ROBOTIC PROGRAMMER: VIDEO INSTRUCTED POLICY CODE GENERATION FOR ROBOTIC MANIPULATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Zero-shot generalization across various robots, tasks and environments remains a significant challenge in robotic manipulation. Policy code generation methods use executable code to connect high-level task descriptions and low-level action sequences, leveraging the generalization capabilities of large language models and atomic skill libraries. In this work, we propose Robotic Programmer (RoboPro), a robotic foundation model, enabling the capability of perceiving visual information and following free-form instructions to perform robotic manipulation with policy code in a zero-shot manner. To address low efficiency and high cost in collecting runtime code data for robotic tasks, we devise Video2Code to synthesize executable code from extensive videos in-the-wild with off-the-shelf vision-language model and code-domain large language model. Extensive experiments show that RoboPro achieves the state-of-the-art zero-shot performance on robotic manipulation in both simulators and real-world environments. Specifically, the zero-shot success rate of RoboPro on RLBench surpasses the state-of-the-art model GPT-40 by 11.6%, which is even comparable to a strong supervised training baseline. Furthermore, RoboPro is robust to different robotic configurations, and demonstrates broad visual understanding in general VQA tasks.

1 INTRODUCTION

031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

A long-term goal of embodied intelligence research is to develop a single model capable of solving any task defined by the user. Recent years have witnessed a trend towards large-scale foundation models on natural language processing tasks (Achiam et al., 2023; Touvron et al., 2023). Scaling up these language models in terms of model size and training tokens significantly improves the few-shot performance on a range of end tasks, even achieving performance comparable to previous state-of-the-art fine-tuning methods. However, for robotic tasks, we have yet to see large-scale pre-trained models that can directly transfer across different robots, tasks and environments without additional fine-tuning.

To improve the zero-shot generalization ability of robotic models, one common approach is to unify 040 different tasks as the next action prediction. This paradigm requires the model to directly generate 041 low-level action sequences. Brohan et al. (2023a); Padalkar et al. (2023); Kim et al. (2024); Niu et al. 042 (2024) collected large amount of trajectories across various robots, tasks and environments. They 043 trained vision-language-action (VLA) models derived from LLMs to map images and task instructions 044 into discrete action tokens. Despite these models achieve better performance and show the capacity to transfer on novel objects and different tasks, fine-tuning is still required when deploying on new robots and environments. Besides, it is extremely expensive to collect trajectories through real-world 046 robots, while using human-built simulators often leads to lack of diversity and introduces additional 047 gap between simulation platform and real-world usages. 048

Another line of research aims to use code as compromise solution for bridging high-level instructions
 and low-level robot execution, leveraging the generalization capabilities of Large Language Models
 (LLMs) and atomic skill libraries. RoboCodeX (Mu et al., 2024) utilizes large vision-language model
 (VLM) to generate tree-of-thought plans and grasp preference. However, it also relies on manually built simulation environment and human-annotated code for data curation, which is expensive and not friendly for scaling up in terms of training data.



Figure 1: Visualization of evaluation tasks and execution results. RoboPro shows impressive zero-shot performance on novel and compositional tasks in RLBench (a), long-termed manipulation tasks in LIBERO (b), and real-world tasks (c). Video demos can be found in our supplementary materials.

In this work, we introduce Robotic Programmer (RoboPro), a robotic foundation model, enabling
 the capability of perceiving visual information and following free-form user instructions to perform
 manipulation tasks without additional fine-tuning. RoboPro generates the executable code to connect
 high-level instructions and low-level action sequences. To address low efficiency and high cost in
 collecting runtime code data for robotic tasks, we devise Video2Code, an automatic data curation
 pipeline for multimodal code generation.

We draw our inspiration from the extensive amount of operational videos in-the-wild that implicitly 087 contain necessary procedural knowledge about how to finish operational tasks. Previous research 088 has focused on utilizing videos for large-scale supervised learning (Brohan et al., 2023a; Kim 089 et al., 2024; Niu et al., 2024) or extracting relevant knowledge (e.g., affordance (Bahl et al., 2023)), 090 while extracting executable policy code from videos is still under-explored. Our data curation 091 pipeline uses the off-the-shelf VLM and Code LLM to synthesize code execution data from videos, 092 which is much more efficient and scalable compared with generating code data from manually-built simulation environments. With Video2Code, we synthesize 115k robot execution code data along 094 with the corresponding scene information and task descriptions from DROID (Khazatsky et al., 2024). Extensive experiments (examples depicted in Figure 1) show that RoboPro achieves the state-095 of-the-art zero-shot performance on robotic manipulation tasks in both simulators and real-world 096 environments. Specifically, the zero-shot success rate of RoboPro on RLBench outperforms the 097 state-of-the-art model GPT-40 by a gain of 11.6%. It is even comparable to a strong supervised 098 training method PerAct (Shridhar et al., 2023). Furthermore, RoboPro is robust to different robotic configurations, and shows broad visual understanding on general VQA tasks. 100

101

103

076

077

079

2 RELATED WORKS

Language-guided robot manipulation. Language-conditioned robot manipulation refers to the
 use of natural language instructions to guide robotic actions. Natural language instructions allow
 non-experts to interact with robots through intuitive commands and enable robots to generalize to
 various tasks based on natural language input (Winograd, 1971). Recent advancements in language-conditioned embodied agents have leveraged Transformers (Vaswani et al., 2017) to enhance perfor-

108 mance on multi-task settings. One category of recent approaches is language-conditioned behavior 109 cloning (BC), where models learn to mimic demonstrated language-conditioned actions and output 110 dense action sequences directly. 3D BC methods (Shridhar et al., 2023; Zhang et al., 2024) trained 111 from scratch perform well on specific environment, while lacking of generalization ability across 112 environments. Vision-language-action (VLA) models (Brohan et al., 2023a; Kim et al., 2024; Niu et al., 2024) built on pre-trained large language models (LLMs) show capacity to transfer on novel 113 objects and task settings, but need additional fine-tuning when being deployed on new environments 114 and robots. Another line is to create high-level planners based on LLMs (Huang et al., 2022; Brohan 115 et al., 2023b; Driess et al., 2023; Huang et al., 2023c), which output step-by-step natural language 116 plans according to human instructions and environmental information. These methods show better 117 generalization ability across environments, leveraging the reasoning and generalization ability of 118 LLMs on language instructions and environments. However, there is still a gap between generated 119 natural language plans and low-level robotic execution, requiring an extra step to score potential 120 actions or decompose plans into relevant policies (Singh et al., 2023).

121

122 Robot-centric policy code generation. Code-as-Policies (Liang et al., 2023) proposes that exe-123 cutable code can serve as a more expressive way to bridge high-level task descriptions and low-level 124 execution. Atomic skills to perceive 3D environments and plan primitive tasks are provided in prede-125 fined API libraries. LLMs process textual inputs and generate executable policy code conditioned on the API libraries (Liang et al., 2023; Huang et al., 2023a;; Xu et al., 2023; Vemprala et al., 2024; 126 Singh et al., 2023). RoboScript (Chen et al., 2024) further suggests that unified interface facilitates 127 LLM's adaptability across different environments and hardware platforms. However, these methods 128 rely solely on linguistic inputs, requiring detailed descriptions of environments and instructions as 129 textual inputs, which limits their generalization and visual reasoning ability across environments. 130 RoboCodeX (Mu et al., 2024) utilizes large vision-language model (VLM) to decompose multimodal 131 information into object-centric units in a tree-of-thought format. Nevertheless, it relies on manually-132 built simulation environments and human-annotated data, which lacks environmental richness and is 133 expensive for scaling up. Different from previous works using language-only LLMs, RoboPro enables 134 visual reasoning ability and follows free-form instructions in a zero-shot manner. Furthermore, an 135 automatic and scalable data curation pipeline Video2Code is developed to synthesize runtime code 136 data from extensive videos in-the-wild in a quite efficient and low-cost fashion.

137 138

3 Method

139 140 141

3.1 PROBLEM STATEMENT

142 We consider language-guided robotic manipulation where each task is described with a free-form language instruction I. Given RGBD data from the wrist camera as the observation space O_t and 143 robot low-dimension state s_t (e.g., gripper pose at current time t), the central problem investigated 144 in this work is how to generate motion trajectories T, where T denotes a sequence of end-effector 145 waypoints to be executed by an Operational Space Controller (Khatib, 1987). However, generating 146 dense motion trajectories at once according to the free-form instruction I is quite challenging, as I 147 can be arbitrarily long-horizon and would require comprehensive contextual understanding. Policy 148 code generation methods map long-horizon instructions to a diverse set of atomic skills, leading 149 to rapid adaptation capabilities across various robotic platforms. With comprehensive contextual 150 understanding and advanced visual grounding capabilities, large vision-language models can function 151 as intelligent planners, translating the task execution process into generated programs due to their 152 robust emergent capabilities.

To prompt vision-language models (VLMs) to generate policy code, we assume a set of parameterized skills with unified interface, which is defined as the API library L_{API} . L_{API} can be categorized into perception module L_{per} and control module L_{con} based on the API's role in task execution process. L_{per} is tasked with segmenting the task-relevant part point cloud Π_I and predicting the physical property ϕ_I of relevant objects, while L_{con} predicts the contact pose of the gripper and generates the motion trajectory T based on the output of L_{per} and the current robot state s_t :

$$L_{API} = \{L_{per}, L_{con}\} \\ \{\Pi_I, \phi_I\} = L_{per}(O_t, I) \\ T = L_{con}(s_t, \{\Pi_I, \phi_I\}).$$
(1)

189

190

191 192

200

201

203



Figure 2: The data curation pipeline of Video2Code. We first use the Draft VLM to extract a brief natural language plan for execution of the user instruction. After that, the Code LLM generates robot-centric code using the provided API library and natural language plan from the first stage.

193 With the visual observation and the language instruction, VLMs generate executable policy code 194 ${\pi_i, p_i}_{i=1}^N$ conditioned on the API library L_{API} , where π_i denotes the *i*-th L_{per} or L_{con} calls and p_i 195 represents corresponding parameters for API calls. Each API call generates a sub-trajectory sequence τ_i of arbitrary length (the length is ≥ 0). All sub-trajectory sequences $\{\tau_i\}_{i=1}^N$ are then concatenated 196 to form the final complete motion trajectory T. The whole generation process is formulated as: 197

$$O_t, I) \stackrel{\text{VLM}}{\Longrightarrow} \{\pi_i, p_i\}_{i=1}^N \Longrightarrow \{\tau_i\}_{i=1}^N. \tag{2}$$

Explainable API calls generated by VLMs connect the observation and high-level instructions to lowlevel execution, enabling the capacity of zero-shot generalization in free-form language instructions 202 and across different environments. Obviously, training such VLMs to perceive environments, follow instructions and generate executable code will inevitably require a vast amount of diverse and 204 well-aligned robot-centric multimodal runtime code data, which poses a significant challenge. 205

206 3.2 VIDEO2CODE: SYNTHESIZE ROBOTIC RUNTIME CODE FROM VIDEOS 207

208 Videos are widely available raw data sources for runtime code data synthesis. Extensive operational 209 videos naturally provide low-level details of performing tasks such as "how to pour tea into a cup", 210 which inherently contain necessary procedural knowledge for runtime code data. Despite their 211 favorable diversity and considerable quantity, it is still an under-explored and challenging problem 212 how to collect executable policy code from demonstration videos efficiently. To this end, we devise 213 Video2Code, a low-cost and automatic data curation pipeline to synthesize high-quality runtime code data from videos in an efficient way. Although open-source or lightweight vision-language 214 models exhibit promising performance on video understanding tasks, a performance gap remains 215 when compared to code-domain large language models in handling complex code generation tasks.


Figure 3: The overview of RoboPro. RoboPro utilizes environmental observation and natural language instruction as multimodal input, then outputs executable policy code. Extendable API library plays a role in mapping policy code into low-level execution sequences.

As depicted in Figure 2, to combine the visual reasoning ability of VLM and coding proficiency of code-domain LLM, Video2Code adopts a two-stage strategy.

Plan extraction. The first stage is to extract robot-centric plans in natural language from in-structional videos. These instructional videos are filtered from DROID (Khazatsky et al., 2024), a large-scale robot manipulation dataset with 350 hours of interaction data across 564 scenes, 86 tasks, and 52 buildings. We extract 50k independent instructional videos with at least one free-form human instruction and further clip each video into 16 key frames. After that, we use Gemini-1.5-Flash (Team et al., 2023) as the Draft VLM to generate a brief list of actions for human instruction with these key frames as reference. As shown in Figure 2, the Draft VLM generates a step-by-step robot-centric plan from an instructional video to "stack the cups together". The generated natural language plans contain knowledge and habit of human to follow free-form embodied instructions, and key visual information is extracted automatically from the instructional video.

Policy code generation. After plan extraction, we use Code LLM DeepSeek-Coder-V2 (Zhu et al., 2024) to "translate" these natural language plans into executable code. A complete prompt fed into the Code LLM includes API definitions, the natural language plan, and auxiliary part containing rules to follow. In the API definitions part, parameterized API functions are classified into two categories as formulated in Sec. 3.1: perception module, and control module. For each of these API functions, we provide API definitions and descriptions to demonstrate their usage. Auxiliary part contains prefix, third party tools, and rules to follow, similar to previous practices in RoboCodeX (Mu et al., 2024). Natural language plans accompanied with original human instructions are attached at the end of the prompt. As shown in Figure 2, step-by-step decomposed natural language plan guides the Code LLM to generate high-quality policy code in a Chain-of-Thought format. As for API implementation, we use GroundingDINO (Liu et al., 2023) and AnyGrasp (Fang et al., 2023) to get the bounding boxes and grasp preferences, respectively. Besides, we provide heuristic implementation for compositional skills. We finally collect 115k runtime code data with task descriptions and environmental observations using Video2Code for supervised fine-tuning.

3.3 ROBOPRO: ROBOTIC FOUNDATION MODEL

Model architecture. As shown in Figure 3, RoboPro has a vision encoder and a pre-trained LLM. They are connected with a lightweight adaptor layer consisting of a two-layer MLP. Specifically, the vision backbone first encodes the image into a sequence of visual tokens. After that, the lightweight adaptor is designed to project visual tokens onto embedding space of the LLM. In addition, we provide the API definitions and the user instruction as the text inputs. The visual and text tokens are directly concatenated and then fed into the LLM, as similarly done in Liu et al. (2024b). The LLM are trained to generate the runtime code based on the visual inputs and task description.

RoboPro is designed to reason on multimodal inputs and generate executable policy code for robotic
manipulation. Thus, two key factors for the choice of its components are the ability of visual
reasoning and the quality of policy code generation. RoboPro adopts SigLIP-L (Zhai et al., 2023)
as the vision encoder, which yields favorable performance on general visual reasoning tasks. For
the base LLM, a code-domain LLM, CodeQwen-1.5 (Bai et al., 2023), is utilized, which shows
state-of-the-art performance among open-source code models. The model architecture and working
process of RoboPro are illustrated in Figure 3.

278 **Training.** The training procedure of RoboPro consists of three stages: visual alignment, pre-training, and supervised fine-tuning (SFT). We first train a lightweight adaptor layer while freezing the vision 279 encoder and LLM with LLaVA-Pretrain (Liu et al., 2024b). Then we pre-train the lightweight 280 adaptor and the LLM on a corpus of high-quality image-text pairs (Chen et al., 2023). For supervised 281 fine-tuning, the 115k runtime code data generated by Video2Code (as noted in Sec. 3.2) are used. To 282 avoid overfitting and enhance visual reasoning ability, a general vision language fine-tuning dataset 283 (LLaVA-1.5 (Liu et al., 2024b)) is also involved during the SFT process. Thus, RoboPro is trained to 284 follow free-form language instructions and perceive visual information to generate executable policy 285 code for robotic manipulation. Meanwhile, it exhibits broad visual understanding to perform general 286 VQA tasks. Our code and model will be released to the public.

288

287

289 290

291

293

294

295

296

4 **EXPERIMENTS**

4.1 ZERO-SHOT ROBOTIC MANIPULATION

Setup. Following PerAct (Shridhar et al., 2023), we select 9 tasks with the requirement of novel instruction understanding or long-horizon reasoning in RLBench (James et al., 2020) for evaluation. Each task is evaluated with 25 episodes scored either 0 or 100 for failure or success in task execution. Detailed experiment settings and task information in RLBench can be found in Appendix A.1.

297 **Baselines.** The baselines can be categorized into two groups. One common approach requires super-298 vised training on the simulation platform, e.g., behavior cloning methods, including PerAct (Shridhar 299 et al., 2023) and LLARVA (Niu et al., 2024). They are either trained from scratch or fine-tuned with hundreds of episodes from RLBench. PerAct is trained on 100 episodes, and LLARVA is fine-tuned 300 on 800 episodes per task in RLBench. The methods from another group do not require additional 301 training. They first output robot-centric policy code, then execute it with provided APIs. We evaluate 302 their zero-shot performance on RLBench. CaP (Liang et al., 2023) equips large language model with 303 the ground-truth textual scene descriptions, containing object names, attributes, and instructions, to 304 generate executable code. Following their paper, we implement CaP with GPT-3.5-Turbo (gpt-3.5-305 turbo-0125). GPT-40 (OpenAI (2024), gpt-40-2024-05-13) is the state-of-the-art multimodal model 306 for various vision-language tasks. For RoboPro and GPT-40, we require the model to directly generate 307 the executable code given the image from the wrist camera, user instructions and API definitions. We 308 also analyze the generalization ability of RoboPro on the formation of API libraries, which is further 309 elaborated in Sec. 4.2. For a fair comparison, we adopt the same API library for these methods (i.e., 310 CaP, GPT-40, and RoboPro). Our API library shares similar design formulation as RoboCodeX (Mu et al., 2024)¹, with detailed implementation in Appendix B. 311

312

Results. We report the average success rate on 25 313 episodes for each task. As shown in Table 1, the zero-314 shot result of RoboPro surpasses language-only pol-315 icy code generation method (CaP) by 19.1%. Besides, 316 our model significantly outperforms the state-of-the-317 art VLM GPT-40 by 11.6% on average success rate. 318 More importantly, the zero-shot success rate of Robo-319 Pro is even comparable with a strong behavior-cloning 320 baseline PerAct that requires supervised training. It



Figure 4: Error breakdown on RLBench.

demonstrates the effectiveness of our model for manipulation tasks. To thoroughly analyze the factors

321

¹RoboCodeX was not compared in our experiments as this work has not released its model checkpoints and its original evaluations were not conducted on publicly available simulation platforms.

³²² 323

Models	Push Buttons	Stack Blocks	Open Drawer	Close Jar	Stack Cups	Sweep Dirt	Slide Block	Screw Bulb	Put in Board	Avg.
Specialists, w/ training on RI	Bench									
LLARVA (Niu et al., 2024)	56	0	60	28	0	84	100	8	0	37.3
PerAct (Shridhar et al., 2023)	48	36	80	60	0	56	72	24	16	43.6
CaP (Liang et al., 2023) GPT-40 (OpenAI, 2024)	KLBenci 72 72	4 20	24 56	40 36	0 4	36 40	4 20	20 20	12 12	23.6 31.1
RoboPro (ours)	68	48	68	44	4	48	60	32	12	42.7
w/ API Renaming	68	40	60	48	4	48	68	36	12	42.7
w/ API Refactoring	68	36	72	44	8	16	80	28	12	40.4

Table 1: Success rate (%) on RLBench Multi-Task setting. Methods greyed on need supervised training on the simulation platform.

Table 2: Success rate (%) on 8 tasks in LIBERO. The result of PerAct is a zero-shot transfer from RLBench to LIBERO.

Models	Turn on Stove	Close Cabinet	Put in Sauce	Put in Butter	Put in Cheese	Place Book	Boil Water	Identify Plate	Avg.
Specialists, RLBench \rightarrow LIB PerAct (Shridhar et al., 2023)	ERO 0	0	0	0	0	0	0	0	0
Generalists CaP (Liang et al., 2023)	0	37	17	13	7	30	7	7	14.8
GPT-40 (OpenAI, 2024)	37	17	63	43	57	43	17	3	35.0
RoboPro (ours)	97	60	67	53	63	43	23	13	52.4

contributing to the performance gap between different methods, we conducted an error breakdown for the policy code generation approaches. In the context of policy code generation methods, the successful execution of manipulation tasks relies on both the accuracy of the policy code and the capabilities of the API library. The main types of errors impacting the quality of robot-centric policy code are logical errors, functional errors, and grounding errors. These errors are associated with challenges in the appropriate selection and utility of APIs, as well as issues related to visual grounding. As depicted in Figure 4, the results show that all these methods perform well on following functional definition of API library, causing a low occupancy of functional error. Compared with linguistic only method CaP, GPT-40 and RoboPro show a noticeable improvement in target object grounding. The main failure cases of CaP and GPT-40 fall in logical error, including API selection and proper order of API calls. In contrast, RoboPro effectively reduces this margin, mainly owing to the procedural knowledge about long-term execution learned in Video2Code. Execution errors maintain a consistent proportional relationship with successful cases, which result from API limitations rather than inaccuracies in the policy code. Detailed illustration of error cases can be found at Appendix A.4.

4.2 ZERO-SHOT GENERALIZATION

In supervised training methods, outstanding performance is often achieved in familiar environments.
However, due to data scarcity, challenges arise in terms of generalizing across different environments and robot configurations in a zero-shot manner. Different environments mainly introduce variations in tasks, context, and objects, while robot configurations refer to changes of degrees of freedom, action spaces in different robotic embodiments. For policy code generation methods, different robot configurations are mainly reflected in variations in the formats of APIs (e.g. input-output format).
Additionally, users may have their own preferences when customizing API libraries for similar functions. We believe that a robust policy code generation model should also demonstrate strong

adaptability to these variations. To evaluate the zero-shot generalization ability of RoboPro, we
 conduct experiments in two aspects: across different environments and across different API libraries.

Generalization across different environments. We use LIBERO (Liu et al., 2024a) as an extra 382 simulator to evaluate the zero-shot generalization across different environments. We choose 8 383 representative tasks from LIBERO-100 as the evaluation set. Each task is evaluated with 30 episodes. 384 These tasks include short-horizon tasks which need scene understanding, and long-horizon tasks 385 which require multi-step implementation. Detailed task descriptions and corresponding examples 386 can be found in Appendix A.2. For behavior cloning method PerAct, we evaluate its model trained 387 on RLBench as described in Sec. 4.1, which will be tested on LIBERO without further fine-tuning. 388 For CaP, GPT-40 and RoboPro, we evaluate their zero-shot performance. As reported in Table 2, 389 PerAct trained on RLBench struggles on the tasks from LIBERO. It indicates that PerAct is difficult 390 to generalize across different environments without additional fine-tuning. Furthermore, RoboPro significantly outperforms GPT-40 by a gain of 17.4% average success rate on 8 LIBERO tasks, 391 which is aligned with the observations from the experiments on RLBench. Compared with GPT-40, 392 RoboPro executes more accurate sequences of actions to complete various manipulation tasks. For 393 instance, when given the task "Turn on the stove", RoboPro consistently approaches the stove knob, 394 grasps it, and rotates it clockwise. In contrast, GPT-40 sometimes misinterprets the knob's affordance, 395 attempting to press it rather than rotate. 396

397

381

Generalization across different API libraries. The formation and definition of pre-defined API 398 library is a key factor that affects the performance of general robotic models, since they are usually 399 deployed across different types of robots. Robustness to the changes of API library implies that 400 the model can understand and internalize the atomic skills under the API interface. To assess 401 the generalization of RoboPro under different level of changes in API library, we designed two 402 representative sets of experiments: the API Renaming set and the API Refactoring set. For renamed 403 APIs, we only change in their names and keep consistent in functional structure (e.g., the type of 404 return values and arguments). For refactored APIs, we change in functional structure but keep their 405 names. Take the control API "get_best_grasp_pose()" as an example. In the API Renaming 406 set, it is renamed as "generate_obj_grasp_pos()" without changes on functionality, and in 407 the API Refactoring set, the inputs, outputs and comments are all changed (e.g., the input format changes from "bbox" to "np.ndarray"). As shown in Table 1, the performance of RoboPro on 408 RLBench is robust to the changes in API formation. The detailed implementations of renamed and 409 refactored APIs can be found in Appendix B. 410

411 412

413

4.3 REAL-WORLD EXPERIMENTS

414 To evaluate the performance of RoboPro in real-world 415 scenarios, we conduct realistic experiments on a Franka 416 Emika robot arm equipped with an Intel RealSense 417 D435i wrist camera. As emphasized in Sec. 3.1, longhorizon task decomposition and visual understanding 418 capabilities are crucial for zero-shot generalization in 419 language-guided robotic manipulation. To assess Robo-420 Pro's performance in these aspects, we carefully de-421 signed 8 tasks, ranging from short-horizon to long-422 horizon tasks, as well as tasks that require visual com-423 prehension. For instance, RoboPro is required to select 424 object with "wipe" affordance from the scene given 425 instruction "wipe the desk". Additionally, to rigorously 426 validate RoboPro's generalization capability across dif-427 ferent real-world scenarios, we ensure that each task

Table 3: Th	e zero-shot s	success rate of	of Robo-
Pro across	8 real-world	manipulatic	on tasks.

Task	# Var	# Test	Succ. %
Move in Direction	2	10	80
Setup Food	2	10	90
Distinct Base	2	10	70
Prepare Meal	2	10	60
Tidy Table	2	10	70
Express Words	4	10	60
Stack on Color	5	10	50
Wipe Desk	2	10	100

428 consists of at least two variations (denoted as "# Var") in terms of object categories and physical
429 properties (10 tests are run for each task). As shown in Table 3, RoboPro is able to achieve 72.5%
430 success rate on average among all 8 tasks, which verifies RoboPro's strong generalization ability in
431 real-world scenarios without any specific fine-tuning. We also observe RoboPro exhibits impressive
emergent ability in visual reasoning. For example, as depicted in Figure 5, when asked to wipe the

 (a) Instruction: say yes
 (b) Instruction: stack the other block on the red block

 (b) Instruction: say yes
 (c) Instruction: wipe the desk

Figure 5: Illustration of execution on visual reasoning tasks in real-world environment. RoboPro presses buttons to express words (a), stacks object in an appropriate order based on visual properties (b), and chooses the best tool to wipe the desk (c).

desk, RoboPro will choose the appropriate tool (the sponge) among irrelevant objects, and grasp it to wipe water on the desk. We also provide detailed real-world setup in Appendix A.3.

4.4 EVALUATION ON GENERAL VQA TASKS

As mentioned in Sec. 3.3, RoboPro can meanwhile 454 exhibit broad visual understanding to perform general 455 visual question answering. To evaluate this ability, we 456 conduct experiments across a range of general VQA 457 tasks. We compare RoboPro with InstructBLIP based 458 on Vicuna 7B v1.1 (Dai et al., 2023), LLaVA-1.5 (Liu 459 et al., 2024b) and ShareGPT4V (Chen et al., 2023). We 460 evaluate the zero-shot performance of these models on 461 perception and reasoning tasks with VQAv2 (Goyal 462 et al., 2017), GQA (Hudson & Manning, 2019) and 463 TextVQA (Singh et al., 2019). As shown in Table 4, RoboPro can not only generate executable code for 464

Table 4: The zero-shot accuracy of RoboPro and the baselines on general VQA tasks.

Models	VQA ^{v2}	GQA	VQA ^T
InstructBLIP	-	49.2	50.1
LLaVA-1.5	78.5	62.0	58.2
ShareGPT4V	80.6	63.3	60.4
RoboPro	80.9	63.9	62.9

robotic control, but perform well on multimodal perception and reasoning tasks. The results indicate that our model demonstrates quite competitive performance on general VQA tasks compared to ShareGPT4V, which is the state-of-the-art vision-language model with similar model size. Experiments on general VQA tasks further confirm RoboPro's capability of comprehensive visual understanding, which is a key factor in its success of manipulation tasks.

- 470 471 4 5
- 472

4.5 ABLATION STUDY

We conduct extensive ablations to evaluate the effectiveness of Video2Code and the contributions of
individual components in RoboPro and Video2Code framework. Specifically, we conduct ablation
studies on the base LLM in RoboPro, as well as the Draft VLM and Code LLM in Video2Code. We
provide detailed ablation results in Appendix A.5.

477

478 Effectiveness of Video2Code. We compare our model trained with and without Video2Code on 479 manipulation and general VQA tasks. For a fair comparison, we only remove Video2Code from 480 the fine-tuning stage for the baseline, that is, the 115k runtime code data are excluded and only the 481 general vision language fine-tuning dataset is used during the SFT process, as described in Sec. 3.3. 482 The first two rows of Table 5 show the comparison of the two settings. It is found that the Video2Code 483 generated data have significantly improved the performance on both RLBench and LIBERO by a gain of 42.3% and 45.4%, respectively, which indicate Video2Code's efficacy in enhancing the ability 484 of skills utility and instruction following. Moreover, our model trained with such code data can also 485 bring slight improvement on general VQA tasks.

9

444

445

446 447 448

449

450 451

452 453

432

IIM	Video2Codo	Manip	ulation	Ge	neral V()A
	vide02Code	RLBench	LIBERO	VQA ^{v2}	GQA	VQA ^T
CodeQwen-1.5-7B	×	0.4	7.0	80.5	63.8	62.1
CodeQwen-1.5-7B	1	42.7	52.4	80.9	63.9	62.9
DeepSeek-Coder-6.7B	1	41.3	48.8	78.3	60.9	59.5

Table 5: Ablations of Video2Code and different base LLMs on manipulation and general VQA tasks.

Table 6: The selection of the Draft VLM and Code LLM for Video2Code.

Method	RLBench	LIBERO
MiniCPM-V + Gemini-1.5-Flash	8.0	21.7
Gemini-1.5-Flash + Gemini-1.5-Flash	22.7	31.7
Gemini-1.5-Flash + DeepSeek-Coder-V2	42.7	52.4

Choice of base LLM. We further compare the performance of RoboPro using different codedomain base LLMs. Specifically, we choose DeepSeek-Coder-6.7B-Instruct (Guo et al., 2024) and CodeQwen-1.5-7B-Chat (Bai et al., 2023) for comparison. As shown in Table 5, RoboPro trained on CodeQwen-1.5-7B-Chat outperforms the version trained on DeepSeeK-Coder-6.7B-Instruct on both manipulation and general VQA tasks. These results demonstrate that employing a more powerful base LLM for code generation task can consequently enhance performance in both tasks.

509 Choice of Draft VLM and Code LLM. The Draft VLM and Code LLM are key components in 510 the design of Video2Code. As stated in Sec. 3.2, we choose Gemini-1.5-Flash as Draft VLM and 511 DeepSeek-Coder-V2 as Code LLM for default configurations. To analyze how the choice of Draft 512 VLM and Code LLM effects the quality of runtime code data, we set three different combinations 513 of Draft VLM and Code LLM for data curation. We choose Gemini-1.5-Flash and a light-weight 514 VLM MiniCPM-V (Yao et al., 2024) for Draft VLM evaluation, while selecting a code domain LLM 515 DeepSeek-Coder-V2 and a general VLM Gemini-1.5-Flash for Code LLM evaluation. All other 516 settings are consistent with those in our main experiment. As shown in Table 6, enhanced visual 517 reasoning capabilities of the Draft VLM, along with stronger code synthesis abilities of the Code 518 LLM, both play a crucial role in curating high-quality runtime code data.

519 520 521

522

486

504

505

506

507

508

5 CONCLUSION AND FUTURE WORK

- 523 In this work, we propose RoboPro, a robotic foundation model, which perceives visual information 524 and follows free-form instructions to perform robotic manipulation in a zero-shot manner. To 525 address low efficiency and high cost for runtime code data synthesis, we propose Video2Code, a scalable and automatic data curation pipeline. Through extensive experiments, with assistance 526 of Video2Code, RoboPro achieves impressive generalization capability compared with training-527 based methods, and exhibits significant improvement on performance compared with other policy 528 code generation methods. These results indicate that incorporating procedural knowledge within 529 operational videos into training process will bring substantially enhanced understanding of skills (i.e., 530 API libraries) and free-form instructions. Beyond the scope of robotic manipulation tasks, policy 531 code generation methods also show potential in many other robotic applications (e.g., navigation). 532 In the future, we would like to expand our method to more application scenarios to provide more 533 comprehensive support for complex real-world robotic deployments.
- 534
- 535

536 REFERENCES

537

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.

540 Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances from human 541 videos as a versatile representation for robotics. In Proceedings of the IEEE/CVF Conference on 542 Computer Vision and Pattern Recognition, pp. 13778–13790, 2023. 543 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, 544 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, 546 Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin 547 Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng 548 Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, 549 Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. arXiv preprint arXiv:2309.16609, 550 2023. 551 552 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action 553 models transfer web knowledge to robotic control. arXiv preprint arXiv:2307.15818, 2023a. 554 555 Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, 556 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, 2023b. 558 559 Junting Chen, Yao Mu, Qiaojun Yu, Tianming Wei, Silang Wu, Zhecheng Yuan, Zhixuan Liang, Chao 560 Yang, Kaipeng Zhang, Wenqi Shao, et al. Roboscript: Code generation for free-form manipulation tasks across real and simulation. arXiv preprint arXiv:2402.14623, 2024. 561 562 Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua 563 Lin. Sharegpt4v: Improving large multi-modal models with better captions. arXiv preprint arXiv:2311.12793, 2023. 565 566 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, 567 Boyang Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose vision-568 language models with instruction tuning. In Advances in Neural Information Processing Systems, 569 2023. 570 Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan 571 Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal 572 language model. In International Conference on Machine Learning, pp. 8469–8488. PMLR, 2023. 573 574 Hao-Shu Fang, Chenxi Wang, Hongjie Fang, Minghao Gou, Jirong Liu, Hengxu Yan, Wenhai Liu, 575 Yichen Xie, and Cewu Lu. Anygrasp: Robust and efficient grasp perception in spatial and temporal 576 domains. IEEE Transactions on Robotics, 2023. 577 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vga 578 matter: Elevating the role of image understanding in visual question answering. In Proceedings of 579 the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6904–6913, 2017. 580 581 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, 582 Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming-the 583 rise of code intelligence. arXiv preprint arXiv:2401.14196, 2024. 584 Siyuan Huang, Zhengkai Jiang, Hao Dong, Yu Qiao, Peng Gao, and Hongsheng Li. Instruct2act: 585 Mapping multi-modality instructions to robotic actions with large language model. arXiv preprint 586 arXiv:2305.11176, 2023a. 587 588 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot 589 planners: Extracting actionable knowledge for embodied agents. In International Conference on 590 Machine Learning, pp. 9118–9147. PMLR, 2022. 591 Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: 592 Composable 3d value maps for robotic manipulation with language models. In Conference on Robot Learning, pp. 540-562. PMLR, 2023b.

594 595 596 597	Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. In <i>Conference on Robot Learning</i> , pp. 1769–1782. PMLR, 2023c.
598 599 600	Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 6700–6709, 2019.
601 602 603	Stephen James, Marc Freese, and Andrew J. Davison. Pyrep: Bringing v-rep to deep robot learning. arXiv preprint arXiv:1906.11176, 2019.
604 605	Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J. Davison. Rlbench: The robot learning benchmark & learning environment. <i>IEEE Robotics and Automation Letters</i> , 2020.
606 607 608	O. Khatib. A unified approach for motion and force control of robot manipulators: The operational space formulation. <i>IEEE Journal on Robotics and Automation</i> , 3(1):43–53, 1987.
609 610 611 612	Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. <i>arXiv preprint arXiv:2403.12945</i> , 2024.
613 614 615 616	Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. <i>arXiv preprint arXiv:2406.09246</i> , 2024.
617 618 619 620	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.
621 622 623	Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. <i>Advances in Neural Information</i> <i>Processing Systems</i> , 36, 2024a.
624 625 626	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 26296–26306, 2024b.
627 628 629 630	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. <i>arXiv preprint arXiv:2303.05499</i> , 2023.
631 632 633	Yao Mu, Junting Chen, Qinglong Zhang, Shoufa Chen, Qiaojun Yu, Chongjian Ge, Runjian Chen, Zhixuan Liang, Mengkang Hu, Chaofan Tao, et al. Robocodex: Multimodal code generation for robotic behavior synthesis. <i>arXiv preprint arXiv:2402.16117</i> , 2024.
634 635 636 637	Dantong Niu, Yuvan Sharma, Giscard Biamby, Jerome Quenum, Yutong Bai, Baifeng Shi, Trevor Darrell, and Roei Herzig. Llarva: Vision-action instruction tuning enhances robot learning. <i>arXiv</i> preprint arXiv:2406.11815, 2024.
638	OpenAI. Hello gpt-4o, May 2024. URL https://openai.com/index/hello-gpt-4o/.
639 640 641 642	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.
643 644 645	Eric Rohmer, Surya PN Singh, and Marc Freese. V-rep: A versatile and scalable robot simulation framework. In <i>2013 IEEE/RSJ international conference on intelligent robots and systems</i> , pp. 1321–1326. IEEE, 2013.
647	Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for

647 Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *Conference on Robot Learning*, pp. 785–799. PMLR, 2023.

648 Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus 649 Rohrbach. Towards vqa models that can read. In Proceedings of the IEEE Conference on Computer 650 Vision and Pattern Recognition, pp. 8317–8326, 2019. 651 Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter 652 Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using 653 large language models. In 2023 IEEE International Conference on Robotics and Automation 654 (ICRA), pp. 11523–11530. IEEE, 2023. 655 656 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu 657 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable 658 multimodal models. arXiv preprint arXiv:2312.11805, 2023. 659 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 660 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and 661 efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023. 662 663 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz 664 Kaiser, and Illia Polosukhin. Attention is all you need. Advances in Neural Information Processing 665 Systems, 30, 2017. 666 Sai H Vemprala, Rogerio Bonatti, Arthur Bucker, and Ashish Kapoor. Chatgpt for robotics: Design 667 principles and model abilities. IEEE Access, 2024. 668 669 Terry Winograd. Procedures as a representation for data in a computer program for understanding 670 natural language. 1971. 671 Mengdi Xu, Peide Huang, Wenhao Yu, Shiqi Liu, Xilun Zhang, Yaru Niu, Tingnan Zhang, Fei Xia, 672 Jie Tan, and Ding Zhao. Creative robot tool use with large language models. arXiv preprint 673 arXiv:2310.13065, 2023. 674 675 Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, 676 Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. arXiv preprint 677 arXiv:2408.01800, 2024. 678 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language 679 image pre-training. In Proceedings of the IEEE/CVF International Conference on Computer Vision, 680 pp. 11975–11986, 2023. 681 682 Junjie Zhang, Chenjia Bai, Haoran He, Zhigang Wang, Bin Zhao, Xiu Li, and Xuelong Li. SAM-683 E: leveraging visual foundation model with sequence imitation for embodied manipulation. In 684 International Conference on Machine Learning, 2024. 685 Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, 686 Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models 687 in code intelligence. arXiv preprint arXiv:2406.11931, 2024. 688 689 690 TASK DETAILS А 691 692 A.1 TASKS IN RLBENCH 693 694 RLBench is a simulation platform set in CoppelaSim (Rohmer et al., 2013) and interfaced through 695 PyRep (James et al., 2019). Robotic models control a 7-dof Franka Panda robot with a parallel gripper

REBENCH is a simulation platform set in Copperasim (Romner et al., 2013) and interfaced unough
 PyRep (James et al., 2019). Robotic models control a 7-dof Franka Panda robot with a parallel gripper
 to complete language-conditioned tasks. RoboPro is evaluated on 9 tasks from RLBench (James
 et al., 2020). Modification on these tasks is consistent with PerAct (Shridhar et al., 2023). Each
 task in RLBench is provided with several variations on language instructions describing the goal. In
 order to validate RoboPro's adaptation ability across various and vague instructions, we pop out an
 instruction from the language template list for each episode during evaluation instead of just using
 the first language template. Detailed descriptions and modification for each task in RLBench are

Push Buttons. Push down colored buttons in a specific order. The task has 20 different variances on the color of buttons, and three variances on the number of buttons to be manipulated. The success metric of this task is to push down specific buttons in correct order.

Close Jar. Put the lid on the table onto the jar with specific color. This task also has 20 different variations on the color of the jars. The success metric is that the lid is on the top of the target jar, and the gripper doesn't grasp anything.

710 Stack Blocks. Stack two to four blocks with specific color onto the green target area. There are 711 always two groups of four blocks with the same color, and this task has 20 variations on the color 712 of the blocks. The success metric has a further requirement that all stacked blocks inside the area 713 of a green platform beyond the original language instruction. We add target prompt to specify the 714 stacking area.

715

729

733

709

Open Drawer. Open specific drawer of a cabinet. there are three different variations on the position of the drawer: top, middle, and bottom. The success metric is a full extension of the target drawer joint. Before execution, we first adjust the gripper position to face the cabinet.

Stack Cups. Stack other two cups onto the cup with specific color. This task has 20 variations on the color of the cups. The success metric of this task is that the other cups are inside the target cup.

722
 723
 724
 725
 Sweep Dirt. Sweep dirt particles to the target dustpan. There are two dustpans specified as a tall dustpan and a short dustpan. The success metric of this task is that all 5 dirt particles are in the target dustpan. This task is modified by PerAct.

Slide Block. Slide the red cube in the scene to the target colored area. There are four areas with different color on each corner of the scene, and the cube cannot be picked up. The success metric is that the cube is inside the area with the target color, which is modified by PerAct.

Screw Bulb. Screw light bulb with the specified base onto the lamp base. There are two bulbs in
the scene at once, and the color of the holders have 20 different variations. The success metric is that
the bulb is inside the lamp stand.

Put in Board. Pick up the specified object and place it into the cupboard above. There are always 9 different objects on the table. The success rate is that the target object is in the cupboard.

736 737 A.2 TASKS IN LIBERO

738 In this section, we provide a detailed description of 8 tasks selected from the LIBERO-100 dataset. 739 Each task is associated with a specific language instruction, with the task ID and corresponding 740 instruction shown in Table 7. The tasks "Turn on Stove" and "Close Cabinet" are taken from LIBERO-741 90, which focuses on testing atomic skills and environmental understanding. The remaining tasks 742 are more complex, requiring multi-step execution, and are selected from LIBERO-10. These 8 743 tasks challenge RoboPro to comprehend diverse visual environments and follow extended language 744 instructions. As illustrated in Figure 6, the tasks encompass a wide range of robotic capabilities, including object selection, spatial reasoning, scene comprehension, and long-term execution. 745

- 746
- 747
- 748 749
- 750
- 751
- 752
- 753
- 754
- 755

Task ID	Task Instruction
Turn on Stove	turn on the stove
Close Cabinet	close the top drawer of the cabine
Put in Sauce	put both the alphabet soup and the tomato sauce in the baske
Put in Butter	put both the cream cheese box and the butter in the basket
Put in Cheese	put both the alphabet soup and the cream cheese box in the basket
Place Book	pick up the book and place it in the back compartment of the caddy
Boil Water	turn on the stove and put the moka pot on i
Identify Plate	put the white mug on the left plate and put the yellow and white mug on the right plate

Table 7: The manipulation tasks selected for the evaluation of zero-shot generalization on LIBERO.



Figure 6: Illustration of the selected tasks from LIBERO benchmark.

A.3 TASKS IN REAL-WORLD EXPERIMENTS

To validate the performance of RoboPro, the real-world experiments are implemented on a Franka Emika Panda robotic arm with a parallel jaw grip-per, as shown in Figure 7. We use an Intel Re-alSense D435i camera to provide RGB-D input signals under the camera-in-hand setting. Easy-handeye ROS package is used to calibrate the ex-trinsics of the camera frame with respect to the robot base frame. For robot control, we use the open-source frankapy package to send real-time position-control commands to robot after receiv-ing the control signals from RoboPro. During test time for each task, natural language instructions, extrinsic matrix, intrinsic matrix, current environ-ment observation in the form of RGB-D image,



Figure 7: The setup for real-world experiments.

and the low dimensional state of the robot are prepared for RoboPro to generate corresponding 6-DOF action trajectories. Examples of all 8 real-world tasks with natural language instructions are illustrated in Figure 8.



A.4 ILLUSTRATION OF ERROR CASES



Figure 9: Illustration of failure and success cases of different policy code generation methods in LIBERO and RLBench.

A.5 ADDITIONAL EXPERIMENT RESULTS

LLM	Video2Code	Push Buttons	Stack Blocks	Open Drawer	Close Jar	Stack Cups	Sweep Dirt	Slide Block	Screw Bulb	Put in Board	Av
CodeQwen-1.5-7B	X	0	0	0	0	0	0	0	4	0	0.
CodeQwen-1.5-7B	1	68	48	68	44	4	48	60	32	12	42
DeepSeek-Coder-6.7B	1	72	32	68	48	0	24	84	32	12	41

LLM	Video2Code	Turn on Stove	Close Cabinet	Put in Sauce	Put in Butter	Put in Cheese	Place Book	Boil Water	Identify Plate	Avg
CodeQwen-1.5-7B	X	0	43	0	0	13	0	0	0	7.0
CodeQwen-1.5-/B DeenSeek-Coder-6 7B		97 97	60 60	67 47	53 53	63 60	43 53	23	13 20	52.4 48 2
B PROMPT AND	API Impl	.EMEN	ΤΑΤΙΟΙ	٩S						
	Listing 1. A	n exam	ole of a	full pro	mnt in	Rohor)ro			
"""You're a vision la	nguage model	control	ling a	grippe:	r to co	mplete	manipu	latior	tasks.	
Combine the image code for the curr You have access to th	es you see wi rent scene. e following	ith the tools:	text ins	structi	ons to	genera	te det	ailed	and work	able
<pre>import numpy as np import torch import math</pre>										
#Perception Modules										
<pre>def get_obj_bbox(desc """get_the 2D boundin</pre>	ription: str gbox of all)->list	[bbox]:	escript	tion W	hen it	comes	to the	specif	ic
parts or orienta	tion of obje	cts, the	descrip	otion s	hould	be deta:	iled.	Like '	handle o	of
microwave', 'lef	t side of sh	elf'.								
veratu: TISt[bbox; ub	.nuarray]									
<pre>def get_best_grasp_po</pre>	s(grasp_bbox	: bbox):								
Return: grasp pose: P	е то grasp s ose"""	pecific	object.							
0 1 -1										
<pre>def get_place_pos(hol """Predict the place</pre>	der_bbox: bb pose for an	ox): obiect r	elative	to a h	nolder					
Args: holder_bbox: bb	ox of target	region	of the 1	holder	,					
Return: place_pose: P	ose"""									
<pre>def get_joint_axis(jo</pre>	int_object_n	ame: str):							
"""Get the joint dire	ction of an	object	oct have	ioin+	avie					
Return: joint_axis: n	p.ndarray"""	or onle	ct nave	JOTHC	dYID'					
dof monorate deduct	+h(ici-+ '	.	0.000							
<pre>dei generate_joint_pa """Generate a gripper</pre>	path of pos	s: np.nc es arour	d the i	o pen: b pint. d	open is	True w	hen ne	ed ope	n conta:	iner
around joint, Fa	lse when clos	se conta	iner.		-			1		
<pre>Keturn: path: list[Po """</pre>	sej									
<pre>def generate_slide_pa</pre>	th(target: O	ptional	[str] =]	None, d	lirecti	on: Opt	ional[np.nda	rray] =	Non
"""Generate path of p	oses to slid	e or pus	h objec	t to ta	arget o	r in sp	ecific	direc	tion.	
Args:		-	1.12							
target: The target lo direction: The direct	cation. If p ion vector t	rovided, o slide	'direc	tion' n ect alc	nust be ong. If	None.	ed, 't	arget'	must b	e Nc
500	-		J		0	1	, .	0		
Return: path: list[Po	se]									
def generate_sweep_pa	th(object: 0	ptional	[str] =]	None, t	arget:	Option	al[<mark>str</mark>] = No	ne, dire	ecti
Uptional[np.ndar """This function is d	rray] = None esigned to g): enerate	movemen	t paths	s for s	weeping	actio	ns usi	ng tools	5 S11
as sweepers, bro	poms. Grasp	the tool	before	sweepi	ng.		~~~~	401		
Args:	he arrest T	f ant t	Nora	the from			mf 0			
object: The object to target: The target ar	pe swept. I ea or locati ion vector f	I Set to on to su	None, veep tow	tne fur ards.]	iction If prov	will pe ided, '	riorm direct 'tar	a gene	ral swee	epin None

```
972
        .....
973
974
        def generate_wipe_path(region: str):
        """This function is designed to generate movement paths for wiping actions using tools such
975
            as towel, sponge. Grasp the tool before wiping.
976
        Args:
977
        region (str): region to be wiped or cleaned.
        Return: path: list[Pose]
978
        11.11.11
979
980
        def generate_pour_path(grasped object: str, target: str):
        """Generate gripper path of poses to pour liquid in grasped object to target.
981
        Return: path: list[Pose]
982
983
        def generate_press_pose(bbox):
984
        """Get best pose to press or push buttons."""
985
986
        #Action Modules
        def move_to_pose(Pose):
987
        """Move the gripper to pose."""
988
989
        def move_in_direction(direction: np.ndarray, distance: float):
        """Move the gripper in the given direction in a straight line by certain distance.
990
        .....
991
992
        def follow_way(path: List[Pose]):
        """Move the gripper to follow a path of poses."""
993
994
        def rotate(angle: float)
995
        """Rotate the gripper clockwise at certain degree while maintaining the original position."""
996
        def open_gripper():
997
        """Open the gripper to release the object, no args"""
998
        def close_gripper():
999
        """Close the gripper to grasp object, no args. Move to best grasp pose before close gripper.
1000
1001
1002
        Rules you have to follow:
1003
        #Directions: right: [0,1,0], left: [0,-1,0], upward or lift object: [0,0,1], forward or move
1004
            away: [1,0,0]
        #Please solve the following instruction step-by-step.
1005
        #You should ONLY implement the main() function and output in the Python-code style. Except
1006
            the code block, output fewer lines.
1007
        Begin to excecute the task:
1008
        #Instruction:
1009
1010
1011
                     Listing 2: An example of a full prompt in RoboPro with API renaming
1012
        """You're a vision language model controlling a gripper to complete manipulation tasks.
1013
             Combine the images you see with the text instructions to generate detailed and workable
             code for the current scene.
1014
        You have access to the following tools:
1015
1016
        import numpy as np
1017
        import torch
1018
        import math
1019
        #Perception Modules
1020
        def detect_bbox(description: str)->list[bbox]:
1021
         ""get the 2D boundingbox of all objects match description. When it comes to the specific
1022
             parts or orientation of objects, the description should be detailed. Like 'handle of
             microwave', 'left side of shelf'.
1023
        Return: list[bbox: np.ndarray]"""
1024
1025
        def generate_obj_grasp_pos(grasp_bbox: bbox):
```

"""get best grasp pose to grasp specific object.

1026 Return: grasp_pose: Pose""" 1027 1028 def best_place_locator(holder_bbox: bbox): """Predict the place pose for an object relative to a holder 1029 Args: holder_bbox: bbox of target region of the holder. 1030 Return: place_pose: Pose""" 1031 def find_axis_of_joint(joint_object_name: str): 1032 """Get the joint direction of an object 1033 Args: joint_object_name: the name of object have joint axis. 1034 Return: joint_axis: np.ndarray""" 1035 def map_joint_path(joint_axis: np.ndarray, open: bool): 1036 """Generate a gripper path of poses around the joint. open is True when need open container around joint, False when close container. 1037 Return: path: list[Pose] 1038 1039 def build_slide_path(target: Optional[str] = None, direction: Optional[np.ndarray] = None): 1040 """Generate path of poses to slide or push object to target or in specific direction. 1041 Args: 1042 target: The target location. If provided, 'direction' must be None. 1043 direction: The direction vector to slide the object along. If provided, 'target' must be None. 1044 Return: path: list[Pose] 1045 11.11.11 1046 def sweep_motion_path(object: Optional[str] = None, target: Optional[str] = None, direction: 1047 Optional[np.ndarray] = None): 1048 """This function is designed to generate movement paths for sweeping actions using tools such 1049 as sweepers, brooms. Grasp the tool before sweeping. Args: 1050 object: The object to be swept. If set to None, the function will perform a general sweeping. 1051 target: The target area or location to sweep towards. If provided, 'direction' must be None. direction: The direction vector for the sweeping motion. If provided, 'target' must be None. 1052 Return: path: list[Pose] 1053 1054 1055 def create_wipe_path(region: str): """This function is designed to generate movement paths for wiping actions using tools such 1056 as towel, sponge. Grasp the tool before wiping. 1057 Args: 1058 region (str): region to be wiped or cleaned. Return: path: list[Pose] 1059 1060 1061 def pour_path_mapper(grasped object: str, target: str): """Generate gripper path of poses to pour liquid in grasped object to target. 1062 Return: path: list[Pose] 1063 1064 def best_press_pos(bbox): 1065 """Get best pose to press or push buttons.""" 1066 #Action Modules 1067 def relocate_to_pose(Pose): 1068 """Move the gripper to pose.""" 1069 1070 def reach_in_direction(direction: np.ndarray, distance: float): ""Move the gripper in the given direction in a straight line by certain distance. 1071 1072 def follow_path(path: List[Pose]): 1073 """Move the gripper to follow a path of poses.""" 1074 1075 def spin_gripper(angle: float) 1076 """Rotate the gripper clockwise at certain degree while maintaining the original position.""" 1077 def open claw(): 1078 """Open the gripper to release the object, no args""" 1079 def clamp_gripper():

"	""Close the gripper to grasp object, no args. Move to best grasp pose before close gripper.
- R # #	<pre>ules you have to follow: Directions: right: [0,1,0], left: [0,-1,0], upward or lift object: [0,0,1], forward or move away: [1,0,0] Please solve the following instruction step-by-step. You should ONLY implement the main() function and output in the Python-code style. Except the code block, output fewer lines.</pre>
B #	egin to excecute the task: Instruction:
_	Listing 3: An example of a full prompt in RoboPro with API refactoring
Ч Ч	<pre>""You're a vision language model controlling a gripper to complete manipulation tasks. Combine the images you see with the text instructions to generate detailed and workable code for the current scene. 'ou have access to the following tools: ""</pre>
i i	mport numpy as np mport torch mport math
#	Percention APIs
# d " A R	<pre>reception APIS ef get_obj_bbox(description: str) -> list[np.ndarray]: ""Get the 2D bounding box of all objects that match the description. The description should be detailed when it comes to specific parts or orientations of objects, such as 'handle of microwave' or 'left side of shelf'. rgs:description (str): The description of the objects to find. eturns:list[np.ndarray]: A list of bounding boxes for the objects matching the description. </pre>
d "	ef get_joint_axis(joint_object_bbox: np.ndarray): ""Get the joint direction of an object
A R	rgs: joint_object_name: the name of object have joint axis. .eturn: joint_axis: np.ndarray"""
#	Control APIs
d " P R	<pre>ief get_best_grasp_pos(grasp_bbox: np.ndarray): ""Calculate the best grasp pose to grasp a specific object. arameters: grasp_bbox (np.ndarray): The bounding box of the object to grasp. eturn: Pose: The best grasp pose for the given object."""</pre>
d "	ef get_place_pos(holder_bbox: np.ndarray): ""Predict the place pose for an object relative to a holder.
P R	arameters: holder_bbox (np.ndarray): The bounding box of the target region of the holder. eturn: Pose: The predicted place pose for the given object."""
d	ef generate_joint_path(joint_axis: np.ndarray, open: bool) -> list[Pose]:
" P	""Generate a gripper path of poses around the joint. arameters: joint_axis (np.ndarray): The axis of the joint. open (bool): True if the container needs to be opened around the joint. False if it needs to be closed
R	eturns: list[Pose]: The generated path of poses around the joint."""
d	ef generate slide nath(target bhoy: nn ndarray) -> list[Pose].
"	""Generate a path of poses to slide or push an object to a target or in a specific direction.
P R	arameters: target_bbox (np.ndarray): bbox of the target location. eturns: List[Pose]: The generated path of poses."""
d	ef generate_sweep_path(target_bbox: np.ndarray) -> List[Pose]:
11	""Generate movement paths for sweeping actions using tools such as sweepers or brooms. Grasp
P R	tne tool before sweeping. arameters: target_bbox (np.ndarray): The target area or location to sweep towards. eturns: List[Pose]: The generated path of poses for the sweeping action."""

def generate_wipe_path(region_bbox: np.ndarray) -> List[Pose]:

1134 """Generate movement paths for wiping actions using tools such as towels or sponges. Grasp 1135 the tool before wiping. 1136 Parameters: region_bbox (np.ndarray): The region to be wiped or cleaned. Return: List[Pose]: The generated path of poses for the wiping action.""" 1137 1138 def generate_pour_path(grasped_object: str, target_bbox: np.ndarray) -> List[Pose]: """Generate a gripper path of poses to pour liquid from a grasped object to a target. 1139 Parameters: grasped_object (str): The object being grasped that contains the liquid. 1140 target_bbox (np.ndarray): The bounding box of the target area where the liquid will be 1141 poured. 1142 Returns: List [Pose]: The generated path of poses for the pouring action.""" 1143 def generate_press_pose(bbox: np.ndarray) -> Pose: 1144 """Get the best pose to press or push buttons. Parameters: bbox (BBox): The bounding box of the button or area to be pressed. 1145 Return: Pose: The best pose for pressing or pushing the button.""" 1146 1147 def move_to_pose(pose: Pose): """Move the gripper to the specified pose. 1148 Parameters: pose (Pose): The target pose to move the gripper to.""" 1149 1150 def move_in_direction(direction: np.ndarray, distance: float): 1151 """Move the gripper in the given direction in a straight line by a certain distance. Parameters: direction (np.ndarray): The direction vector to move the gripper along. distance 1152 (float): The distance to move the gripper.""" 1153 def follow_way(path: List[Pose]) -> None: 1154 ""Move the gripper to follow a path of poses. 1155 Parameters: path (List[Pose]): The list of poses that defines the path to follow.""" 1156 1157 def rotate(angle: float) -> None: """ Rotate the gripper clockwise by a certain angle while maintaining the original position. 1158 Parameters: angle (float): The angle in degrees to rotate the gripper.""" 1159 def open_gripper() -> None: 1160 """Open the gripper to release the object. 1161 Parameters: None 1162 Returns: None"" 1163 def close_gripper() -> None: 1164 """Close the gripper to grasp an object. Move to the best grasp pose before closing the 1165 gripper. Parameters: None 1166 Returns: None""" 1167 1168 _____ Rules you have to follow: 1169 #Directions: right: [0,1,0], left: [0,-1,0], upward or lift object: [0,0,1], forward or move 1170 away: [1,0,0] 1171 #Please solve the following instruction step-by-step. #You should ONLY implement the main() function and output in the Python-code style. Except 1172 the code block, output fewer lines. 1173 1174 Begin to excecute the task: #Instruction: 1175 1176 1177 1178 1179 1180 1181 1182 1183 1184 1185 1186