's not achieved in any of apisodes. Suggestion**: Rewriting or significantly adjusting this component is necessary. al Reward (total_reward)** Observation**: The total reward has shown improvement over time but doesn't lead to success. In fluctuates with no episodes achieving success. Suggestion**: Improve the component scales to better balance the total reward for successful pisodes. Points for Improvement le or rewrite dist_garlic_to_bowl for better contribution. or introduce final_position_reward to enhance the likelihood of success.
<pre>### Anal 1. **Dis - ** - *</pre>	Figure 7: An example of the generated code for the task Insert.

Figure 8: An example of the reward analysis for the task Throw.

810 **Examples of Code Evolutions** Additionally, we also provide an example of modifying the reward 811 function based on the reward logs during training. During the reward iteration stage, we provide 812 examples of how humans adapt and refine the reward functions based on results. Also, we add the 813 logs of each reward component during training as prompts for the next iteration. These logs will be 814 categorized into two groups: from successful trajectories and from failure trajectories. This approach allows the GPT-40 to identify which reward components are most significant in driving success. If 815 the training consistently results in failure without any successful outcomes, we will prompt GPT-40 816 to redesign the reward function to improve performance and outcomes. The analysis example by the 817 GPT-40 is shown in Fig. 8, and the newly generated reward function example is shown in Fig. 9. 818

819

820 ### Improved Reward Function 821 2 Based on the analysis and suggestions, here's an improved reward function: torch.jit.script 822 4 **def** reward function (actions, states): # type: (Tensor, Dict[str, Tensor]) -> Tuple[Tensor, Dict[str, torch.Tensor]]
""" 823 824 # Chain of Thoughts for reward designs: 7 1. Grasp the garlic. 825 ° Throw the garlic towards the bowl.
 Ensure the garlic reaches the bowl. 826 10 # Reward Components: The distance between the gripper and the garlic.
 The garlic's distance to the bowl. 827 12 13 828 14 3. The velocity of the garlic towards the bowl. 4. The final position of the garlic being inside the bowl. 15 829 16 batch_size = states['eef_pos'].shape[0] 17 830 18 rewards = torch.zeros(batch_size, device=states['eef_pos'].device, dtype=torch.float32) 831 19 reward_components = {} interval_components = (from gripper to garlic reward dist_eef_to_garlic = torch.norm(states['garlic_pos'] = states['eef_pos'], dim=-1) dist_grasp_reward = 1 = torch.tanh(8.0 * dist_eef_to_garlic) # Adjusting TANH temperature reward_components['dist_grasp_reward'] = dist_grasp_reward * 0.15 # Updated scale for emphasis # 2. Increased temperature for distance between garlic and bowl 20 832 21 22 833 23 834 24 dist_garlic_to_bowl = 1 - torch.tanh(12.0 * states['dist_garlic_to_bowl']) # Increased temperature for 835 reward_components['dist_garlic_to_bowl'] = dist_garlic_to_bowl * 0.25 # Increased scale to 0.25 26 836 reward_components['dist_garlic_to_bowl'] = dist_garlic_to_bowl * 0.25 * increased scale to 0. # 3. Velocity reward (slight temperature increase) garlic_vel = states['garlic_vel_linear'] garlic_to_bowl_dir = states['bowl_pos'] - states['garlic_to_bowl_dir, dim=-1, keepdim=True) vel_towards_bowl = torch.sum(garlic_vel * garlic_to_bowl_dir, dim=-1) vel_towards_bowl = torch.tamb(4 0 * vel_towards_bowl) # Slight_adjustment in temperature 27 28 837 838 30 31 839 vel_reward = torch.tank(4.0 * vel_towards_bowl) # Slight adjustment in temperature reward_components['vel_reward'] = vel_reward * 0.30 # Increased scale to 0.30 for stronger influence # 4. Final Position Reward (revised and scaled) dist_xy_to_center = torch.norm(states['garlic_pos'][:, :2] - states['bowl_pos'][:, :2], dim=-1) 32 33 840 34 841 35 36 842 37 38 843 844 39 reward_components['final_position_reward'] = final_position_reward 40 845 41 Total Reward including penalties 42 total_reward = dist_grasp_reward * 0.15 + dist_garlic_to_bowl * 0.25 + vel_reward * 0.30 + 846 ← final_position_reward * 0.70
step_penalty = torch.ones_like(total_reward) * 0.01 847 43 rewards = total_reward - step_penalty reward_components['total_reward'] = total_reward reward_components['step_penalty'] = step_penalty 848 45 46 849 47 return rewards, reward_components



853

854

Figure 9: An example of the newly generated reward function for the task Throw.

A.2 MORE EXPERIMENTS

Robustness analysis for each component. To quantify the robustness of our method, we evaluate 855 vision models used in our pipeline individually by testing their reconstruction success rates. For 856 instance, when assessing the grounding accuracy of DINO, we sample 20 videos. Similarly, for 857 evaluating the segmentation accuracy of SAM-2, we sample 20 successful bounding boxes generated 858 by DINO. Following this approach, we systematically test the robustness of each module in the 859 pipeline. If the reconstruction is broken or identifies the wrong object, it is classified as a failure 860 case. The results are in Tab. 4. We can find that the segmentation and mesh reconstruction parts are 861 more robust than the others. 862

863 Moreover, we also conduct experiments to evaluate the noise effects of the depth estimation component. For one thing, we use the depth-aware Realsense Camera to get the ground-truth depth of the

880 881

882

883 884 885

887 888

889

895

896 897

899900901902903904

905



Figure 10: Video2Policy achieves better performance across iteration, which demonstrates better superior Pareto optimality.

Table 4: Failure rates (smaller is better) of the vision models in our pipeline. The average failure rate is 42%, and SAM-2 is the most robust model.

Failure Rates	Grounding DINO	SAM-2	InstantMesh	FoundationPose	Avg.
SSv2 Videos	0.60	0.15	0.35	0.55	0.42

Table 5: D1 distance between the predicted depth by Unidep (Piccinelli et al., 2024) and the ground-truth depth in Sliding video.

Full size	Center region (0.8x crop)	Object bounding box
d1 distance 67.3%	86.9%	93.4%

Table 6: D1 distance between the predicted depth by Unidep (Piccinelli et al., 2024) and the ground-truth depth in Sliding video.

Object	Remote Control	Mouse
Predicted Size (l, w, h) Ground-truth Size (l, w, h)	(0.18, 0.04, 0.02) (0.21, 0.05, 0.02)	(0.17, 0.10, 0.05) (0.19, 0.08, 0.03)
Delta size (l, w, h)	(0.03, 0.01, 0.00)	(0.02, 0.02, 0.02)

self-recorded video Sliding and evaluate the d1 metric (higher is better) for the video (Piccinelli at al. 2024) d1 = $\frac{\text{Number of pixels where } \frac{|d_{pred} - d_{gl}|}{d_{gt}} > 0.1\text{m}}{T}$ The results are in Tab. 5. The depth estimation

et al., 2024).d1 = $\frac{1}{\text{Total number of pixels mather } \frac{d_{\text{gt}}}{d_{\text{gt}}}$. The results are in Tab. 5. The depth estimation of the object region is accurate. And the size error under the depth prediction is as shown in Tab. 6. The error of the size prediction is small. Furthermore, we resize the objects to the GT size and apply the previously trained model to study how the noise affects the final performance. We keep the same state inputs. The performance drops from 97% to 83%, but the model can still solve the task at a high success rate.

912 Performance analysis for iterative generation. Following Eureka (Ma et al., 2023), we visualize
913 the performance of our method and baselines after each evolution iteration in Fig. 10. And we
914 also conduct an ablation study, Video2Policy w.o. Evolution (16 Samples), which only performs
915 the initial reward generation step without iterative improvement. The results are based on the tasks
916 lifting, uncover, throw, which follow the setting in Sec. 5.3. This study examines whether,
917 given a fixed reward function budget, it is more effective to allocate resources toward iterative evolution or simply generate more first-attempt rewards. The results demonstrate that our method signifi-

cantly outperforms the other baselines across multiple iterations. Our method demonstrates superior Pareto optimality, effectively balancing multiple objectives to achieve optimal trade-offs compared to other approaches, shown in Fig. 10.

A.3 MORE ABLATION RESULTS

All the ablation experiments follow the setting in Sec. 5.3.

Tracking prompts help for policy learning. Here we investigate the tolerance of 6D position errors for the tracking part. We conduct an ablation study that removes the 6D position tracking information obtained from FoundationPose in Tab. 7. As shown, after removing the tracking information, the performance drops for the tasks with multiple objects.

Table 7: After removing the tracking information in prompts, the performance drops from 87% to 75%, which is still superior to Eureka.

	Lifting Up	Uncover	Throw	Average
Video2Policy Video2Policy, w.o. tracking info Eureka	$\begin{array}{c} \textbf{0.93} \pm 0.05 \\ 0.90 \pm 0.08 \\ 0.83 \pm 0.13 \end{array}$	$\begin{array}{c} \textbf{0.97} \pm 0.05 \\ 0.77 \pm 0.13 \\ 0.63 \pm 0.38 \end{array}$	$\begin{array}{c} \textbf{0.70} \pm 0.36 \\ 0.57 \pm 0.17 \\ 0.37 \pm 0.29 \end{array}$	0.87 0.75 0.61

	wels 1.
def	pre : compute observation key(self, states):
	obs_keys = super().compute_observation_key(states)
	return obs_keys
def	<pre>get_observation_space(self): bb_observation_space(self):</pre>
	obs_space = super().get_observation_space()
def	<pre>compute_success(self, states):</pre>
	<pre>card_lifted = (states['card_pos'][:, 2] - states['table_height']) > (torch.max(states['card_size']) /</pre>
	success = card_lifted
	return success
Exan	mple 2:
def	<pre>compute_observation_key(self, states):</pre>
	obs_keys = super().compute_observation_key(states)
	aist_gripper_card = torch.norm(states("card_pos") - states("eet_pos"), dim=-1)
	obs_keys.append("dist_gripper_card")
	return obs_keys
def	<pre>get_observation_space(self): bb_observation_space(self):</pre>
	obs_space = super().get_observation_space()
def	<pre>compute_success(self, states):</pre>
	<pre>card_height_above_table = states['card_pos'][:, 2] - states['table_height']</pre>
	card_lifted = card_height_above_table > 0.2
	success - card_fifted & (states[card_vel_fifted].nofm(dim=-1) < 0.01)
Exan	ple 2 is more reasonable and better for several reasons:
±. 1	This is beneficial because it provides an additional critical feature that can help the agent understam
	\leftrightarrow its relative position to the target object. This added information is highly useful for learning
	\hookrightarrow tasks involving object manipulation.
2 1	
2. E	The second implementation refines the success criteria by adding a condition that the card's linear
	↔ velocity should be minimal (states['card_vel_linear'].norm(dim=-1) < 0.01), ensuring that the card ↔ is not just lifted but also stable. This is a more precise definition of success for manipulation
	↔ tasks.
з. п	ynamic and Informative Observations:
4.0	Observation Space Adjustment:



Hallucination issue can be alleviated by picking under GPT-40. To evaluate the validity of the generated task code, we have some instructions. For example, for correctness, we inform the GPT-40 to read and analyze the success part; for reasonability, we inform the GPT-40 to avoid picking the code that assumes some scalar value or states. To quantify the frequency of the hallucination,

972 we run Video2Policy (16 samples, 1 iteration) for the 3 tasks (Lifting, Uncover, Throw). Here we 973 remove the while-loop generation so that not all samples are runnable. (Previously, if one sample 974 fails, we regenerate again until all 8 samples are runnable.) Here the hallucination samples include 975 non-runnable samples and zero-score samples. We can find that after picking by GPT-40, the hallu-976 cination problem alleviates in a degree. Here is an example of how GPT-40 picks better task codes, shown in Fig. 11. 977

Table 8: Querying GPT-40 for choosing across multiple samples helps alleviate hallucination.

	Without Picking	Picking by GPT-40
Hallucination	0.40	0.19
Hallucination - Non-runnable Hallucination - Runnable	0.25 0.15	0.13 0.07

The generated codes do not reset cheating. Since the task codes are generated under different 987 prompts for all the methods, we conduct ablation studies by choosing the same 'reset' function from 988 our method for the baselines. For the code-as-policy, we generate the policy code and execute it 989 in the Issac Gym. Thus, the 'reset' function is the same as the Video2Policy since they share the 990 task code. We choose 3 tasks, one of a single object and two of multiple objects. The results are illustrated in Tab. 9. For tasks with a single object, the results are the same because the generated 992 'reset' function calls the base reset function. For tasks with multiple objects, the results have limited 993 changes. The reason is that when generating the reset function, the LLM introduces certain con-994 stants. These constants may vary. However, the variance has limited effects on the final results. It 995 indicates that there is no reset cheating.

Table 9: Using the same 'reset' function as the Video2Policy, the baselines have limited changes 997 for the evaluation results. This proves that the better results of our method do not come from 'reset' 998 cheating. 999

	Lifting Up	Uncover	Throw	Average
Video2Policy Code-as-Policy	$\begin{array}{c} \textbf{0.93} \pm 0.05 \\ 0.33 \pm 0.21 \end{array}$	$\begin{array}{c} \textbf{0.97} \pm 0.05 \\ 0.10 \pm 0.08 \end{array}$	$\begin{array}{c} \textbf{0.70} \pm 0.36 \\ 0.00 \pm 0.00 \end{array}$	0.87 0.14
RoboGen RoboGen (same reset function)	$\begin{array}{c} 0.28 \pm 0.09 \\ 0.28 \pm 0.09 \end{array}$	$\begin{array}{c} 0.67 \pm 0.26 \\ 0.60 \pm 0.28 \end{array}$	$\begin{array}{c} 0.03 \pm 0.05 \\ 0.03 \pm 0.05 \end{array}$	0.33 0.30
Eureka Eureka (same reset function)	$\begin{array}{c} 0.83 \pm 0.13 \\ 0.83 \pm 0.13 \end{array}$	$\begin{array}{c} 0.63 \pm 0.38 \\ 0.67 \pm 0.21 \end{array}$	$\begin{array}{c} 0.37 \pm 0.29 \\ 0.37 \pm 0.29 \end{array}$	0.61 0.62

1008 1009 1010

1011

1007

978

991

996

A.4 EXPERIMENTS OF SIM2REAL

1012 Although this work aims to leverage internet videos for simulation policy learning, we also conduct 1013 some sim2real experiments to verify the effectiveness. Specifically, following Sec. 5.2, we collect 1014 200 trajectories from each reconstructed scene in simulation for 100 lifting tasks. Then we 1015 train a general policy via imitation learning and subsequently deployed the policy in the real-world 1016 setting.

1017 To alleviate the sim-to-real gap, we employed the following two strategies: 1018

Input Representation and Network Architecture. We take as input the 0/1 segmentation masks 1019 of the image, the robot's end-effector (EEF) state, and the gripper state. SAM-2 is adopted for 1020 segmentation, where the pixel position of the target object is provided manually as input in the 1021 first frame, shown in Fig. 12. We stack 2 frames and add an additional multi-layer perceptron 1022 (MLP) layer to map the robot state into a 256-dimensional feature vector. Furthermore, the rotation 1023 component of the action is scaled by a factor of 0.2 for better stability. 1024

Domain Randomization. During data collection in the simulation, randomization is applied to the 1025 actions with noise levels of 0.02 and a random delay of 0.01-0.02 seconds. Moreover, the physical

19

Figure 12: Examples of the 0-1 Segmentation Mask Observation between the simulation and the real, which can better bridge the sim2real gap.

properties of the object, such as size and weight, are also randomized. We ensure consistency between the camera poses in simulation and the real world.

In real-world experiments, the object's position varies within a 10 cm range. The image input had a resolution of 256x256. In terms of the setup, we use Franka robotic arms, Robotiq grippers, and Stereolabs cameras. We evaluate the performance of the policy towards lifting a mouse, a cup, and a piece of bread. Notably, while there are some mouse and bottle objects in the simulation, the bread is absent in the collected simulation dataset and is soft.

Here the general lifting policy achieves a success rate of 72% across 10 novel objects in simula-tion. The sim2real results are as shown in Tab. 10. Compared to the 72% success rate in simulation, it achieves 47% success rate in the real world. It proves the efficacy of our pipeline, which builds the pipeline of internet videos to policies. We notice that gripping slippery objects, such as a cup, pose challenges for the robot, resulting in a relatively low success rate. For the bread, the task was relatively easier despite the object being unseen in the simulation. This can be attributed to the seg-mentation mask observation and the bread's relatively large surface area, which facilitates successful manipulation.

Overall, these experiments demonstrate that the general policy trained in simulation possesses effective sim-to-real transfer capabilities. Additionally, the results highlight the potential of the proposed Video2Policy pipeline, underscoring its effectiveness in enabling good performance, scalability, and deployment in real-world scenarios.

	Mouse	Cup	Bread	Average
Succ.	0.50	0.40	0.50	0.47

Table 10: The success rates of the learned policy for lifting tasks on real robots.