

# TOWARDS END-TO-END EMBODIED DECISION MAKING VIA MULTI-MODAL LARGE LANGUAGE MODEL: EXPLORATIONS WITH GPT4-VISION AND BEYOND

Liang Chen<sup>1</sup>, Yichi Zhang<sup>1</sup>, Shuhuai Ren<sup>1</sup>, Haozhe Zhao<sup>1</sup>, Zefan Cai<sup>1</sup>, Yuchi Wang<sup>1</sup>  
Peiyi Wang<sup>1</sup>, Tianyu Liu<sup>2</sup>, Baobao Chang<sup>1\*</sup>

<sup>1</sup> Peking University, <sup>2</sup> Tencent Cloud AI

{leo.liang.chen, yczhang, shuhuai\_ren}@stu.pku.edu.cn

{hanszhao, zefan, wangyuchi, wangpeiyyi}@stu.pku.edu.cn

{tianyu0421, chbb}@pku.edu.cn

<https://github.com/pkunlp-icler/PCA-EVAL>

## ABSTRACT

In this study, we explore the potential of Multimodal Large Language Models (MLLMs) in improving embodied decision-making processes for agents. While Large Language Models (LLMs) have been widely used due to their advanced reasoning skills and vast world knowledge, MLLMs like GPT4-Vision offer enhanced visual understanding and reasoning capabilities. We investigate whether state-of-the-art MLLMs can handle embodied decision-making in an end-to-end manner and whether collaborations between LLMs and MLLMs can enhance decision-making. To address these questions, we introduce a new benchmark called **PCA-EVAL**, which evaluates embodied decision-making from the perspectives of **P**erception, **C**ognition, and **A**ction. Additionally, we propose **HOLMES**, a multi-agent cooperation framework that allows LLMs to leverage MLLMs and APIs to gather multimodal information for informed decision-making. We compare end-to-end embodied decision-making and HOLMES on our benchmark and find that the GPT4-Vision model demonstrates strong end-to-end embodied decision-making abilities, outperforming GPT4-HOLMES in terms of average decision accuracy (+3%). However, this performance is exclusive to the latest GPT4-Vision model, surpassing the open-source state-of-the-art MLLM by 26%. Our results indicate that powerful MLLMs like GPT4-Vision hold promise for decision-making in embodied agents, offering new avenues for MLLM research.

## 1 INTRODUCTION

The capacity to make well-informed decisions is essential for the survival and success of living organisms in their respective environments. Similarly, a major goal in embodied artificial intelligence is to develop agents, like robots, with sophisticated decision-making abilities. This could enable artificial agents to intelligently interact with their surroundings and efficiently accomplish a variety of real-world tasks such as autonomous driving (Hu et al., 2023; Wayve, 2023), domestic assistance (Kolve et al., 2017; Shridhar et al., 2020; Huang et al., 2022b), and game playing (Fan et al., 2022; Wang et al., 2023a; Zhu et al., 2023b). Recently, there has been a notable increase in leveraging exceptional reasoning capabilities and world knowledge of Large Language Models (LLMs) to enhance decision making in agents. However, LLMs are primarily designed to process textual context, creating a modality gap (Liang et al., 2022; Ren et al., 2023a) for the LLM-powered agent when dealing with multimodal observations in real-world scenarios.

To bridge this modality gap, a common approach is converting multimodal observations into text using various APIs (Wu et al., 2023; Yang et al., 2023). However, this conversion can result in information loss during the transition from multimodal to unimodal text. At the same time, recent advances in Multimodal Large Language Models (MLLMs), particularly Visual Large Language

---

\*Corresponding author.

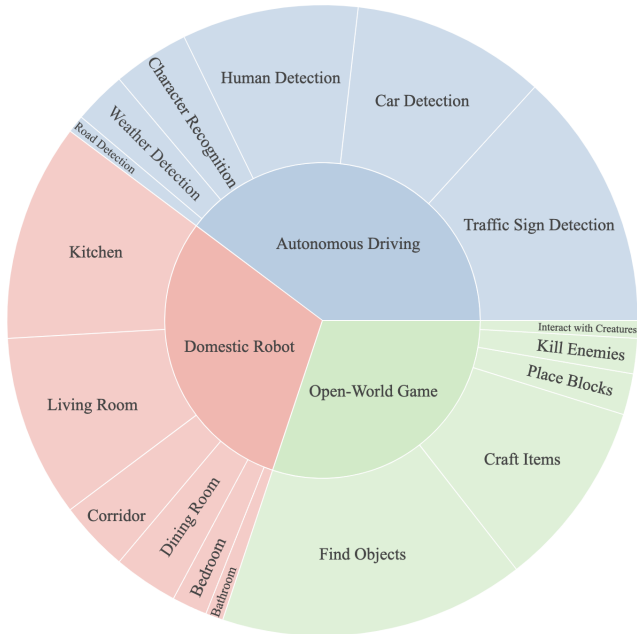


Figure 1: Domain and required ability distribution of PCA-EVAL.

Models (VLLMs) like GPT4-Vision (OpenAI, 2023a), have showcased impressive general-purpose visual understanding and reasoning abilities (Zhu et al., 2023a; Dai et al., 2023; Liu et al., 2023a; Li et al., 2023b; Zhao et al., 2023). These VLLMs can directly perceive the visual information rather than relying on textual intermediaries, potentially enabling more sophisticated reasoning and decision making for embodied agents operating in complex real-world environments. Considering these developments, two research questions naturally arise: (1) Can current state-of-the-art VLLMs perform various embodied decision making tasks in an end-to-end manner? What are the current strengths and limitations when compared to LLM-powered agents? (2) Can LLMs and VLLMs collaborate to enhance embodied decision-making capabilities?

However, addressing these questions is challenging due to the absence of an existing evaluation benchmark that satisfies the following criteria: (1) supporting end-to-end embodied decision making by providing agents with direct multimodal observations; (2) enabling multi-dimensional evaluation of the decision-making process, encompassing perception, reasoning, and action perspectives, rather than relying solely on final rewards or success rate; and (3) covering diverse domains, drawing from different areas of embodied AI. The development of more comprehensive benchmarks that meet these desiderata could substantially advance research on decision making in embodied systems.

In this paper, we propose a new benchmark, **PCA-EVAL**, for evaluating the embodied decision-making ability of agents from three perspectives, i.e., **P**erception, **C**ognition, and **A**ction. Our benchmark covers three domains as illustrated in Figure 1: autonomous driving, domestic assistance, and game-playing. The corresponding data are collected from real-world transportation scenes (Zhu et al., 2016), domestic housekeeper environment based on ALFRED (Shridhar et al., 2020), and Open-world environment Minedojo (Fan et al., 2022) based on the famous game Minecraft. This diverse set of domains allows for a comprehensive assessment of embodied decision-making capabilities across various contexts. Distinct from the MDP-based evaluation that solely focuses on maximizing cumulative rewards, we divide the sequential decision making process into multiple one-step decision problems based on a task-specific topology graph. Each instance in the benchmark consists of a 6-element tuple:  $\langle image, question, action\ candidates, answer, reason, key\ concept \rangle$ . Adopting this approach offers two major advantages: (1) It enables a more comprehensive evaluation of the decision-making process, with each decision step being assessed in terms of perception, cognition, and action. (2) The evaluation can be conducted outside complex simulation environments, simplifying the process of evaluating different agents and their performance.

With the proposed benchmark, we conduct two series of evaluation: (1) We examine multiple state-of-the-art VLLMs, like InstructBLIP (Dai et al., 2023), MMICL (Zhao et al., 2023), QwenVL-Chat (Bai et al., 2023) and the latest GPT4-Vision (OpenAI, 2023a), in an end-to-end decision making context. (2) We introduce **HOLMES**,<sup>1</sup> a multi-agent cooperation framework. In this framework, we provide large language models, such as ChatGPT (OpenAI, 2022), GPT4 (OpenAI, 2023b), and Vicuna (Chiang et al., 2023), with descriptions of vision models like image captioning, object detection, Optical Character Recognition (OCR), and traffic sign detection models. Additionally, we supply descriptions of valid APIs within the simulated environment. The large language model subsequently initiates a search for clues pertaining to the question by engaging in a multi-turn conversation. This process involves alternating between invoking models or APIs to find clues and analyzing the discovered clues to facilitate informed decision making.

From our experimental results, we discerned that within the end-to-end framework, GPT4-Vision significantly outshines the contemporary state-of-the-art vision-language model, MMICL, boasting an average action accuracy improvement of 26%. Notably, GPT4-Vision can furnish a detailed rationale behind its embodied decision-making process, a feature absent in present open-source VLLMs. When assessing HOLMES models, GPT4 consistently emerges superior across all three domains. Drawing a comparison between GPT4-Vision and HOLMES, we observed that GPT4-Vision surpasses GPT4-HOLMES with multiple expert visual APIs in terms of cognition and action scores. This underscores its broad adaptability across a spectrum of visual tasks and its good fusion of visual understanding, world knowledge, and embodied decision making.

In summary, we introduce three key contributions in this study:

1. We propose PCA-EVAL, a novel evaluation benchmark for multi-domain embodied decision making that evaluates performance in perception, cognition, and action.
2. We present HOLMES, a multi-agent cooperation framework designed to tackle various embodied decision-making tasks that include multimodal observations. It mimics the process of playing a detective game in which the LLM uncovers clues by utilizing various multimodal models or APIs supplied by the environment.
3. We conducted a systematic comparison of two embodied decision-making methods: end2end and HOLMES, across various models. Our findings suggest that when utilizing MLLM with the end2end method, it not only achieves decision accuracy better than the top-performing model (GPT-4) in HOLMES but also secures a superior cognition score. However, this level of performance is exclusive to the latest GPT4-Vision model, which significantly outpaces the open-source state-of-the-art VLLMs.

We believe that powerful MLLMs like GPT4-Vision pave a new and promising way toward decision making in embodied agents using LLMs. It enables decisions across diverse domains to be made and justified seamlessly in an end-to-end manner. PCA-EVAL serves as an effective metric for evaluating the embodied decision-making capabilities of both end-to-end and HOLMES-based models.

## 2 RELATED WORK

**Embodied Decision Making.** Research on embodied decision-making is an emerging trend for artificial intelligent agents to interact with their surroundings and accomplish numerous tasks. This necessitates proficiency in vision perception, world knowledge, and commonsense reasoning, areas where a large language model can provide some level of expertise. We group prior work on embodied decision-making with LLM into two main trends. The first trend is to transform multimodal information, including object and scenery identification, the current states of AI agents, and the feedback from the environments, to texts. Text-based LLMs can then reason over the textual clues to determine the next action towards completing a designated task (Huang et al., 2022a; Li et al., 2022; Huang et al., 2022b; Chen et al., 2023). This line of research divides the entire decision-making process into two phases: (1) information seeking, usually involving VLLMs 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. The other line of research uses multimodal LLMs directly for end-to-end decision making,

<sup>1</sup>The system is aptly named after the renowned detective, Sherlock Holmes.

such as PALM-E (Driess et al., 2023b). The end-to-end decision making poses greater challenges to multimodal LLMs as it requires the combination of different functionalities including perception, cognition, and action, whereas decision making without explicit multiple steps mitigates the error propagation between information seeking and reasoning.

**LLM-Powered Agents.** Large language models pre-trained on large-scale multimodal (including text, image, video, etc.) corpus demonstrate impressive emergent abilities and immense popularity (Brown et al., 2020; Wei et al., 2022), and have seen tremendous success across various domains covering various natural language processing and computer vision tasks (Radford et al., 2019; Chowdhery et al., 2022; Touvron et al., 2023; Alayrac et al., 2022; Zhu et al., 2023a; Li et al., 2023a). Consequently, using LLMs to empower the AI agents (Xi et al., 2023; Liu et al., 2023b; Park et al., 2023; Wang et al., 2023d; Yuan et al., 2023) 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., 2023c; Song et al., 2023), expanding their perception and action space through strategies like tool utilization (Schick et al., 2023; Qin et al., 2023; Lu et al., 2023). 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 visual information or other modalities can offer agents a broader context and a more precise understanding (Driess et al., 2023a), enhancing their environmental perception. However, no evaluation protocol or benchmark is currently available to evaluate decision making within the multimodal context.

### 3 PCA-EVAL

In this section, we propose to evaluate the decision-making ability of embodied agents from three perspectives: perception, cognition, and action. Accordingly, we present a novel benchmark named PCA-EVAL. Our PCA-EVAL benchmark consists of 300 multimodal multiple-choice questions with diverse embodied topics and annotations of their answers with corresponding explanations.

As shown in Figure 5, each instance in the benchmark consists of 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, like transportation scenes, housekeeper environments, and game worlds in Minecraft. Questions, action candidates, and answers are derived from real tasks within the corresponding environment. The reasoning explains why the answer is the best choice for the current image, while the key concept highlights the most question-related aspect in the image.

Unlike traditional visual question-answering datasets that emphasize visual perception (e.g., VQA (Goyal et al., 2017)), visual reasoning (e.g., NLVR (Suh et al., 2017)), or world knowledge (e.g., OKVQA (Marino et al., 2019)), the most distinctive characteristic of PCA-EVAL is its grounding in embodied actions. Compared to embodied simulation environments like ALFRED (Shridhar et al., 2020) and Minedojo (Fan et al., 2022), PCA-EVAL proves to be more effective in evaluating various LLM-based agents. This is primarily due to PCA-EVAL’s provision of high-level actions that can be readily implemented or programmed using the low-level actions in the corresponding domains. The high-level actions are more comprehensible for LLMs than the direct low-level actions like robotic movements 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.

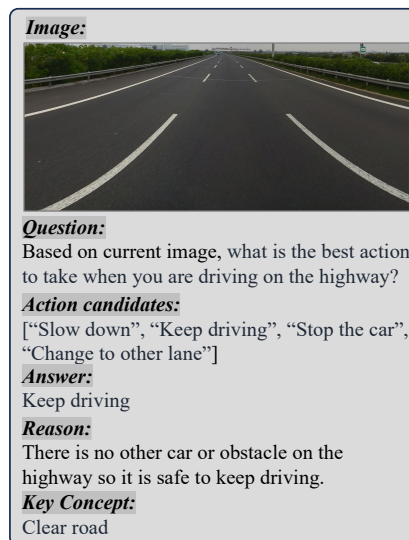


Figure 2: An instance of PCA-EVAL.

To answer a question in PCA-EVAL, 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 inevitably result in an incorrect answer, posing a significant challenge to the more complex capabilities of embodied agents. Although challenging, all the aforementioned abilities are essential for the decision-making process in embodied environments.

### 3.1 EVALUATION METRICS

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

**Perception Score.** The Perception Score (P-Score) measures the model’s ability to accurately perceive and interpret the observation. It is computed based on whether the agent’s output image description  $d$  includes the key concept of the instance. If the agent accurately describes the question-related key concept in the image, the P-score is assigned a value of 1; otherwise, it is assigned a value of 0. For the instance in Figure 5, 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.

**Cognition Score.** The 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 5, the agent should link the “clear road” to the action “keep driving” based on transportation commonsense to get the score.

**Action Score.** The 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.

The final Perception, Cognition, and Action scores of the agents are obtained by averaging the scores across all instances and domains in our PCA-EVAL dataset.

### 3.2 AUTOMATIC EVALUATION

Recent advancements have seen researchers harnessing powerful LLMs for the evaluation of 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. (2023c;b). 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 annotations and human annotator results. This is substantiated by Pearson correlation coefficients of 0.8, 0.9, and 0.95 for perception, cognition, and action evaluations, respectively. To facilitate ongoing and future research endeavors, we share our automatic evaluation script<sup>2</sup> for seamless adoption, which could also be improved in the future. For a detailed description of our evaluation methodology, kindly refer to Appendix C

### 3.3 DATASET OVERVIEW

The PCA-EVAL benchmark currently comprises three domains, with a total of 300 instances, including 100 instances per domain. In our preliminary study, we find that the annotation process requires proactive thinking of the questions, actions, and corresponding answers, which makes quality control difficult. In order to ensure the quality of PCA-Eval, every single test case has been verified by at least three authors of this paper. Although challenging, we would keep scaling this

<sup>2</sup>[https://github.com/pkunlp-icler/PCA-EVAL/blob/main/pca-eval/evaluation/pca\\_auto\\_scoring.py](https://github.com/pkunlp-icler/PCA-EVAL/blob/main/pca-eval/evaluation/pca_auto_scoring.py)

benchmark in order to advocate further attention to end-to-end decision-making. 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.

**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 in accordance with 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 and examples for each domain in Appendix A.

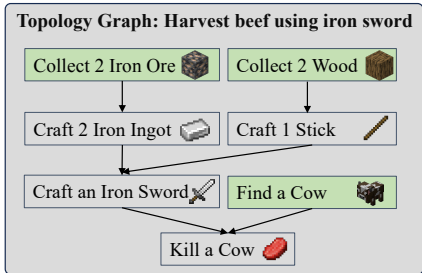


Figure 3: Illustration of task topology graph. Events in green represent the leaf nodes of the graph.

### 3.4 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.

## 4 METHODS

### 4.1 END2END DECISION MAKING VIA VLLMS

In this subsection, we detail the evaluation process for assessing state-of-the-art VLLMs, e.g., InstructBLIP, MMICL, and GPT4-Vision, on end-to-end embodied decision-making using the proposed PCA-EVAL benchmark. End2End embodied decision making is straightforward since we can directly feed the visual observation and the textual question to the multi-modal agent. As illustrated in Figure 5, the agent is prompted to output the image description and reasoning process before giving the final action.



Figure 4: Three examples of HOLMES solving questions from different domains of PCA-EVAL.

#### 4.2 HOLMES: MULTI-AGENT COOPERATION

Different from End2End embodied decision making, within HOLMES, we prompt large language models like ChatGPT-3.5 (OpenAI, 2022), GPT4 (OpenAI, 2023b) to call different visual models or APIs to gather information about the environment.

We provide these models with descriptions of the input and output for different visual models such as the image caption model based on InstructBLIP, the object detection model based on POMP (Ren et al., 2023b), and the traffic sign detection model based on YOLO (Redmon & Farhadi, 2018). Additionally, we supply descriptions of valid APIs within the simulated environment, such as `list_nearby_mobs_in_minecraft()` to tell what creatures can current player see and `list_items_at_hand_in_alfred()` to tell what item the robot is holding in hand. Full API description files for each domain are shown in Appendix B.

These integrations enable the large language model to initiate a search for clues pertaining to a given question through a multi-turn conversation. As shown in Figure 4, the process involves alternating between invoking models or APIs to gather relevant information and analyzing the discovered clues to facilitate informed decision making. The HOLMES framework is designed to enhance cooperation and coordination among multiple agents in dynamic and complex environments.

In HOLMES, there are four key components as depicted in Figure 4: the image, the user, the LLM, and the Model/API Hub. Initially, the user poses a question about the optimal action to take based on the environment shown in the image, providing potential action choices. As the LLM cannot directly view the image, it’s briefed with descriptions of available visual models and APIs supplied by the simulation environment. It’s then tasked with gathering relevant data via these models and APIs to determine the appropriate action. When the LLM responds, the system checks if it has invoked a legitimate model or API, subsequently relaying the results from the invoked API. This feedback is logged into the dialogue history, allowing the LLM to analyze and form subsequent responses.

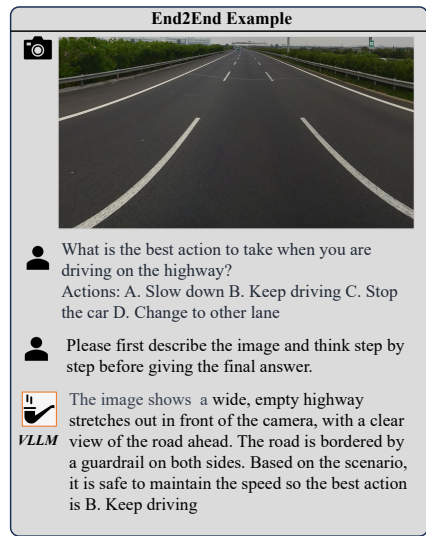


Figure 5: An example of end-to-end decision making.

Once equipped with sufficient information, the LLM proposes the final action, accompanied by its underlying rationale. HOLMES emulates the detective game process, where one alternates between searching for clues using various tools and analyzing them before arriving at a conclusion.

## 5 EXPERIMENTS

### 5.1 CONFIGURATIONS

**End2End.** Under this setting same image and prompts are provided to different VLLMs. Additionally, the non-visual information “items in hand” and “items in inventory” for domestic and game domains are directly given to the models in the prompt since these information is hard to perceive from the image and is easy to obtain from the simulation environments. We would also make the prompts we use open-source for fair and convenient evaluation.

We compare four different models, InstructBLIP-Vicuna-13B<sup>3</sup>, MMICL-FLANT5XXL<sup>4</sup>, QwenVL-Chat<sup>5</sup> and GPT4-Vision<sup>6</sup>. We apply default inference configurations for the corresponding models.

**HOLMES.** In HOLMES framework, the LLM is required to continuously invoke various APIs and retrieve their return information. To streamline the evaluation process, we initially execute all APIs for every instance in PCA-EVAL, storing the result for each instance. This approach allows us to directly access the specific result of a given API without the need to run the model each time an evaluation is conducted. We would also make the API results open-source together with the benchmark. The description and implementation details of the APIs are listed in Appendix B.

We compare three LLMs: Vicuna<sup>7</sup>, ChatGPT-3.5-Turbo and GPT4<sup>8</sup>. However we found Vicuna models lack the capability to call various APIs for information gathering, thus we have only reported the results for ChatGPT and GPT4. We anticipate supplementing these results as soon as open-source models become available, which can understand API descriptions and correspondingly call different APIs.

### 5.2 EVALUATION

PCA-Eval assesses embodied decision-making through three distinct lenses: perception, cognition, and action. The scores we reported in Table 1 rely on the consensus score from three human evaluators. We compute the average kappa correlation coefficient for these evaluators, resulting in 0.91 for the Perception Score and 0.88 for the Cognition Score. These figures indicate a good consistency in the evaluation process.

### 5.3 MAIN RESULTS

We evaluate various methods and models on the PCA-EVAL benchmark, as shown in Table 1.

In the upper block concerning End2End-VLLMs, the recently unveiled closed-source model, GPT-4V, outperforms existing open-source models by achieving the highest scores of 0.84, 0.74, and 0.74 in the perception, cognition, and action dimensions respectively. This performance represents a 26% action score improvement over its open-source counterpart, MMICL. The impressive performance of GPT-4V is primarily attributed to its exceptional ability to perceive visual information across different domains, particularly in the challenging game domain.

We also assessed the performance of embodied decision making using our HOLMES system.

As shown in the bottom block of the table, the HOLMES system, based on GPT4, achieves an Action Score of 0.71, matching the performance of GPT-4V (0.74). This suggests that the HOLMES system

<sup>3</sup><https://github.com/salesforce/LAVIS/tree/main/projects/instructblip>

<sup>4</sup><https://huggingface.co/BleachNick/MMICL-Instructblip-T5-xxl>

<sup>5</sup><https://huggingface.co/Qwen/Qwen-VL-Chat>

<sup>6</sup><https://chat.openai.com>

<sup>7</sup><https://huggingface.co/lmsys>

<sup>8</sup><https://platform.openai.com>



Method	Model	Traffic			Domestic			Game			Average		
		P	C	A	P	C	A	P	C	A	P	C	A
End2End	InstructBLIP <sup>†</sup>	-	-	0.42	-	-	0.41	-	-	0.24	-	-	0.36
	MMICL <sup>†</sup>	-	-	0.63	-	-	0.51	-	-	0.29	-	-	0.48
	QwenVL-Chat <sup>†</sup>	-	-	0.59	-	-	0.55	-	-	0.24	-	-	0.46
	GPT-4V <sup>‡</sup>	0.75	0.73	0.78	0.81	<b>0.69</b>	<b>0.67</b>	<b>0.95</b>	<b>0.79</b>	<b>0.77</b>	0.84	<b>0.74</b>	<b>0.74</b>
HOLMES	ChatGPT <sup>‡</sup>	0.75	0.68	0.66	<b>0.88</b>	0.52	0.50	0.78	0.40	0.36	0.80	0.53	0.51
	GPT4 <sup>‡</sup>	<b>0.87</b>	<b>0.82</b>	<b>0.82</b>	0.85	0.61	0.56	0.91	0.77	0.74	<b>0.88</b>	0.73	0.71

Table 1: Main results on PCA-EVAL. Models with <sup>†</sup> are fully open-source. Models with <sup>‡</sup> only provide API to access. P, C, and A represent Perception, Cognition, and Action Scores, respectively. For the open-source models in End2End setting, we find it hard to prompt them to output correct cross-modal reasoning information, so their Perception and Cognition scores are not reported.

is proficient in understanding the task goal, breaking down the larger goal into multiple smaller steps, and accurately invoking the relevant APIs to accomplish each step.

Specifically, the GPT4-HOLMES system can identify key concepts in an image through the results returned by APIs such as *list\_nearby\_mobs\_in\_minecraft()*. As a result, the system achieves an average Perception Score of 0.88, surpassing GPT-4V’s 0.84. However, when compared to End2End methods, HOLMES relies on multi-step reasoning for the final decision. This approach can lead to the accumulation of reasoning errors, resulting in a lower Cognition Score in both Domestic and Game domains.

## 6 DISCUSSION

### 6.1 COMPARISON BETWEEN END2END AND HOLMES

We conduct an analysis and comparison of the outputs generated by the End2End method with GPT4-Vision, as well as the HOLMES method with GPT4. Our findings indicate that the End2End method effectively mitigates information loss during the modality conversion process. As illustrated in Figure 6a, an image depicts a road with several nearby cars. GPT4-Vision is capable of discerning that these cars are situated in a safe space, thereby suggesting that the driver can continue driving.

Conversely, GPT4, while aware of the number of cars, lacks information about their spatial relation, leading it to recommend slowing down. 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-Vision 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.

### 6.2 ALIGNMENT BETWEEN AGENT DECISIONS AND HUMAN VALUES

We have observed instances where the decisions made by the agent contradict human values. For instance, consider the scenario depicted in Figure 6b. The image illustrates a crosswalk devoid of pedestrians. The appropriate response in this situation would be to slow 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 unoccupied, 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. We believe it is crucial for embodied agents to make decisions that are in harmony with human values, rather than solely focusing on maximizing their advantage.

### 6.3 LIMITATION AND FUTURE WORK

The current scope of PCA-EVAL is confined to merely three domains, with a cap of 100 instances per domain. One of our future work aims to broaden this scope to encompass more domains and em-



(a) A Comparison between GPT4-V and GPT4-HOLMES

(b) An example of decision misalignment.

Figure 6: Case studies.

bodied environments where MLLMs could keep getting feedback. Furthermore, we plan to increase the number of instances for both the existing and newly introduced domains.

## 7 CONCLUSION

In this study, we present PCA-EVAL, a comprehensive evaluation benchmark for embodied decision-making that gauges performance in perception, cognition, and action, thereby offering an all-encompassing measure for various embodied agents. We conduct a systematic comparison between End2End embodied decision-making and HOLMES, a multi-agent cooperation framework developed by us. Our findings reveal that MLLM, when applied with the end2end method, surpasses the top-performing model in HOLMES, GPT-4, in terms of decision accuracy and cognition score. However, it is crucial to underscore that this superior performance is specific to the GPT4-Vision model, which significantly outperforms the open-source state-of-the-art VLLMs. These results and subsequent analysis underscore the necessity for ongoing exploration in embodied decision-making and the development of open-source MLLMs to ensure wider accessibility and progress in the field.

## REFERENCES

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- Xiaoyu Chen, Shenao Zhang, Pushi Zhang, Li Zhao, and Jianyu Chen. Asking before action: Gather information in embodied decision making with language models. *arXiv preprint arXiv:2305.15695*, 2023.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.
- Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model, 2023a.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023b.
- Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. URL [https://openreview.net/forum?id=rc8o\\_j8I8PX](https://openreview.net/forum?id=rc8o_j8I8PX).
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 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*, pp. 6325–6334, 2017.
- Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, Wenhai Wang, Lewei Lu, Xiaosong Jia, Qiang Liu, Jifeng Dai, Yu Qiao, and Hongyang Li. Planning-oriented autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pp. 9118–9147. PMLR, 2022a.

- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Noah Brown, Tomas Jackson, Linda Luu, Sergey Levine, Karol Hausman, and Brian Ichter. Inner monologue: Embodied reasoning through planning with language models. In arXiv preprint arXiv:2207.05608, 2022b.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. Advances in neural information processing systems, 35:22199–22213, 2022.
- Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. AI2-THOR: An Interactive 3D Environment for Visual AI. arXiv, 2017.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597, 2023a.
- 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. M<sup>3</sup>it: A large-scale dataset towards multi-modal multilingual instruction tuning. arXiv preprint arXiv:2306.04387, 2023b.
- Minghao Li, Feifan Song, Bowen Yu, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. Api-bank: A benchmark for tool-augmented llms, 2023c.
- Shuang Li, Xavier Puig, Chris Paxton, Yilun Du, Clinton Wang, Linxi Fan, Tao Chen, De-An Huang, Ekin Akyürek, Anima Anandkumar, et al. Pre-trained language models for interactive decision-making. Advances in Neural Information Processing Systems, 35:31199–31212, 2022.
- Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Y. Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. ArXiv, abs/2203.02053, 2022. URL <https://api.semanticscholar.org/CorpusID:247244904>.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. 2023a.
- Ruibo Liu, Ruixin Yang, Chenyan Jia, Ge Zhang, Denny Zhou, Andrew M Dai, Diyi Yang, and Soroush Vosoughi. Training socially aligned language models in simulated human society. arXiv preprint arXiv:2305.16960, 2023b.
- Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. arXiv preprint arXiv:2304.09842, 2023.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 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, pp. 3195–3204, 2019.
- 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. Webgpt: Browser-assisted question-answering with human feedback, 2022.
- OpenAI. 2022. URL <https://chat.openai.com/>.
- OpenAI. 2023a. URL <https://openai.com/blog/chatgpt-can-now-see-hear-and-speak>.
- OpenAI. Gpt-4 technical report, 2023b.
- Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. arXiv preprint arXiv:2304.03442, 2023.

- Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, et al. Tool learning with foundation models. [arXiv preprint arXiv:2304.08354](#), 2023.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. [OpenAI blog](#), 1(8):9, 2019.
- Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. [arXiv](#), 2018.
- Shuhuai Ren, Lei Li, Xuancheng Ren, Guangxiang Zhao, and Xu Sun. Delving into the openness of CLIP. In [Findings of the Association for Computational Linguistics: ACL 2023](#). Association for Computational Linguistics, July 2023a. URL <https://aclanthology.org/2023.findings-acl.610>.
- Shuhuai Ren, Aston Zhang, Yi Zhu, Shuai Zhang, Shuai Zheng, Mu Li, Alex Smola, and Xu Sun. Prompt pre-training with twenty-thousand classes for open-vocabulary visual recognition. [arXiv preprint arXiv:2304.04704](#), 2023b.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools, 2023.
- Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks. In [The IEEE Conference on Computer Vision and Pattern Recognition \(CVPR\)](#), 2020. URL <https://arxiv.org/abs/1912.01734>.
- 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. Restgpt: Connecting large language models with real-world restful apis, 2023.
- Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. A corpus of natural language for visual reasoning. In [Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics \(Volume 2: Short Papers\)](#), pp. 217–223, 2017.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. [arXiv preprint arXiv:2302.13971](#), 2023.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. [arXiv preprint arXiv:2305.16291](#), 2023a.
- Peiyi Wang, Lei Li, Liang Chen, Feifan Song, Binghuai Lin, Yunbo Cao, Tianyu Liu, and Zhifang Sui. Making large language models better reasoners with alignment. [arXiv preprint arXiv:2309.02144](#), 2023b.
- Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. Large language models are not fair evaluators. [arXiv preprint arXiv:2305.17926](#), 2023c.
- Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. [ArXiv, abs/2302.01560](#), 2023d. URL <https://api.semanticscholar.org/CorpusID:256598146>.
- Wayve. Lingo. 2023. URL <https://wayve.ai/thinking/lingo-natural-language-autonomous-driving/>.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models, 2022.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
- Chenfei Wu, Sheng-Kai Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *ArXiv*, abs/2303.04671, 2023. URL <https://api.semanticscholar.org/CorpusID:257404891>.
- 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. The rise and potential of large language model based agents: A survey, 2023.
- Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *ArXiv*, abs/2303.11381, 2023. URL <https://api.semanticscholar.org/CorpusID:257637012>.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In The Eleventh International Conference on Learning Representations, 2022.
- Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh Murthy, Zeyuan Chen, Jianguo Zhang, Devansh Arpit, et al. Retroformer: Retrospective large language agents with policy gradient optimization. *arXiv preprint arXiv:2308.02151*, 2023.
- Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and Zongqing Lu. Plan4MC: Skill reinforcement learning and planning for open-world Minecraft tasks. *arXiv preprint arXiv:2303.16563*, 2023.
- Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng Wang, Wenjuan Han, and Baobao Chang. Mmicl: Empowering vision-language model with multi-modal in-context learning. *arXiv preprint arXiv:2309.07915*, 2023.
- 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. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023a.
- Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, Yu Qiao, Zhaoxiang Zhang, and Jifeng Dai. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory. *arXiv preprint arXiv:2305.17144*, 2023b.
- Zhe Zhu, Dun Liang, Songhai Zhang, Xiaolei Huang, Baoli Li, and Shimin Hu. Traffic-sign detection and classification in the wild. In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.

## A EXAMPLES OF PCA-EVAL

### A.1 DATA DISTRIBUTION

The PCA-EVAL benchmark data distribution across various domains is outlined in Table 2.

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.

Table 2: Data Distribution in the PCA-EVAL Benchmark

<b>Domain</b>	<b>Task Type/Location</b>	<b>Instances</b>
Autonomous Driving	Traffic Sign Detection	44
	Car Detection	33
	Human Detection	30
	Weather Detection	9
	Road Detection	3
	Character Recognition	13
Domestic Robot	Living Room	31
	Dining Room	11
	Bedroom	6
	Bathroom	3
	Kitchen	37
	Corridor	12
Open-World Game	Find Objects	52
	Kill Enemies	6
	Craft Items	32
	Place Blocks	7
	Interact with Creatures	3

A.2 PCA-EVAL EXAMPLES

We list three examples of each domain from PCA-EVAL, as shown in Figure 7, Figure 8 and Figure 9.

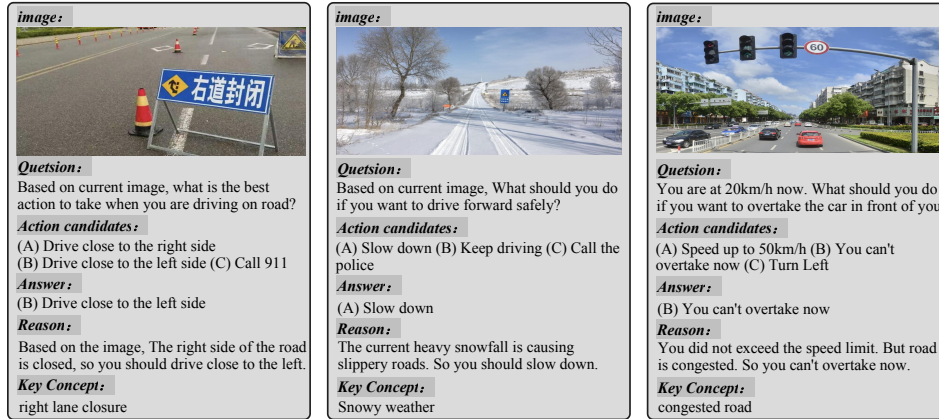


Figure 7: Three examples of PCA-EVAL in the autonomous driving domain.

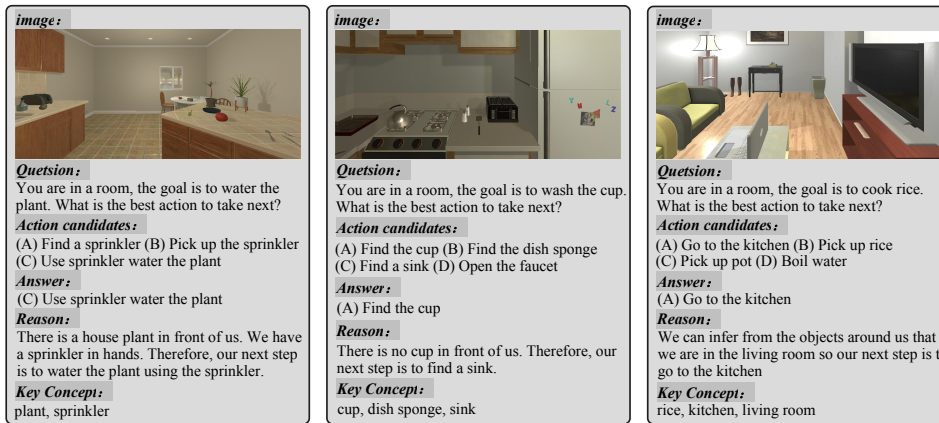


Figure 8: Three examples of PCA-EVAL in the domestic robot domain.



Figure 9: Three examples of PCA-EVAL in the open-world game domain.



## B API DESCRIPTION AND IMPLEMENTATION OF HOLMES

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

```

- *detect\_traffic\_sign()*: The detection of road traffic signs model utilize YOLO (Redmon & Farhadi, 2018) which trained on the Tsinghua-Tencent 100K dataset (Zhu et al., 2016). TT100K 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., 2023b), 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 PaddleOCR<sup>9</sup> 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<sup>10</sup> 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.

<sup>9</sup><https://github.com/PaddlePaddle/PaddleOCR/tree/release/2.7>

<sup>10</sup><https://huggingface.co/Salesforce/blip2-flan-t5-xl>

- *weather\_detection()*: Weather detection leverages a pre-trained ResNet50 model<sup>11</sup>, 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 Domain
2 def object\_detection():
3     """
4     Detects objects in current view, which you don't need to find.
5     :return: list of detected objects, e.g. ['chair', 'table']
6     """
7     pass
8
9 def list\_items\_in\_hands():
10    """
11    Lists items in your hand, which you don't need to pick up
12    :return: list of items in hand, e.g. ['coffee cup', 'milk']
13    """
14    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     """
4     Lists nearby mobs in Minecraft.
5     :return: list of nearby mobs, e.g. ['creeper', 'pig']
6     """
7     pass
8
9 def list\_inventory\_information():
10    """
11    Lists inventory information of the player in Minecraft.
12    :return: list of inventory information with number, e.g. [('diamond',
13    64), ('iron', 32)]
14    """
15    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.

<sup>11</sup><https://github.com/mengxianglong123/weather-recognition>

## C AUTOMATIC EVALUATION

<p>[Question]: {question}                  [Action Choices]: {actions}                  [Agent Answer]: {model_output}                  [Correct Action]: {true_action}                  [Key Concepts]: {key_concept}                  [Reference Reasoning Process]: {reason}                  [System]</p> <p>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.</p> <p>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]:</p> <ol style="list-style-type: none"> <li>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.</li> <li>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.</li> <li>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.</li> </ol> <p>Please note that there are only scores of 0 and 1.</p> <p>You should carefully compare the [Agent Answer] with the [Correct Action], [Key Concepts] and [Reference Reasoning Process] to give your assessment.</p> <p>You need first to give your assessment evidence and then the scores.</p> <p>Your output MUST contain 6 lines with the following format:</p> <p>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)</p>
--

Table 3: We utilize the template 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.