

PCA-Bench: Evaluating Multimodal Large Language Models in Perception-Cognition-Action Chain

Liang Chen¹, Yichi Zhang¹, Shuhuai Ren¹, Haozhe Zhao¹, Zefan Cai¹, Yuchi Wang¹, Peiyi Wang¹, Xiangdi Meng¹, Tianyu Liu², Baobao Chang¹

¹ National Key Laboratory for Multimedia Information Processing, Peking University

² Alibaba Group

{leo.liang.chen, yczhang, shuhuai_ren}@stu.pku.edu.cn
tianyu0421@alibaba-inc.com, chbb@pku.edu.cn

PCA-EVAL 😊 PCA-Bench-V1

Abstract

We present PCA-Bench, a multimodal decision-making benchmark for evaluating the integrated capabilities of Multimodal Large Language Models (MLLMs). Departing from previous benchmarks focusing on simplistic tasks and individual model capability, PCA-Bench introduces three complex scenarios: autonomous driving, domestic robotics, and open-world games. Given task instructions and diverse contexts, the model is required to seamlessly integrate multiple capabilities of Perception, Cognition, and Action in a reasoning chain to make accurate decisions. Moreover, PCA-Bench features error localization capabilities, scrutinizing model inaccuracies in areas such as perception, knowledge, or reasoning. This enhances the reliability of deploying MLLMs. To balance accuracy and efficiency in evaluation, we propose PCA-Eval, an automatic evaluation protocol, and assess 10 prevalent MLLMs. The results reveal significant performance disparities between open-source models and powerful proprietary models like GPT-4 Vision. To address this, we introduce Embodied-Instruction-Evolution (EIE), an automatic framework for synthesizing instruction tuning examples in multimodal embodied environments. EIE generates 7,510 training examples in PCA-Bench and enhances the performance of open-source MLLMs, occasionally surpassing GPT-4 Vision (+3% in decision accuracy), thereby validating the effectiveness of EIE. Our findings suggest that robust MLLMs like GPT4-Vision show promise for decision-making in embodied agents, opening new avenues for MLLM research.

1 Introduction

Multimodal Large Language Models (MLLMs) have demonstrated remarkable capabilities in tackling complex tasks that necessitate a chain of integrated skills, including visual perception, world knowledge, reasoning, action, and more (OpenAI,

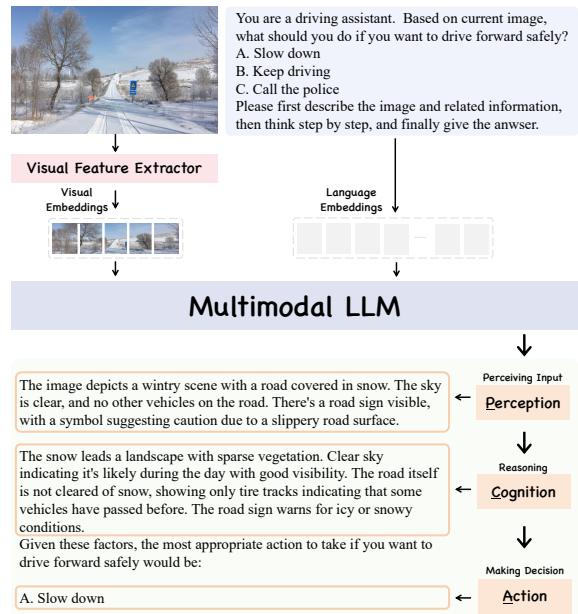


Figure 1: Example of decision making with MLLMs in the Perception-Cognition-Action Chain.

2023; Dai et al., 2023a; Liu et al., 2023b; Li et al., 2023c; Zhao et al., 2023).

However, current MLLM benchmarks often evaluate these capabilities individually (Fu et al., 2023; Liu et al., 2023e), overlooking the significant integrated potential that Large Language Models (LLMs) contribute to multimodal models. While some benchmarks like MMMU (Yue et al., 2023) and MathVista (Lu et al., 2023a) require abilities from both the vision and language part, they lack error localization techniques beyond accuracy assessments. This complicates identifying which part of the MLLM malfunctioned when making mistakes—whether it was the visual or the language component—and determines which aspect requires enhancement to enhance overall performance.

To address the challenges of insufficient integrated benchmarking and error localization problems, we introduce **PCA-Bench**. It arises with

MLLM’s applications in embodied AI and decision making, where models called agents need to first process multimodal observation from different environments, reason with the current situation and goal, and finally make an action from a given action space. The abilities in the complex decision making process can be abstracted to Perception, Cognition and Action according to the Perception-Action loop (Fuster, 2004) in Cognitive Science, a fundamental concept that describes how organisms process sensory information to interact with their environment through actions, offering a comprehensive framework for assessment. Figure 1 shows how MLLMs make decisions in the PCA chain.

The instances in PCA-Bench are from three influential domains in embodied decision-making: autonomous driving, domestic robotics, and open-world gaming. As shown in Figure 2, each instance is annotated by human annotators with a 6-element tuple: $\langle \text{image}, \text{question}, \text{action candidates}, \text{answer}, \text{reason}, \text{key concept} \rangle$. The last three elements serve as anchors for error localization for Action, Cognition and Perception, correspondingly.

PCA-Eval is an anchor-based evaluation protocol, designed to automatically conduct error localization utilizing the powerful semantic parsing ability of LLMs and the anchor information in data annotation. In the past, such localization was both labor-intensive and time-consuming. PCA-Eval with strong LLMs like GPT4 demonstrates a strong kappa correlation with human assessments, reaching 0.8+ average kappa coefficients for perception, cognition, and action scores. The anchor-based evaluation provides the LLMs with groundtruth answers for each sub-score, preventing the systematic bias of LLM evaluators, such as position bias (Wang et al., 2023b; Zheng et al., 2023) in the pair-wise evaluation and verbosity bias (Zheng et al., 2023) in simple preference evaluation. We also compared open state-of-the-art LLMs in PCA-Eval. Though they lag behind close ones in alignment with human assessments, we see large improvement when the model scales up. We believe that with specific training for error localization and improved general ability of open LLMs in the future, they would be more suitable evaluation tools for the reproducible and transparent characteristics.

Aiming at scaling up PCA-Bench, using LLM to synthesize training examples is an increasingly popular method for enhancing models without additional human involvement. We expand this approach to generate more samples following the

Autonomous Driving Image: 	Question: Based on current image, what is the best action to take when you are driving on the highway? Action candidates: ["Slow down", "Keep driving", "Stop the car", "Change to other lane"] Answer: Keep driving Reason: There is no other car or obstacle on the highway so it is safe to keep driving. Key Concept: Clear Road
Domestic Robot Image: 	Question: Fill the bathtub with water. Action candidates: ["Go to the bathroom", "Find the bathtub", "Get in the tub", "Switch on the bathtub faucet"] Answer: switch on the bathtub faucet Reason: You are already in the bathroom and there is bathtub in front of you. To fill the bathtub with water, you need to switch on the faucet of the bathtub. Key Concept: Bathroom, Bathtub
Open World Game Image: 	Question: Craft a glass bottle. Action candidates: ["Craft glass bottle", "Find wood", "Craft crafting table"] Answer: Find wood Reason: To craft a glass bottle, you need 3 glass blocks. You have enough glass to make the bottle, but you don't have a crafting table to craft it. So you need to find wood to craft one. Key Concept: Have glass, No crafting table

Figure 2: Instances of PCA-Bench in 3 domains.

PCA guideline. Unlike text-based instruction generation methods like Self-Instruct (Wang et al., 2023c), generating instructions in embodied environments poses distinct challenges. It demands not only the creation of textual instructions but also the generation of corresponding precise observations. To address these challenges, we propose **Embodied Instruction Evolution (EIE)**, which integrates external environments with LLMs, thereby extending the LLMs’ ability to data synthesize across various embodied environments, contributing to 7,510 training data in PCA-Bench.

We conduct comprehensive experiments and analysis on PCA-Bench, our findings are summarized as follows:

1. Visual perception and reasoning with world knowledge are two core abilities for an MLLM to make correct decisions in PCA-Bench. GPT4-Vision shows strong zero-shot cross-modal reasoning ability for embodied decision-making tasks, surpassing open-source MLLMs and even Tool-Using LLM-agent.

2. EIE could generate training samples significantly enhancing the performance of open-source MLLMs (surpassing GPT-4V at some scores), validating the effectiveness of the method.

3. PCA-Eval serves as a good error locator. Above the high average kappa coefficient (0.8+) with human assessments and its ability to pinpoint the error source, it can effectively distinguish whether a model’s correct decisions are fluky or through genuine understanding. This leads to a better ensemble metric for MLLM evaluation named Genuine PCA Score.

2 PCA-Bench

2.1 Problem Definition

Multimodal decision-making problems are commonly formalized with a partially observable Markov decision process. For MLLMs \mathcal{F} tested in PCA-Bench, we care about given the multi-modal observation $o \in O$, the goal description g , a subset of candidates actions $A_C \subseteq A$, whether the model could make correct action $a \in A_C$ and give proper reasoning process r .

$$\mathcal{F}(g, o, A_C) = (a, r) \quad (1)$$

As shown in Figure 2, each instance in the benchmark is a 6-element tuple: $\langle \text{image}, \text{question}, \text{action candidates}, \text{answer}, \text{reason}, \text{key concept} \rangle$. The image is collected from various embodied environments, including transportation scenes, house-keeper environments, and Minecraft. Questions, action candidates, and answers are derived from real tasks within the corresponding environment. The reasons explain why the answer is the best choice for the current image, while the key concept highlights the most question-related aspect of the image.

Unlike traditional visual question-answering datasets that emphasize visual perception (e.g., VQA (Goyal et al., 2017)) or visual reasoning (e.g., NLVR (Suhr et al., 2017)), PCA-Bench mandates accurate observation perception, complex task decomposition, and understanding the outcomes of various actions simultaneously. Compared to embodied simulation environments such as ALFRED (Shridhar et al., 2020) and Minedojo (Fan et al., 2022), PCA-Bench stands out for its focus on high-level actions, proving to be more effective for evaluating MLLMs. This is because high-level actions, which can be readily translated or programmed into low-level actions within their respective domains, are inherently more accessible to LLMs. The high-level actions are more comprehensible for LLMs than the direct low-level actions like action vectors in the simulation environments because (1) the high-level actions are in the form of natural languages, making it easier for LLMs to understand the meaning and connect with world knowledge. (2) LLMs are not grounded with low-level actions during the pretraining or finetuning stage, making it hard for LLMs to understand the consequences of executing an action.

To answer a question in PCA-Bench, the agent

must possess the following abilities: (1) **Perception**: Accurately identify the concept related to the question within the image; (2) **Cognition**: Engage in reasoning based on image perception and worldly knowledge; (3) **Action**: Comprehend the potential actions, selecting the one that best aligns with the outcome of the reasoning process. A deficiency in any of these abilities would possibly result in an incorrect answer, posing a significant challenge to the more integrated capabilities of MLLMs.

2.2 PCA-Eval

For each instance, we prompt the model to deliver an answer comprising a reasoning process r , and a final action a , represented as $\langle r, a \rangle$. By comparing the model prediction with the ground truth answer, we can obtain a fine-grained diagnosis of the decision making process as follows:

Perception Score (P-Score) measures the model’s accuracy in perceiving the observation. It is computed based on whether the agent’s reasoning process r includes the key concept of the instance. A score of 1 is assigned if at least one question-related key concept is described by the agent; otherwise, it is 0. For the top example in Figure 2, the agent should output “clear road” or “no car visible” or other semantically equivalent concepts in its description of the image to get the perception score.

Parsing the model’s output and determining whether it entails the key concept using shallow features of the sentence is not trivial. We leverage LLM to conduct entailment detection, which turns out to have a high alignment with human judgment.

Cognition Score (C-Score) assesses the model’s ability to reason, comprehend, and make informed decisions based on the perceived input data and world knowledge. The score is 1 if the reasoning process is correct, otherwise the score is 0. For the instance in Figure 2, the agent should link the “clear road” to the action “keep driving” based on transportation commonsense to get the score.

Action Score (A-Score) measures the model’s ability to generate appropriate and effective responses or actions based on the perceived input data and the cognitive understanding of the context. The score is assigned a value of 1 if the agent selects the correct action; otherwise, the score is set to 0.

2.3 Automatic Evaluation

Recent advancements have seen researchers harnessing powerful LLMs for the evaluation of the

output of language models. Studies have revealed that the outcomes from LLMs could exhibit remarkable alignment with human judgments (Zheng et al., 2023; Wang et al., 2023b,a). In our investigation, we employed GPT-4 to automatically evaluate perception, cognition, and action scores based on the model’s outputs. Our findings underscore a significant agreement between GPT-4 scoring and human evaluation results. This is substantiated by Cohen-Kappa coefficients of 0.71, 0.82, and 0.94 for perception, cognition, and action evaluations, respectively. Experiments of human evaluation and comparison of open LLMs are in section 4.1. For a detailed description of our evaluation tool, kindly refer to Appendix D.

2.4 Benchmark Dataset Overview

For the test set, the examples are written by 3 human experts for each domain. There are no overlapped environmental observations between the training and test sets. The details of the human annotation pipeline can be found in Appendix B. We introduce the three domains encompassed by our dataset as follows:

Autonomous Driving. In the autonomous driving domain, instances are derived from real-world transportation scenes, which requires the agent to have particular abilities such as traffic sign recognition, obstacle detection, and decision-making at intersections. The dataset aims to evaluate an agent’s ability to perceive and interpret visual information while making safe and efficient driving decisions. The images are collected from TT100K (Zhu et al., 2016) dataset and annotators are instructed to propose an image-conditioned question that is grounded with real actions of vehicles.

Domestic Robot. The domestic assistance domain features instances from the ALFRED (Shridhar et al., 2020; Kolve et al., 2017) environment, which simulates a housekeeper robot performing tasks within a household setting. These tasks may include object manipulation, navigation, and interaction with various appliances. The environment assesses an agent’s ability to understand and execute complex instructions while navigating and interacting with a dynamic environment. Annotators are asked to select one image from the randomly generated scenes in the environment, propose a question related to the items on the scene, and annotate the full information of the instance.

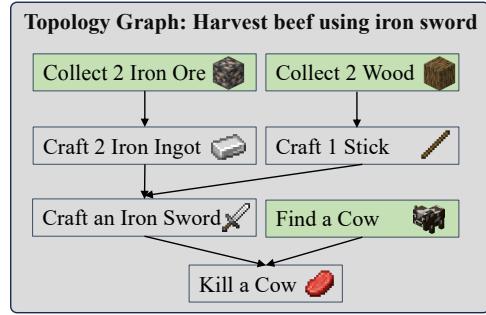


Figure 3: Illustration of task topology graph. Events in green represent the leaf nodes of the graph.

Open-World Game. In the open-world game domain, instances are sourced from the Minecraft environment, where agents are tasked with exploring, crafting, and surviving in a procedurally generated world. This dataset evaluates an agent’s ability to reason and plan actions within a complex, open-ended environment, which often requires long-term strategizing and adaptability. Annotators receive predefined tasks from MineDojo (Fan et al., 2022) as a reference during the task generation phase. For each task, we instruct the annotator to sketch a task topology graph, exemplified in Figure 3. The task should be completed under the topological order of the graph, where the event located in the leaf nodes should be finished first. Each node in the task topology graph can be viewed as a step in the sequential decision. We list the in-domain task distribution in Appendix A.

2.5 Embodied Instruction Evolution

The PCA-Bench benchmark also includes subset of automatic generated samples by Embodied Instruction Evolution(EIE), which is used as training set in our experiment.

The annotation of PCA-Bench examples is a labor-intensive task. As illustrated in Figure 4, we introduce Embodied Instruction Evolution (EIE), a method for automatically augmenting examples in the PCA-Bench format using Large Language Models, such as ChatGPT. This process involves four key steps:

1) Setup of Programmable Interface: Establish a programmable interface with a corresponding template, ensuring that observations in the embodied environment can be generated based on specific parameters.

2) Generation of Seed Tasks: Create initial seed tasks for each environment. These tasks are representative of the general challenges an agent

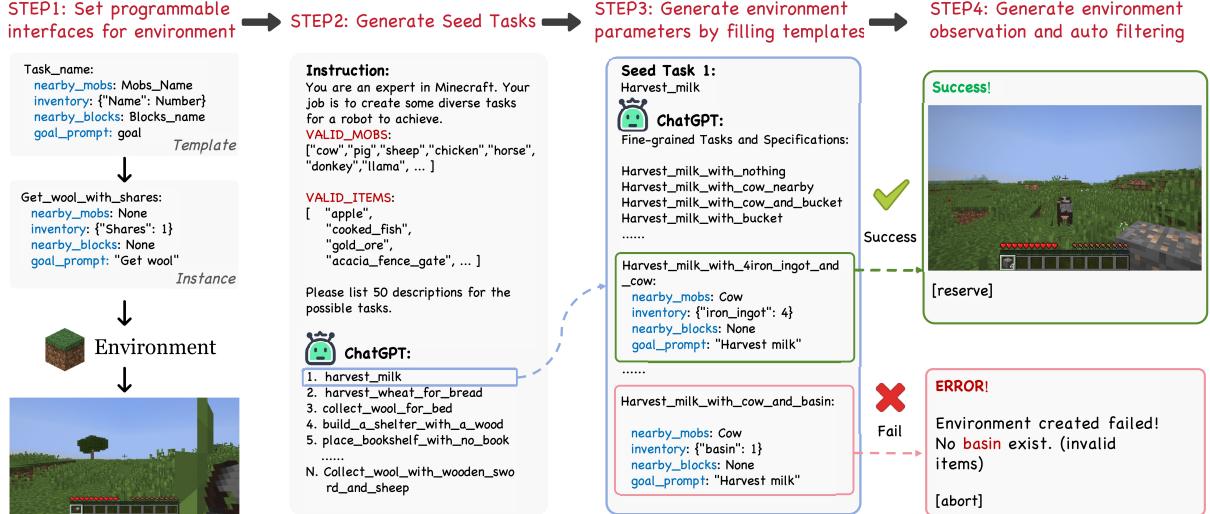


Figure 4: Pipeline of the Embodied Instruction Evolution method.

might encounter. We provide ChatGPT with sample tasks and enable it to generate additional seed tasks.

3) Task Specification and Template Filling:

For each seed task, we instruct ChatGPT to break down the task into multiple subtasks, following its event topology graph (as seen in Figure 3). This approach mimics the multi-step decision-making process. After determining the subtask names, we use the LLM to populate the environment parameter templates created in Step 1 for each subtask.

4) Observation Generation and Filtering:

Generate observations for the environment and implement an automatic process to filter out invalid instances. The filled templates may contain errors, such as incorrect creature names or impossible items, leading to errors during environment creation. When such errors occur, the affected templates are automatically filtered out. For domains without programmable environments (autonomous driving), step 1 and step 4 are not needed, we collect real traffic images and utilize GPT4-Vision to generate seed task based on the image content.

EIE leverages the capabilities of Large Language Models to reduce manual labor and improve the diversity and scalability of PCA-Bench.

3 Experiments

3.1 Tracks

Zero Shot End-to-End. The test set of PCA-Bench serves as an effective tool for comparing the embodied decision-making and cross-modal reasoning capabilities of various Multimodal Lan-

guage Learning Models (MLLMs). In this evaluation, the same images and prompts are provided to each model under test. Additionally, to address the challenge of perceiving certain non-visual information from images, details such as “items in hand” and “items in inventory”, particularly relevant in domestic and gaming domains, are directly included in the question prompts.

In our analysis, we benchmark the performance of the most recently open-sourced models, including LLaVA1.5 and Qwen-VL-Chat, as well as the API-only GPT4-V model. All models are evaluated using their default inference configurations to ensure a fair and standardized comparison.

Finetuning with EIE. In this track, we extend the capabilities of open-source MLLMs by finetuning them with the training set generated through our Embodied Instruction Evolution (EIE) method. After the fine-tuning process, these trained models are subjected to the test set of PCA-Bench. We finetune the LLaVA-7b/13b, MMICL and Qwen-VL-Chat models on the training set for 5 epochs. The training details are in Appendix E.

Zero Shot Modality Conversion. In this track, we introduce and compare a new baseline, termed HOLMES, which utilizes LLM without multi-modal perception capabilities. Instead, HOLMES relies on modality conversion APIs for embodied decision-making processes. Within the HOLMES framework, the LLM must continuously invoke various APIs, retrieving and processing return information about the environment. The HOLMES method is illustrated in Figure 7 from Appendix.

Model	Size	Traffic			Domestic			Game			Average		
		P	C	A	P	C	A	P	C	A	P	C	A
MiniGPT4 (Zhu et al., 2023) [†]	7B	0.45	0.37	0.48	0.81	0.38	0.38	0.38	0.14	0.27	0.55	0.30	0.38
LLaVA1.5 (Liu et al., 2023b) [†]	7B	0.44	0.44	0.53	0.92	0.48	0.44	0.8	0.35	0.39	0.72	0.42	0.45
Qwen-VL-Chat (Bai et al., 2023) [†]	7B	0.53	0.36	0.62	0.77	0.41	0.44	0.39	0.18	0.25	0.56	0.33	0.44
MiniGPT4 (Zhu et al., 2023) [†]	13B	0.41	0.37	0.5	0.85	0.35	0.33	0.41	0.22	0.33	0.56	0.31	0.39
InstructBLIP (Dai et al., 2023b) [†]	13B	0.36	0.41	0.42	0.90	0.44	0.39	0.33	0.25	0.24	0.53	0.37	0.35
MMICL (Zhao et al., 2023) [†]	13B	0.31	0.49	0.47	0.81	0.3	0.33	0.41	0.18	0.27	0.51	0.32	0.36
SPHINX-v1 (Lin et al., 2023) [†]	13B	0.46	0.48	0.61	0.95	0.55	0.31	0.71	0.35	0.43	0.71	0.46	0.45
LLaVA1.5 (Liu et al., 2023b) [†]	13B	0.49	0.56	0.61	0.95	0.62	0.46	0.74	0.45	0.51	0.73	0.54	0.53
Qwen-VL-Chat-PLUS (Bai et al., 2023) [‡]	UNK	0.57	0.56	0.65	0.86	0.44	0.43	0.68	0.47	0.49	0.70	0.49	0.52
GPT-4V (OpenAI, 2023) [‡]	UNK	0.73	0.72	0.74	0.96	0.66	0.62	0.88	0.72	0.69	0.86	0.7	0.68

Table 1: Zero Shot results on the full test set of PCA-Bench. Highest scores in each line are **bold** while second highest scores are underlined. Models with [†] are fully open-source. Models with [‡] only provide API to access. P, C, and A represent Perception, Cognition, and Action Scores, respectively.

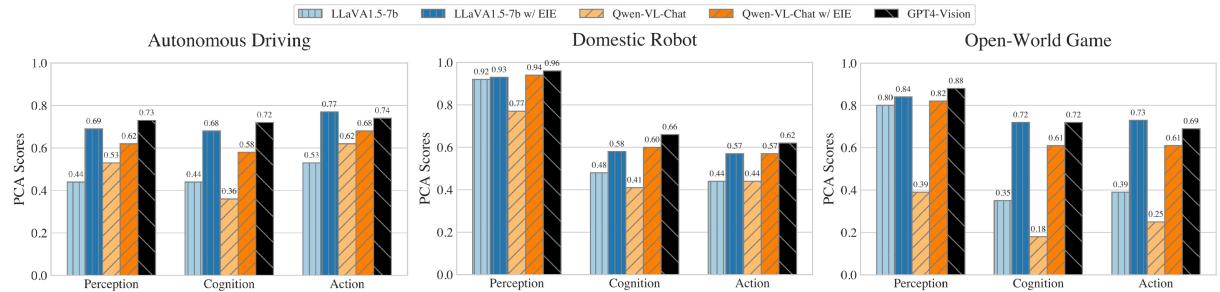


Figure 5: Performance comparsion between models' zero-shot results and models' finetuned results with the data generated by Embodied-Instruct-Evolution (EIE) method. EIE improves the performance on all domains for both LLaVA1.5-7b and Qwen-VL-Chat models. Results for LLaVA1.5-13B and MMICL are in Figure 13 from appendix.

We evaluate two LLMs in this track: ChatGPT-3.5-Turbo and GPT-4-0613, comparing their performances against the advanced GPT-4-Vision. Implementation details of the HOLMES framework and the APIs are provided in Appendix C.

3.2 Evaluation and Metrics

We use our PCA-Eval evaluation tool proposed in Section 2.3 to automatically assess the output of different models through three lenses: perception (P-Score), cognition (C-Score), and action (A-Score).

3.3 Main Results

Zero Shot Results. The results of the zero-shot end-to-end track are shown in Table 1. Among all MLLMs, GPT4-V, outperforms existing open-source models by achieving the highest scores of 0.86, 0.7, and 0.68 in the perception, cognition, and action dimensions respectively. This performance represents a 15% action score improvement over its strongest open-source counterpart, LLaVA1.5-13B. The impressive performance of GPT4-V is primarily attributed to its exceptional ability to perceive visual information across different domains and the world knowledge in the language model,

particularly in the challenging game domain.

Impact of Finetuning with EIE. The results of the fine-tuning track are illustrated in Figure 5. Our EIE method has been found to significantly enhance the general decision-making abilities of various models, encompassing perception, cognition, and action. Notably, it has led to an average increase of 0.24 and 0.19 in action scores for the LLaVA1.5-7b and Qwen-VL-Chat models, respectively. Results for LLaVA1.5-13b and MMICL are illustrated in Figure 13, also showing improved performance when trained with EIE. We note that there exist reasoning or perception errors in some of the generated sample due to the hallucination problem of LLM generated content, however they do not influence the overall performance. In some cases, these sub-scores have matched or even surpassed those of the GPT4-V model, demonstrating the potential of the EIE to scale up and apply to different environments.

Comparison Between End-to-End and Modality Conversion Method

In the zero-shot modality conversion track, we conduct an analysis and comparison of the outputs generated by the End2End

Method	Model	Traffic			Domestic			Game			Average		
		P	C	A	P	C	A	P	C	A	P	C	A
End-to-End	GPT-4V	0.75	0.73	0.78	0.81	0.69	0.67	0.95	0.79	0.77	0.84	0.74	0.74
HOLMES	ChatGPT	0.75	0.68	0.66	0.88	0.52	0.50	0.78	0.40	0.36	0.80	0.53	0.51
	GPT4	0.87	0.82	0.82	0.85	0.61	0.56	0.91	0.77	0.74	0.88	0.73	0.71

Table 2: Comparison between End-to-End (MLLM) and HOLMES (LLM+API) methods on a subset of PCA-Bench with API annotation.

method with GPT4-V, as well as the HOLMES method with GPT4 and ChatGPT-3.5 in Table 2.

The results show that the HOLMES system based on GPT4 achieves 0.71 Action Score, which is on par with GPT4-V’s performance (0.74). This indicates that, overall, the HOLMES system is able to accurately understand the task goal, split the larger goal into multiple smaller steps, and correctly invoke the relevant APIs to accomplish each step. Specifically, the HOLMES system based on GPT4 can recognize the key concepts in a task, and perceive the state and environment of these concepts through the results returned by APIs. Consequently, the system achieves an average Perception Score of 0.88, which even outperforms GPT4-V’s 0.84. However, compared to End2End methods, HOLMES relies on multi-step reasoning for the final decision, in which reasoning errors tend to accumulate, and thus achieves a lower Cognition Score in both Domestic and Game domains.

On the other hand, we also find that the End2End method effectively mitigates information loss during the modality conversion process. As illustrated in Figure 8 from Appendix, an image depicts a road with several nearby cars. GPT4-V is capable of discerning that the street is not crowded, thereby suggesting that the driver can continue driving.

Conversely, GPT4-HOLMES, while being aware of the number of cars, lacks information about their spatial relation, leading it to recommend slowing down because of the existence of 14 cars. This suggests that the End2End method is superior in perceiving certain visual features that are not captured by the APIs. Conversely, some specialized APIs, such as traffic sign detection, outperform GPT4-V in tasks like traffic sign detection, as they are specifically trained for this task. This could enable the HOLMES method to gather more accurate information than the End2End model.

Evaluator Model	Kappa Coefficients		
	P	C	A
GPT4 [†]	0.71	0.82	0.94
Qwen1.5-72B-Chat [†]	0.30	0.49	0.60
Qwen1.5-14B-Chat [†]	0.16	0.24	0.16
Qwen1.5-7B-Chat [†]	0.20	0.11	0.06

Table 3: Comparison of Open[†] and Close[†] LLMs as Evaluators. Kappa coefficients of Qwens increase when the model scales up.

4 Discussion

4.1 Strong LLMs are Good Error Locators.

As shown in Table 3, we compare the scoring kappa coefficients with human assessments for different LLMs. We randomly select 300 model outputs equally from different domains and ask 3 human experts to give perception, cognition, and action scores. The final result is based on the majority of three annotators. The result underscores a significant agreement between GPT-4 scoring and human evaluation results. This is substantiated by Cohen-Kappa coefficients of 0.71, 0.82, and 0.94 for perception, cognition, and action evaluations.

We also compare open models as evaluators. We choose one of the best open LLMs, Qwen1.5¹ series from 7B, 14B to 72B version. *Currently open LLMs tend to give wrongly high judgments in all sub-scores.* Although currently trailing behind GPT-4 in performance, we anticipate that with targeted training focused on error identification and enhancements in the overall capabilities of open LLMs, these models will become more effective evaluation tools compared to closed models. This is primarily due to the reproducible and transparent nature of open models, which offer significant advantages in the development of evaluation tools.

¹<https://huggingface.co/collections/Qwen>

4.2 Genuine PCA Score

PCA-Eval could pinpoint cases where the MLLM gets the correct answer by a fluke where perception or cognition score is 0 but the action score is 1. It explains why for some models, the action score is higher than perception and cognition scores. For instance, a model might opt for a conservative action, such as slowing down, even without accurately recognizing snowy weather in the image, resulting in a fluky correct action. In another scenario, if the model exhibits a preference for a specific choice index, it will attain a high action score provided that the evaluation dataset contains a substantial number of correct choices matching the preferred index, a phenomenon attributable to the positional biases inherent in both the model and the dataset. To overcome the mentioned bias when evaluating the genuine ability of MLLM, we propose a new metric **Genuine PCA Score**. It is equal to one if the perception, cognition and action scores are all 1 for one model’s response to a question. We find that for all models, there exists significant gap ($>10\%$) between the action score and genuine PCA score in average, revealing that relying on single metric such as choice accuracy is very problematic when conducting model evaluation. In our [online leaderboard](#), both average action score and average genuine PCA score are considered when ranking the candidate models.

4.3 Alignment between Agent Decisions and Human Values

We have observed instances where the decisions made by the agent contradict human values. Consider the scenario depicted in Figure 9 from Appendix. The image illustrates a crosswalk without pedestrians. The appropriate response would be slowing down, as caution is paramount when approaching a crosswalk, regardless of the presence or absence of pedestrians. However, upon processing the information that the crosswalk is empty, ChatGPT suggests that maintaining the current speed is the optimal action, arguing that the absence of pedestrians eliminates the need to slow down. The rationale provided by ChatGPT is logical, yet it does not align with human values.

5 Related Work

MLLM Benchmark. In recent times, there have been several benchmarks built for evaluating MLLMs, such as MMBench, MME, Seed-Bench,

POPE (Liu et al., 2023e; Fu et al., 2023; Li et al., 2023a,e) that assess MLLMs performance from multiple fine-grained dimensions. Visit-Bench, LVLM-eHub, M3IT (Bitton et al., 2023; Xu et al., 2023; Li et al., 2023c) focus on the general instruction following ability. General VQA tasks like OKVQA, VQAv2, Vizwiz, ScienceQA, VSR and IconQA (Marino et al., 2019; Agrawal et al., 2015; Gurari et al., 2018; Lu et al., 2022; Liu et al., 2023a; Lu et al., 2021) focus on visual understanding. MMMU, MathVista, LLaVA-benchmark and MM-Vet (Yue et al., 2023; Lu et al., 2023a; Liu et al., 2023c; Yu et al., 2023) require abilities from the vision part and specific knowledge in the language part. A lack of error localization techniques beyond accuracy assessments is among current benchmarks. This complicates identifying which part of the MLLM malfunctioned when making mistakes. Unlike prior work, PCA-Bench is more relevant to evaluate MLLMs’ ability to utilize integrated abilities to solve one task and make explainable decisions via error localization.

LLM Agent and Embodied Decision Making.

Using LLMs to empower the AI agents (Xi et al., 2023; Liu et al., 2023d; Park et al., 2023; Wang et al., 2023d) becomes more and more promising. Specifically, we can employ LLMs to enhance the decision making ability of the agents (Nakano et al., 2022; Yao et al., 2022; Li et al., 2023d; Song et al., 2023; Li et al., 2023b), expanding their perception and action space through strategies like tool utilization (Schick et al., 2023; Qin et al., 2023; Lu et al., 2023b). This line of research divides the entire decision-making process into two phases: (1) information seeking, usually involving MLLMs to verbalize the current status of AI agents in the vision-based environment with natural language; (2) reasoning and planning with text-based LLMs to decide what the AI agent should do in the next step with textual clues. Although LLM-based agents demonstrate reasoning and planning abilities through techniques like Chain of Thought or problem decomposition (Wei et al., 2023; Yao et al., 2023; Kojima et al., 2022), they inherently lack visual perception, and are limited to the discrete textual content. Therefore, integrating multimodal information can offer agents a broader context and a more precise understanding, such as PaLM-E (Driess et al., 2023), enhancing their environmental perception. However, there is still large gap deploying MLLM in various embodied

environments due to the lack of appropriate benchmark and interface linking those two domains while PCA-Bench is an attempt towards that goal.

6 Conclusion

In this paper, we introduce PCA-Bench, a multimodal benchmark designed to assess the integrated decision-making capabilities of MLLMs. This benchmark features PCA-EVAL, a novel fine-grained automatic evaluation tool that diagnoses decision making processes from three critical perspectives: perception, cognition, and action. To enhance the decision making ability from data perspective, we propose the Embodied Instruction Evolution method to automatically synthesize instruction examples from different environments, which has been proven effective in our main experiments. We believe that powerful MLLMs pave a new and promising way toward decision making in embodied environments and we hope PCA-Bench could serve as a good benchmark in evaluation and error localization for MLLMs' development.

7 Limitations

The current scope of PCA-Bench is confined to merely three domains in static environments. One of our future works aims to broaden this scope to encompass more domains and dynamic embodied environments where MLLMs could keep getting feedback, which is closer to real embodied AI scenarios. We do not apply different inference enhancement methods like In-Context Learning and Reflection in the decision making process of MLLMs. We just use the simplest prompting method and leave the exploration of a better cross-modal Chain-of-Thought method for future studies. Currently, PCA-Eval shows the best consistency with human evaluators when using powerful close LLM GPT4, which would bring additional cost to the user of PCA-Eval. We plan to develop and release an open error locator for error localization in the benchmark in the future.

References

Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick, Devi Parikh, and Dhruv Batra. 2015. Vqa: Visual question answering. *International Journal of Computer Vision*, 123:4 – 31.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,

and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *ArXiv*, abs/2308.12966.

Yonatan Bitton, Hritik Bansal, Jack Hessel, Rulin Shao, Wanrong Zhu, Anas Awadalla, Josh Gardner, Rohan Taori, and Ludwig Schmidt. 2023. Visit-bench: A benchmark for vision-language instruction following inspired by real-world use.

Liang Chen, Yichi Zhang, Shuhuai Ren, Haozhe Zhao, Zefan Cai, Yuchi Wang, Peiyi Wang, Tianyu Liu, and Baobao Chang. 2023. Towards end-to-end embodied decision making via multi-modal large language model: Explorations with gpt4-vision and beyond. *ArXiv*.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models.

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023a. Instructblip: Towards general-purpose vision-language models with instruction tuning.

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Albert Li, Pascale Fung, and Steven C. H. Hoi. 2023b. Instructblip: Towards general-purpose vision-language models with instruction tuning. *ArXiv*, abs/2305.06500.

Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. 2023. Palm-e: An embodied multimodal language model.

Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. 2022. Minedojo: Building open-ended embodied agents with internet-scale knowledge. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, et al. 2023. Mme: A compreh-

- hensive evaluation benchmark for multimodal large language models. [arXiv preprint arXiv:2306.13394](#).
- Joaquin M. Fuster. 2004. [Upper processing stages of the perception–action cycle](#). *Trends in Cognitive Sciences*, 8(4):143–145.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In [2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017](#), pages 6325–6334.
- Danna Gurari, Qing Li, Abigale Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. In [2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition](#), pages 3608–3617.
- Seungone Kim, Se June Joo, Doyoung Kim, Joel Jang, Seonghyeon Ye, Jamin Shin, and Minjoon Seo. 2023. [The cot collection: Improving zero-shot and few-shot learning of language models via chain-of-thought fine-tuning](#).
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.
- Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli Van derBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. 2017. AI2-THOR: An Interactive 3D Environment for Visual AI. [arXiv](#).
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023a. Seed-bench: Benchmarking multimodal llms with generative comprehension. [ArXiv](#), abs/2307.16125.
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbulin, and Bernard Ghanem. 2023b. [Camel: Communicative agents for "mind" exploration of large language model society](#).
- Lei Li, Yuwei Yin, Shicheng Li, Liang Chen, Peiyi Wang, Shuhuai Ren, Mukai Li, Yazheng Yang, Jingjing Xu, Xu Sun, Lingpeng Kong, and Qi Liu. 2023c. M³it: A large-scale dataset towards multi-modal multilingual instruction tuning. [arXiv preprint arXiv:2306.04387](#).
- Minghao Li, Feifan Song, Bowen Yu, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023d. [Api-bank: A benchmark for tool-augmented llms](#).
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023e. [Evaluating object hallucination in large vision-language models](#).
- Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, Jiaming Han, Siyuan Huang, Yichi Zhang, Xuming He, Hongsheng Li, and Yu Qiao. 2023. [Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models](#).
- Fangyu Liu, Guy Edward Toh Emerson, and Nigel Collier. 2023a. Visual spatial reasoning. [Transactions of the Association for Computational Linguistics](#).
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. Visual instruction tuning. [arXiv preprint arXiv:2304.08485](#).
- Ruibo Liu, Ruixin Yang, Chenyan Jia, Ge Zhang, Denny Zhou, Andrew M Dai, Diyi Yang, and Soroush Vosoughi. 2023d. Training socially aligned language models in simulated human society. [arXiv preprint arXiv:2305.16960](#).
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023e. Mmbench: Is your multi-modal model an all-around player? [arXiv preprint arXiv:2307.06281](#).
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chun yue Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023a. Mathvista: Evaluating math reasoning in visual contexts with gpt-4v, bard, and other large multimodal models. [ArXiv](#), abs/2310.02255.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multi-modal reasoning via thought chains for science question answering. In [The 36th Conference on Neural Information Processing Systems \(NeurIPS\)](#).
- Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2023b. Chameleon: Plug-and-play compositional reasoning with large language models. [arXiv preprint arXiv:2304.09842](#).
- Pan Lu, Liang Qiu, Jiaqi Chen, Tony Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-Chun Zhu. 2021. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. In [The 35th Conference on Neural Information Processing Systems \(NeurIPS\) Track on Datasets and Benchmarks](#).
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. OK-VQA: A visual question answering benchmark requiring external knowledge. In [IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019](#), pages 3195–3204.

- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. 2022. [Webgpt: Browser-assisted question-answering with human feedback](#).
- OpenAI. 2023. Gpt-4v(ision) system card.
- Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. [arXiv preprint arXiv:2304.03442](#).
- Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, et al. 2023. Tool learning with foundation models. [arXiv preprint arXiv:2304.08354](#).
- Joseph Redmon and Ali Farhadi. 2018. Yolov3: An incremental improvement. [arXiv](#).
- Shuhuai Ren, Aston Zhang, Yi Zhu, Shuai Zhang, Shuai Zheng, Mu Li, Alex Smola, and Xu Sun. 2023. Prompt pre-training with twenty-thousand classes for open-vocabulary visual recognition. [arXiv preprint arXiv:2304.04704](#).
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. [Toolformer: Language models can teach themselves to use tools](#).
- Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. [ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks](#). In [The IEEE Conference on Computer Vision and Pattern Recognition \(CVPR\)](#).
- Yifan Song, Weimin Xiong, Dawei Zhu, Wenhao Wu, Han Qian, Mingbo Song, Hailiang Huang, Cheng Li, Ke Wang, Rong Yao, Ye Tian, and Sujian Li. 2023. [Restgpt: Connecting large language models with real-world restful apis](#).
- Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. 2017. A corpus of natural language for visual reasoning. In [Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics \(Volume 2: Short Papers\)](#), pages 217–223.
- Peiyi Wang, Lei Li, Liang Chen, Feifan Song, Binghuai Lin, Yunbo Cao, Tianyu Liu, and Zhifang Sui. 2023a. Making large language models better reasoners with alignment. [arXiv preprint arXiv:2309.02144](#).
- Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023b. Large language models are not fair evaluators. [arXiv preprint arXiv:2305.17926](#).
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023c. [Self-instruct: Aligning language models with self-generated instructions](#).
- Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. 2023d. [Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents](#). [arXiv, abs/2302.01560](#).
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-thought prompting elicits reasoning in large language models](#).
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. 2023. [The rise and potential of large language model based agents: A survey](#).
- Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. 2023. [Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models](#). [arXiv preprint arXiv:2306.09265](#).
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. In [The Eleventh International Conference on Learning Representations](#).
- Weirao Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh Murthy, Zeyuan Chen, Jianguo Zhang, Devansh Arpit, et al. 2023. Retroformer: Retrospective large language agents with policy gradient optimization. [arXiv preprint arXiv:2308.02151](#).
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. [arXiv, abs/2308.02490](#).
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 2023. [Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi](#).
- Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023. Multi-modal chain-of-thought reasoning in language models. [arXiv preprint arXiv:2302.00923](#).

Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng Wang, Wenjuan Han, and Baobao Chang. 2023. Mmicl: Empowering vision-language model with multi-modal in-context learning. [arXiv preprint arXiv:2309.07915](#).

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#).

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. [arXiv preprint arXiv:2304.10592](#).

Zhe Zhu, Dun Liang, Songhai Zhang, Xiaolei Huang, Baoli Li, and Shimin Hu. 2016. Traffic-sign detection and classification in the wild. In [The IEEE Conference on Computer Vision and Pattern Recognition \(CVPR\)](#).

A Examples of PCA-Bench

A.1 Data Distribution

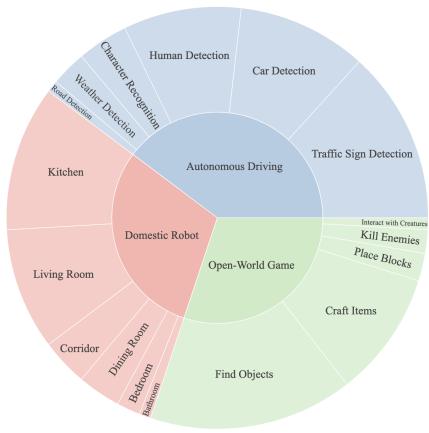


Figure 6: Domain and required ability distribution of PCA-Bench.

The PCA-Bench’s data distribution across various domains is outlined in Figure 6. For the Autonomous Driving domain, instances are grouped by their respective task types. In the Domestic Robot domain, instances are grouped by their locations. In the Open-World Game domain, instances are grouped by the tasks they aim to accomplish.

B Human Annotation Pipelines

The annotation process consists of two stages: (1) Dataset Annotation, and (2) Dataset Refinement. During the initial stage, three annotators are assigned to each domain, adhering strictly to the respective annotation guidelines. They first pinpoint the source images from each domain that are informative and meaningful so that they can write questions for each image. The annotators have the responsibility to ensure every question has only one correct answer and accurate rationales. In the subsequent stage, annotators are instructed to scrutinize the output actions and rationales presented by ChatGPT and check the annotations. This process aims to address the challenge of multiple correct answers, as ChatGPT can furnish comprehensive explanations for its actions. These explanations assist annotators in assessing the acceptability of ChatGPT’s response, particularly when it deviates from the established ground truth answer. This enables annotators to refine annotations to ensure the presence of a single correct answer.

B.1 PCA-Bench Examples

We list three examples of each domain from PCA-Bench, as shown in Figure 10, 11, and 12.

C Zero Shot Modality Conversion: HOLMES

To optimize the evaluation process of HOLMES² method, we pre-execute all relevant APIs for each instance within a selected subset of 300 instances from the PCA-Bench test set, recording the results for individual instances. This method enables immediate access to specific API results, eliminating the need to rerun the model for each evaluation instance.

Traffic Domain. Below is the API description for the traffic domain.

```

1 # API Description for Traffic Domain:
2 def detect_traffic_sign():
3     """
4         Detects traffic signs in the image.
5         :return: list of detected traffic
6             signs and coordinates, e.g. ['stop',
7                 'max speed limit']
8         """
9     pass
10
11 def object_detection():
12     """
13         Detects objects in the image.
14         :return: dict of detected objects
15             and number of the objects, e.g. {'car':10, 'person':1}
16     """
17     pass
18
19 def ocr():
20     """
21         Performs OCR on the image.
22         :return: list of detected text, e.g.
23             ['Changjiang road', 'Right lane
24                 closure']
25     """
26     pass
27
28 def image_caption():
29     """
30         Generates a caption for the image.
31         :return: caption, e.g. 'A red car
32             driving down the street'
33     """
34     pass
35
36 def weather_detection():
37     """
38         Detect current weather.
39         :return: weather, e.g. 'rainy' or
40             'clear'
41     """
42     pass

```

- *detect_traffic_sign()*: The detection of road traffic signs model utilize YOLO (Redmon and Farhadi, 2018) which trained on the Tsinghua-Tencent 100K dataset (Zhu et al., 2016). TT100K

²Originally proposed in an early version of this paper (Chen et al., 2023)

Autonomous Driving

What is the best action to take when you are driving on the highway?
Actions: A. Slow down B. Keep driving C. Stop the car D. Change to other lane

I will first call the `detect_traffic_sign()` API to check if there are any traffic signs in the image that could inform the best action.
`detect_traffic_sign() = None`

Since there are no traffic signs detected, I will now call the `object_detection()` API to check for any objects in the image that could inform the best action.
`object_detection() = None`

Since there are no objects detected, I will now call the `caption()` API to get a general description of the image, which might help in determining the best action.
`caption() = The image features a long, empty highway.`

Based on the provided caption, it seems that the highway is not congested. Therefore, the best action to take when driving on the highway in this situation would be: **B. Keep driving**.

Domestic Robot

You are a domestic robot that helps me do housework. Your goal is to heat a cup of milk. What is the best action to take next?
Actions: A. Find a microwave. B. Put a cup of milk into the microwave. C. Find milk. D. Pick up a cup of milk.

I will first call the `object_detection_alfred()` API to check for any relevant objects in the current view.
`object_detection_alfred() = [StoveTopDoor!, 'Lettuce', 'StoveBurner', 'SaltShaker', 'CounterTop', 'Spatula', 'Pot', 'StoveKnob', 'StoveBase!', 'Cabinet', 'DishSponge', 'Microwave', 'Pan', 'PepperShaker', 'Ladle', 'Floor', 'SoapBottle', 'Dishwasher', 'Sink', 'Ceiling', 'Bread', 'Mug']`

Now that I have detected a *microwave* and a mug in the current view, I will call `list_inventory_alfred()` to check if I already have milk in my inventory.
`list_items_in_inventory_alfred() = ["cup of milk"]`

Since I already have a cup of milk in my inventory, the best action to take next is:
B. Put a cup of milk into the microwave.

Open-World Game

You are a powerful game assistant in Minecraft. Your goal is to Harvest wool. Based on current condition, what is the best action to do next?
Actions: A. Find sheep B. Shear sheep C. Craft shears D. Find iron ore E. Craft iron ingot

In order to answer the question, I would first need to check if I already have shears in the inventory. So I would call `list_inventory_information()`.

`list_inventory_information() = [{"shears": 1}]`

Since the inventory contains shears, then the next action would be dependent on the presence of sheep. I would then call `list_nearby_mobs()` to check if there are sheep nearby.
`list_nearby_mobs() = None`

Given that the inventory does have shears but no nearby sheep are detected, the best action to do next would be: **A. Find Sheep**

Image
 User
 LLM
 Model/API Hub

Figure 7: Three examples of HOLMES solving questions from different domains of PCA-Bench.

comprises 100,000 images encompassing 30,000 instances of traffic signs. The end-to-end YOLO enables simultaneous detection and classification of traffic signs.

- *object_detection()*: Objects demanding attention during vehicle operation primarily encompass cars, pedestrians, and bicycles. A surfeit of vehicles can lead to traffic congestion, while the presence of pedestrians or bicycles ahead necessitates cars to decelerate and proceed cautiously. Hence, the *object_detection()* API predominantly identifies three key object categories: cars, pedestrians, and bicycles. We utilize PMOP (Ren et al., 2023), a model trained on vision-language models through the prompt pre-training method, which enables the detection and counting of the three mentioned objectives by modifying specific class names.

- *ocr()*: We employ PaddleOCR³ to extract textual information from images, providing crucial road data for real-time navigation.

- *image_caption()*: To initially streamline the road information within the image, we employ the BLIP2-flan-t5-xl to generate an initial caption for the picture. This caption, derived from basic image data, is then utilized as input for the model to facilitate decision-making.

- *weather_detection()*: Weather detection leverages a pre-trained ResNet50 model⁴, derived from a dataset of more than 70,000 weather records. This model extracts weather information from provided images to inform decision-making.

Domestic Robot Domain. Below is the API description for the Domestic Robot domain.

```

1 #API Description for Domestic Robot
2     Domain
3 def object\_detection():
4     """
5         Detects objects in current view,
6         which you don't need to find.
7         :return: list of detected objects, e
8             g. ['chair', 'table']
9     """
10    pass
11
12 def list_items_in_hands():
13     """
14         Lists items in your hand, which you
15         don't need to pick up
16         :return: list of items in hand, e.g.
17             ['coffee cup', 'milk']
18     """
19    pass

```

Game Domain. Below is the API description for the Game domain (Minedojo).

```

1 #API Description for Game Domain
2 def list_nearby_mobs_in_minecraft():

```

³<https://github.com/PaddlePaddle/PaddleOCR/tree/release/2.7>

⁴<https://github.com/mengxianglong123/weather-recognition>

```

3     """
4     Lists nearby mobs in Minecraft.
5     :return: list of nearby mobs, e.g.
6     ['creeper', 'pig']
7     """
8
9 def list_inventory_information():
10    """
11    Lists inventory information of the
12    player in Minecraft.
13    :return: list of inventory
14    information with number, e.g. [('
15    diamond', 64), ('iron', 32)]
16    """
17
18 pass

```

Note that within the Domestic Robot Domain and Game Domain, APIs can be directly accessed within the virtual environment, allowing for the perception of the surrounding objects and the current picture context.



Figure 8: A Comparison between GPT4-V and GPT4-HOLMES

Autonomous Driving

What is the best action to take when you are driving?
Actions: A. Slow down B. Keep driving C. Stop the car D. Change to other lane

object_detection() = {'car':0, 'human':0, 'bicycle':0}

caption() = An empty crosswalk.

ChatGPT Based on the provided caption and the result of object detection, the image depicts an empty crosswalk with no humans or cars. Therefore, the best action is **B: Keep driving**, as it's not necessary to slow down.

Figure 9: A case showing the value mis-alignment between agent and human's decision.

D Automatic Evaluation

We utilize the template as shown in Table 4 to query GPT-4, aiming to evaluate its responses and assign scores for perception, cognition, and action. By feeding both the agent's output and the ground truth answer to GPT-4, based on this template, we can then extract the three distinct scores from the conclusion of GPT-4's response.

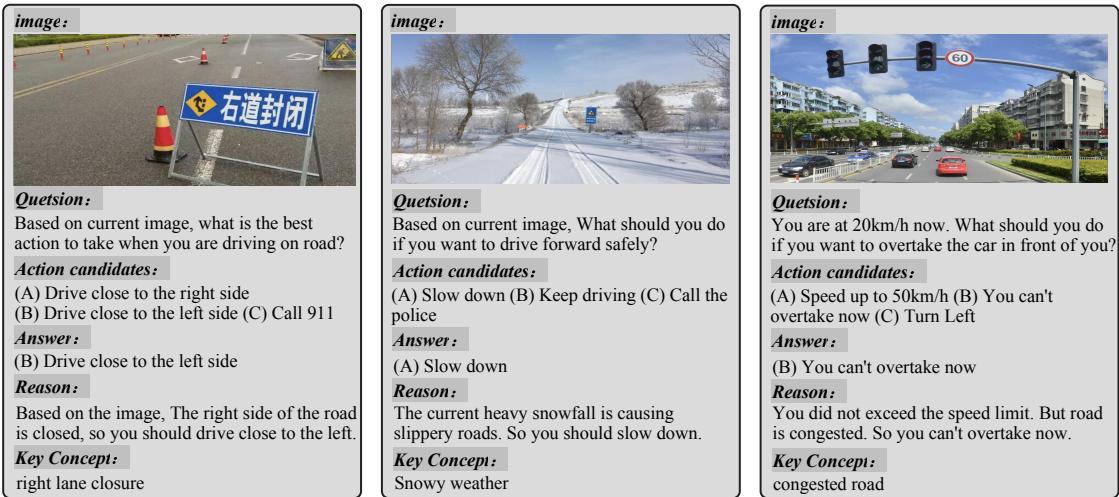


Figure 10: Three examples of PCA-Bench in the autonomous driving domain.

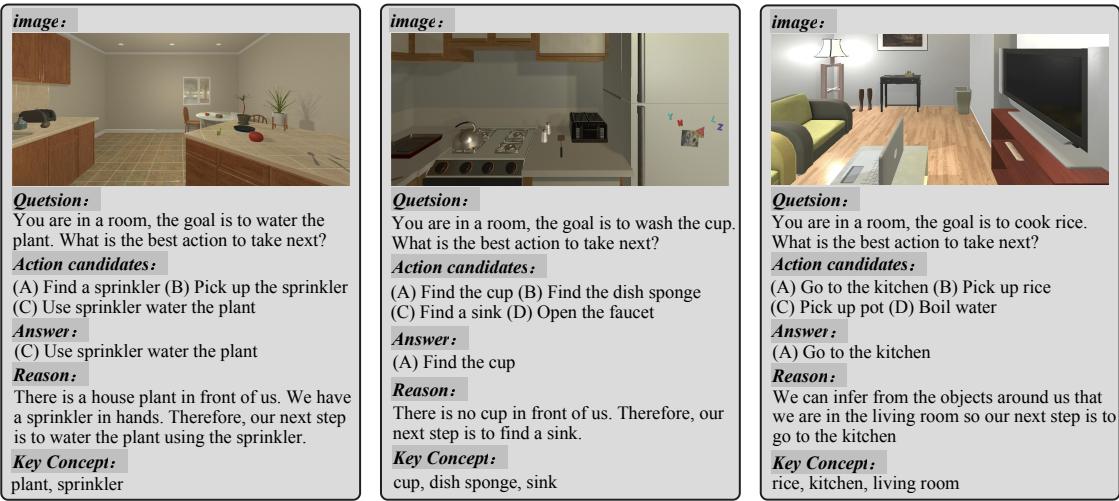


Figure 11: Three examples of PCA-Bench in the domestic robot domain.

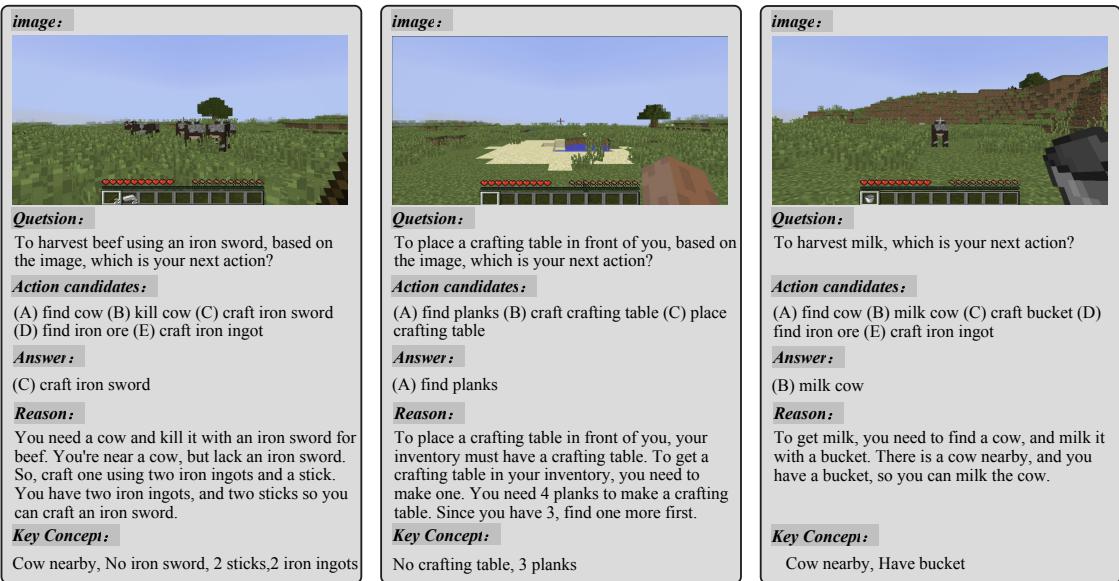


Figure 12: Three examples of PCA-Bench in the open-world game domain.

[Question]: {question}
[Action Choices]: {actions}
[Agent Answer]: {model_output}
[Correct Action]: {true_action}
[Key Concepts]: {key_concept}
[Reference Reasoning Process]: {reason}
[System]

We would like you to access the agent’s performance in the multimodal reasoning task about domain. In this task, the agent is given an image, a [Question], and several candidate [Action Choices], and is asked to give an [Agent Answer] for the [Question]. The [Agent Answer] encapsulates the agent’s perception of the image’s [Key Concepts], the agent’s cognition reasoning process and the final selected action.

We request you to give three types of scores for the agent’s [Agent Answer] in comparison to the given [Key Concepts], [Reference Reasoning Process] and [Correct Action]:

1. action score: If the selected action in the [Agent Answer] matches that of the [Correct Action], the action score is 1; otherwise, it is 0.
2. perception score: This score evaluates the model’s capability to perceive and interpret observations. It is contingent on whether the [Agent Answer] includes any of the [Key Concepts] of the instance. If it accurately describes any one of the [Key Concepts], the score is 1; otherwise, it is 0.
3. cognition score: This score gauges the model’s ability to reason, comprehend, and make informed decisions based on perceived input data and world knowledge. If the reasoning process in the [Agent Answer] aligns with the [Reference Reasoning Process], the score is 1; otherwise, it is 0.

Please note that there are only scores of 0 and 1.

You should carefully compare the [Agent Answer] with the [Correct Action], [Key Concepts] and [Reference Reasoning Process] to give your assessment.

You need first to give your assessment evidence and then the scores.

Your output MUST contain 6 lines with the following format:

action assessment evidence: (assessment evidence here)

action score: (score here)

perception assessment evidence: (assessment evidence here)

perception score: (score here)

cognition assessment evidence: (assessment evidence here)

cognition score: (score here)

Table 4: The template of querying GPT-4.

E Training Details

Table 5 shows the specific parameters used for fine-tuning in different models. The PCA results on the three domains of PCA bench before and after fine-tuning different models are shown in Figure 13.

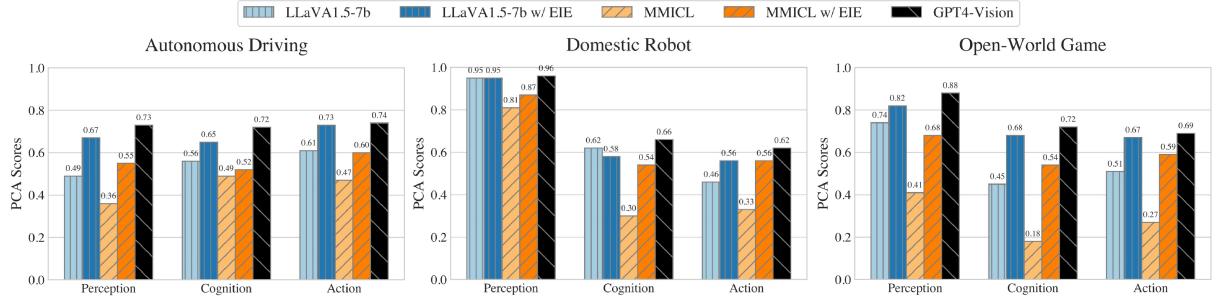


Figure 13: Performance comparsion between models’ zero-shot results and models’ finetuned results with the data generated by Embodied-Instruct-Evolution (EIE) method. EIE improves the performance on all domains for both LLaVA1.5-13b and MMICL models.

Model	Parameter	Value
Qwen-VL-Chat/LLaVA1.5-7/13b	Learning Rate	2e-4
	Use Lora Finetuning?	Yes
	Lora Rank	8
	Lora Alpha	32
	Global Batchsize	20
	Weight Decay	0
	Train Epochs	5
MMICL	Lr Scheduler Type	Cosine
	Warmup Ratio	0.03
MMICL	Learning Rate	5e-4
	Use Lora Finetuning?	No
	Global Batchsize	20
	Weight Decay	5e-4
	Train Epochs	5
	Lr Scheduler Type	Linear
	Warmup Ratio	0.2

Table 5: Training details for different models with EIE.

F Does Chain-of-Thought Finetuning Improve Cross-modal Reasoning?

Unlike vanilla finetuning, which solely focuses on delivering direct answers, Chain-of-Thought Finetuning necessitates the model to first articulate its reasoning before presenting the answer. This approach has been demonstrated to be a highly effective instruction tuning paradigm for LLMs (Chung et al., 2022; Kim et al., 2023). We have incorporated this methodology in our previous finetuning experiments.

To further evaluate its impact, we conducted an ablation study where the reasoning process was omitted from the target output during the training of MLLMs. We then assessed the variations in action scores on the test set. As depicted in Figure 14, to our surprise, the figures suggest that Chain-of-Thought finetuning exerts a relatively minor influence when compared to conventional label finetuning. We have noticed that similar phenomena has been identified by Zhang et al. (2023) that standard CoT finetuning does not work for MLLMs in their explorations.

We think there are three potential explanations: 1) Task Variation: Contrary to mathematics datasets like GSM8K, the current task doesn't require multi-step complex reasoning to arrive at the final answer and the automatic generated CoTs have noise. 2) Modality Discrepancy: The CoT capability, inherent in LLMs, is only moderately adjusted for visual input for current open-source MLLMs. This adaptation process could potentially impair the reasoning ability. 3) Short Cut in Pretraining: We think a deeper reason might lie in the short-cut during pretraining period of current open-source MLLMs, which are pretrained on simple image caption task in a large scale. Those captions are usually short and lose a lot of information about the original image. What's more important is that the reasoning ability of LLM is not utilized during the pretraining stage, which might hurt the reasoning ability of LLM during the SFT period. We defer to future research how to effectively harness the CoT capabilities of LLMs to enhance embodied decision-making processes.

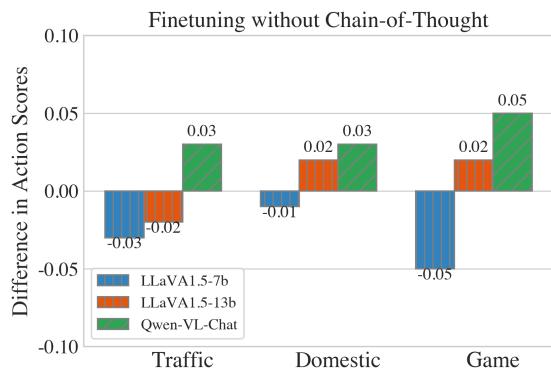


Figure 14: Action scores changes when training without reasoning process for different models. The benefit of CoT finetuning is not consistent among models. Blue means difference of action score for LLaVA1.5-7b, Orange means difference of action score for LLaVA1.5-13b and green means difference of action score for Qwen-VL-Chat.