

# MOTIVEBENCH: How Far Are We From Human-Like Motivational Reasoning in Large Language Models?

Anonymous ACL submission

## Abstract

Large language models (LLMs) have been widely adopted as the core of agent frameworks in various scenarios, such as social simulations and AI companions. However, the extent to which they can replicate human-like motivations remains an underexplored question. Existing benchmarks are constrained by simplistic scenarios and the absence of character identities, resulting in an information asymmetry with real-world situations. To address this gap, we propose MOTIVEBENCH, which consists of 200 rich contextual scenarios and 600 reasoning tasks covering multiple levels of motivation. Using MOTIVEBENCH, we conduct extensive experiments on seven popular model families, comparing different scales and versions within each family. Our analysis reveals several notable findings, such as the difficulty LLMs face in reasoning about "love & belonging" motivations and the tendencies of LLMs toward excessive rationality and idealism. These insights highlight a promising direction for future research on the humanization of LLMs.

## 1 Introduction

Motivation is commonly conceptualized as an internal drive or psychological force that influences individuals to initiate and sustain goal-oriented activities (Hagger and Chatzisarantis, 2005; Brehm, 2014). It serves as a key explanatory factor for understanding why people initiate, continue, or terminate specific behaviors at any given scenarios (Kazdin et al., 2000). Types of motivation include intrinsic motivation, driven by internal factors like values or preferences, and extrinsic motivation, influenced by external rewards or punishments (Ryan and Deci, 2000; Radel et al., 2016).

Mimicking human behavior in specific scenarios has been a crucial task for autonomous agents, forming the foundation for various applications such as problem-solving, testing, and simulation (Schatzmann et al., 2007). Previous studies em-

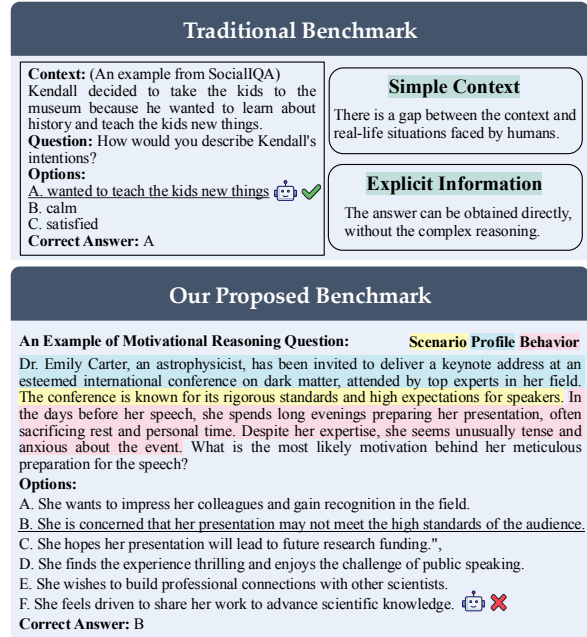


Figure 1: The difference between the existing traditional benchmark, such as SOCIALQA from Sap et al. (2019) and our proposed MOTIVEBENCH.

ploy either rule-based (Keizer et al., 2010; Wilkins, 2014) or machine learning approaches (Asri et al., 2016; Kreyszig et al., 2018) to replicate human interactions in isolated and controlled environments. With the advent of large language models (LLMs) like GPT-4 (Achiam et al., 2023) and LLaMA (Dubey et al., 2024), an increasing number of studies have been adopting LLM-based agents (Aher et al., 2023; Argyle et al., 2023; Boiko et al., 2023; Kang and Kim, 2023; Mehta et al., 2023; Hong et al., 2023; Wang et al., 2024c) due to the remarkable capabilities of LLMs in general problem-solving, reasoning, and autonomous action-taking. However, a critical question remains underexplored: **Can current LLMs truly understand and exhibit human-like motivations and behaviors?** The complexity of human behavior dynamics poses new challenges for LLMs, which are distinct from the challenges in understanding and

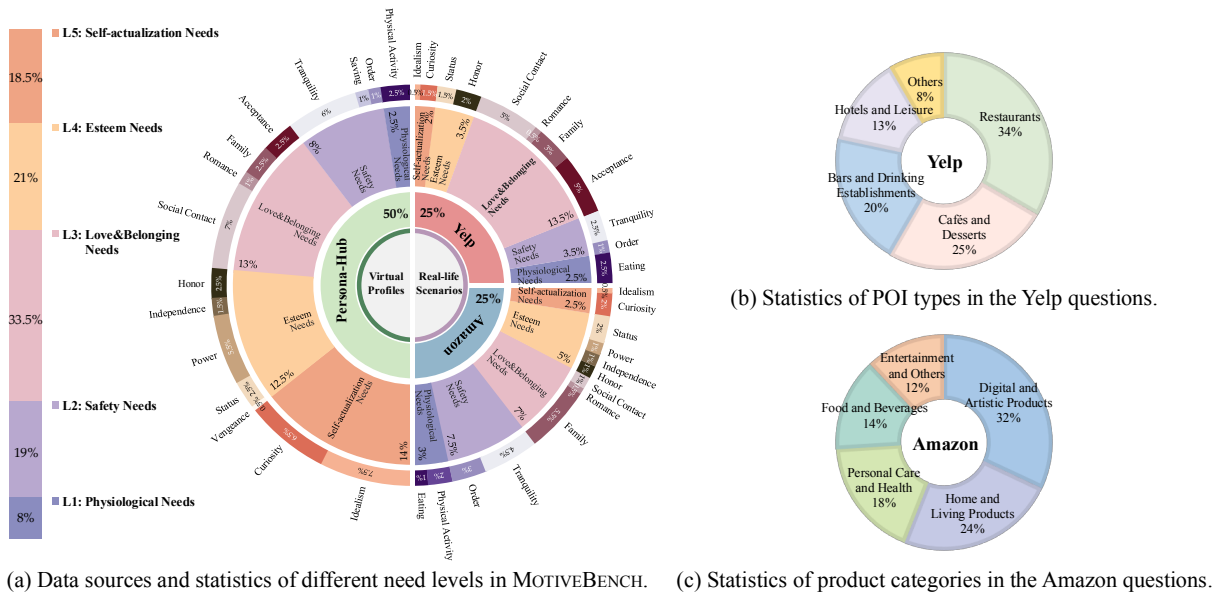


Figure 2: The statistical overview of MOTIVEBENCH. It contains 200 diverse profiles and real-world scenarios, along with 600 motivational and behavioral reasoning questions, covering multiple finely-grained levels of needs.

manipulating the physical world. Delving deeply into this research question can lay a reliable foundation for LLM-based social simulations.

Recent literature has attempted to explore the gap between LLMs and real humans from a narrow behavioral perspective. Xie et al. (2024) found the feasibility of simulating human trust behavior with LLM agents in Trust Game (Berg et al., 1995). Zhou et al. (2023) constructed an interactive environment to analyze the differences between current LLMs and humans in social interactions (Lee et al., 2024). These methods focus on specific subdomains and fail to provide insights into the broader human-like behavior at a macro level. In the study of motivations or behaviors in LLMs, Sap et al. (2019) introduced SOCIALIQA, a benchmark focused on commonsense reasoning in social contexts. While it includes some basic reasoning tasks about the intentions behind behaviors, it lacks more detailed, comprehensive, and challenging assessments, as shown in Figure 1. Existing similar benchmarks (Rashkin et al., 2018a,b; Talmor et al., 2018) also exhibit several clear limitations: (1) **Simplistic contexts**, lacking detailed scenarios and character profiles, leading to information asymmetry compared to real-world situations; (2) **Overly explicit information**, with tasks solvable through basic pattern matching without requiring human-like reasoning; and (3) **Limited theoretical grounding**, failing to systematically capture the multi-level nature of human motivation.

To address the above challenges, we pro-

pose MOTIVEBENCH, a comprehensive evaluation benchmark, consisting of 200 diverse profiles and scenarios, along with 600 motivational and behavioral reasoning questions. Figure 2 shows an overview of MOTIVEBENCH. We strive to cover and balance the proportions of different levels of needs in the benchmark, ultimately formulating the questions. Specifically, it has the following advantages: (1) **Diverse scenarios, profiles, motivations, and behaviors**. We utilize diverse profiles from the Persona-Hub (Ge et al., 2024) dataset, along with real-world motivation and behavior data from platforms like Amazon and Yelp, as the basis for question generation. (2) **Human-in-the-loop multi-agent framework to enhance efficiency and quality**. We propose a multi-agent collaboration framework that efficiently generates high-quality questions across a range of difficulties, requiring minimal human effort to ensure validity. (3) **Grounded in authoritative psychological theories to ensure comprehensive evaluation**. Our test questions cover the five levels of Maslow’s Hierarchy of Needs (Maslow, 1943), as well as the 16 basic desires of human nature from the Reiss Motivation Profile (Reiss, 2004).

To the best of our knowledge, MOTIVEBENCH is the first systematic benchmark designed to evaluate the human-like motivation-behavior reasoning capabilities of LLMs. We conduct comprehensive experiments to draw conclusions. Beyond the quantified results, our experiments reveal several novel insights, such as significant differences in the mo-

tivational reasoning processes of LLMs compared to humans, and the limitations of them for data annotation in the context of human social behaviors. We hope that our research provides practical guidelines for applying LLMs in various social simulations and contributes to future improvements in the humanization of LLMs. Dataset, benchmark, and code will be released upon acceptance of this paper to benefit research in this direction.

## 2 MOTIVEBENCH Preliminaries

### 2.1 Three Types of Reasoning Tasks

Generally, in a specific scenario, an individual with a certain profile will perform a behavior based on a particular motivation. We define this as a complete behavioral quadruple: **Scenario, Profile, Motivation, and Behavior**.

In this quadruple, the scenario provides the context and external triggers (Yang et al., 2009). Profile shapes an individual’s understanding of the behavior and the way they act (Bandura, 2001; Eagly and Wood, 2012). Motivation is the internal driving force behind an individual’s actions, based on their needs, goals, or emotional state (Deci and Ryan, 2012; Shayganfar et al., 2016). Behavior is the result of the interaction between scenario, profile, and motivation (Graham, 1991). Therefore, we define three types of reasoning tasks:

1) **Motivational Reasoning Question**. Given a specific scenario, profile, and detailed behavior information, the task is to infer the motivation behind the individual’s behavior.

2) **Behavioral Reasoning Question**. Given a specific scenario, profile, and detailed motivation information, the task is to infer the most likely behavior the individual would perform.

3) **Motive&Behavior Reasoning Question**. This more challenging reasoning task closely aligns with the ultimate goal of autonomous agents, which is to infer the most reasonable motivation and corresponding behavior when only the scenario and profile are provided.

### 2.2 Fine-Grained Needs Hierarchy

Maslow’s hierarchy of needs theory (Maslow, 1943) suggests that human actions are driven by various needs, which are divided into five levels: Physiological Needs, Safety Needs, Esteem Needs, Love and Belonging Needs, and Self-actualization Needs. When lower-level needs are met, individuals will seek to fulfill higher-level needs (Jr, 1991).

Furthermore, Reiss (2004) proposes 16 more granular categories to provide a broader and more informative range of motivations. These include Curiosity, Idealism, Honor, Independence, Power, Status, Vengeance, Acceptance, Family, Romance, Social Contact, Order, Saving, Tranquility, Eating, and Physical Activity. Although the Reiss theory offers more detailed insights into motivation, the broader range of abstract concepts can be difficult to manage. Inspired by Rashkin et al. (2018a), we adopt a hybrid method in which the Reiss Motive Profile labels are categorized as sub-categories within Maslow’s framework.

## 3 MOTIVEBENCH Construction

We construct MOTIVEBENCH from scratch through collaboration between LLMs and humans. We deliberately avoid relying on existing scales or test items from psychology or sociology to prevent potential data leakage or contamination issues. Figure 3 illustrates the overall process, with each step introduced in the subsequent subsections.

### 3.1 Data Collection and Pre-processing

To obtain diverse scenarios and profiles, we collect data from Persona-Hub proposed by Ge et al. (2024), as well as real-life platforms like Amazon (Hou et al., 2024) and Yelp<sup>1</sup>.

Persona-Hub contains diverse profiles, such as “A fellow astrophysicist who specializes in the study of dark matter and provides valuable insights and critiques to the author’s research.” Based on them, we can synthesize a diverse range of scenarios, motivations, and behaviors using the fine-grained hierarchy of needs outlined in Section 2.2.

In addition, review texts on platforms like Amazon and Yelp contain abundant real user intent data, such as their reasons for visiting a specific point of interest (POI) or purchasing a product. This provides a valuable reference for collecting data grounded in real-world scenarios. Therefore, we first collect review data from a wide range of domains, including 33 product domains from Amazon and 22 business categories from Yelp, ensuring data diversity. We then employ LLMs, such as LLaMA3.1-70B and Qwen2.5-72B, to filter high-quality reviews that align with Maslow’s hierarchy of needs. To capture deeper-level needs, we focus on motivations rooted in real-life contexts, rather than the product’s inherent features. For example,

<sup>1</sup><https://www.yelp.com/dataset>

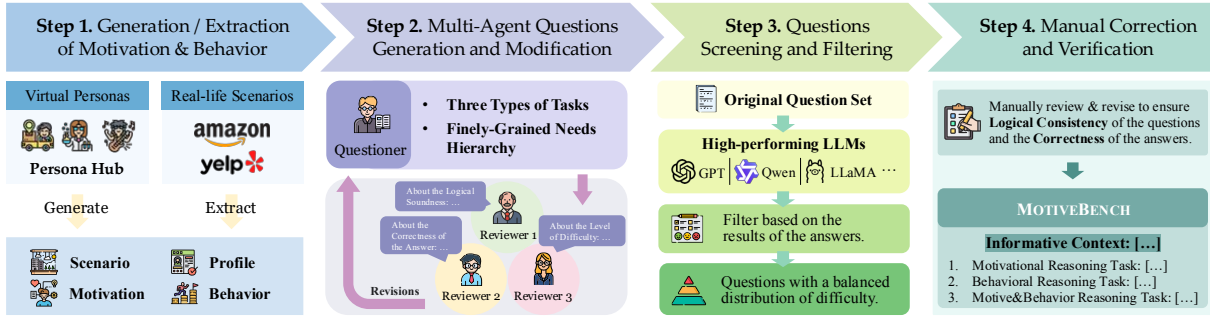


Figure 3: A step-by-step questions generation and correction pipeline using AI-Human collaboration framework.

in reviews of smartphones, we prioritize motivations such as purchasing for better communication with family, reflecting "social needs", or selecting a phone with a long battery life for convenience during travel, reflecting "safety needs", rather than only focusing on product attributes (such as performance or appearance) as the motivation for purchase. An example is shown in Figure 4.

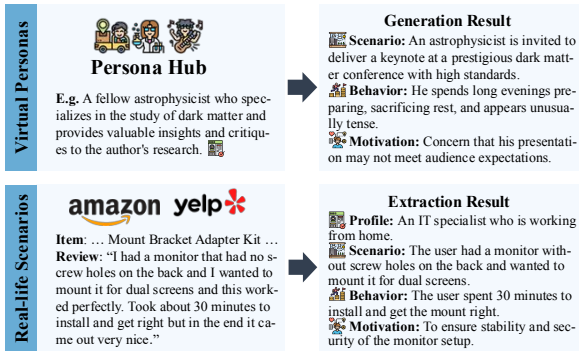


Figure 4: Generation and extraction of motivations and behaviors from virtual personas and real-life scenarios.

### 3.2 Questions Generation and Modification

Existing benchmarks mostly rely on manual construction, which requires significant effort (Sabour et al., 2024; Chen et al., 2024). To reduce financial and labor costs, we employ a multi-agent framework to formulate questions and generate options.

Specifically, we begin by using an LLM-based questioner to formulate complete question content based on the quadruples obtained in the previous stage. Next, three LLM-based reviewers provide feedback from: 1) the logical soundness of the question, 2) the correctness of the answer, and 3) the level of difficulty. The feedback is then compiled and sent back to the questioner for revisions. This process is iterated until no further suggestions are made, or the iteration threshold is reached. To minimize the hidden biases introduced by a single LLM, we use a diverse set of models, such as LLaMA3.1-70B and Qwen2.5-72B, for question

creation and modification. The detailed prompts of our agents can be found in the Appendix E.

### 3.3 Questions Screening and Filtering

A well-rounded benchmark should include questions of varying difficulty levels to comprehensively evaluate the performance of different LLMs. Controlling difficulty during the initial question generation phase is challenging. Therefore, we use high-performing LLMs (e.g., GPT-4o, LLaMA3.1-70B, and Qwen2.5-72B) to answer the questions. By analyzing their responses, we label the difficulty level of each question. Based on this, we categorize the questions into difficult, medium, and easy levels, ensuring a balanced distribution of all three difficulty levels in the final benchmark. Ultimately, We curate a subset of high-quality, diverse scenarios from the original dataset, retaining 100 scenarios for Person-Hub, 50 for Amazon, and 50 for Yelp. These scenarios span different difficulty levels and provide balanced coverage.

### 3.4 Manual Correction

After filtering the questions, we authors manually review and revise each question to ensure high quality of our MOTIVEBENCH. We carefully refine the logical consistency of questions and the accuracy of options, it takes about 6 minutes per question. Due to the hallucination issues (Huang et al., 2023) in LLMs that often introduce logical or factual errors, this step is essential. Our analysis shows that generated questions often struggle with ensuring clear distinctions between correct and incorrect answers. This likely reflects the limitations of LLMs in capturing the nuanced and complex nature of human reasoning. The final statistics of MOTIVEBENCH are shown in Table 1, where the context length is significantly longer than that of previous benchmarks (e.g., SOCIALIQA in Sap et al. (2019): 20.16 tokens per context, and COMMONSENSEQA in Talmor et al. (2018): 13.41 tokens per context).

	#S	#Q	Average Token Count		
			Context	Cor.Opt	Fal.Opt
<b>MOTIVEBENCH</b>	<b>200</b>	<b>600</b>	<b>96.45</b>	<b>13.83</b>	<b>13.31</b>
Persona	100	300	97.13	15.55	14.60
Amazon	50	150	87.19	12.56	12.51
Yelp	50	150	104.35	11.64	11.51

Table 1: Data statistics of our proposed MOTIVEBENCH. (#S: number of scenarios, #Q: number of questions, Cor.Opt: correct options, and Fal.Opt: false options.)

## 4 Experiments

### 4.1 Experimental Setup

For each scenario, there are three types of questions: motivational reasoning, behavioral reasoning, and motive&behavior reasoning. Each question is presented in the form of MCQ. The scenario is only considered correct when all three questions within the same scenario are answered correctly. The specific format is shown in Appendix B.

We evaluate LLMs in two settings: one using vanilla prompting with task instructions (Base), and the other employing chain-of-thought reasoning (CoT) (Wei et al., 2022). The prompts we used are detailed in Appendix E. Given that LLMs have been shown to exhibit a bias towards the order of choices (Zheng et al., 2023), we introduce random variations in the choice order by generating 6 permutations. This ensures the correct option appears in all possible positions, while the incorrect options are randomly shuffled each time. We report the average of the 6 results as the final performance.

We evaluate 29 popular LLMs, as listed in Appendix C. For all open-source models, we use the vLLM<sup>2</sup> inference framework and set the temperature parameter to 0 to ensure result stability. For closed-source models, we access them through the Azure OpenAI API<sup>3</sup>.

### 4.2 Main Results

Table 2 summarizes the performance of various LLMs across three domains, with detailed task-specific results in Appendix D. Below, we analyze the results and highlight several key findings.

First of all, among all the LLMs we evaluated, GPT-4o demonstrates the strongest capability in MOTIVEBENCH. Notably, within the open-source model series, the Qwen2.5 series demonstrates strong performance, with smaller models (14B, 32B) achieving capabilities comparable to the 72B

<sup>2</sup><https://docs.vllm.ai>

<sup>3</sup><https://azure.microsoft.com>

model and even GPT-4o. Similarly, LLaMA 3.1-70B also shows good results. However, other series, especially Baichuan2, exhibit weaker reasoning capabilities for motivation and behavior tasks. From the perspective of model size, small-scale models (<10B) achieve an average accuracy of 55.37%, medium-scale models (10B-34B) 68.24%, and large-scale models (>34B) 76.42%. These findings suggest a clear trend of improved motivational and behavioral reasoning ability as model size increases. The pattern is better visualized in Figure 5.

Secondly, CoT does not enhance the performance of LLMs in MOTIVEBENCH. In fact, results from most models indicate that CoT may lead to a decrease in performance. This effect is particularly pronounced in models with smaller parameter sizes ( $\leq 34B$ ), where accuracy drops by 7.55%, compared to a 1.05% decrease in larger models ( $\geq 70B$ ). This decline may occur because CoT simplifies tasks, but when the model’s reasoning diverges from human cognitive patterns on motivations and behaviors, it hampers alignment with human cognition, reducing performance. In contrast, upgrading models and increasing their size significantly improve human-like motivational reasoning, as illustrated in Figure 5.

Another interesting finding is that, LLMs perform poorly in understanding “love & belonging” needs, which are related to emotional aspects. In Table 3 we break the overall score into five motivation aspects, as introduced in Section 2.2. While some studies suggest that LLMs can provide emotional value comparable to or even surpassing that of humans—for instance, the Replica chatbot reduced suicidal ideation for 3 percent of users (Maples et al., 2024)—they still exhibit limitations in emotional understanding and reasoning (Sabour et al., 2024). This may be attributed to: (1) LLMs excel in providing emotional value by mimicking surface-level language patterns, creating a sense of understanding without deep causal reasoning. (2) The expression of “love & belonging” needs in the text data is often implicit or ambiguous, as it primarily involves internal processes. Since the model lacks direct exposure to comprehensive and real-world social psychology case data, it struggles to handle such issues effectively.

### 4.3 In-Depth Analysis

**Insight 1: Comparison with Existing Benchmarks.** We aim to introduce a new evaluation dimension—Motive. To study the difference be-

Motive&Behavior Reasoning Ability	Vitual Profiles Persona-Hub		Real-life Scenarios				Overall	
	Persona-Hub		Amazon		Yelp		Overall	
LLMs	Base	CoT	Base	CoT	Base	CoT	Base	CoT
Baichuan2-7B-Chat	35.33	<b>32.50</b>	37.33	<b>35.00</b>	30.00	25.00	34.50	31.25
Baichuan2-13B-Chat	<b>43.83</b>	31.83	<b>50.00</b>	32.00	<b>47.00</b>	<b>31.00</b>	<b>46.17</b>	<b>31.67</b>
ChatGLM3-6B	43.17	33.67	42.67	27.67	39.67	33.33	42.17	32.09
GLM4-9B-Chat	<b>60.33</b>	<b>60.33</b>	<b>75.67</b>	<b>74.67</b>	<b>66.33</b>	<b>68.67</b>	<b>65.67</b>	<b>66.00</b>
Yi1.5-6B-Chat	36.67	47.33	47.00	54.33	45.33	54.33	41.42	50.83
Yi1.5-9B-Chat	61.17	55.33	68.00	64.00	64.67	62.33	63.75	59.25
Yi1.5-34b-Chat	<b>66.50</b>	<b>63.33</b>	<b>72.33</b>	<b>71.00</b>	<b>80.67</b>	<b>70.00</b>	<b>71.50</b>	<b>66.92</b>
Phi3-mini-4k-Instruct	59.83	46.50	69.00	47.67	59.67	49.33	62.08	47.50
Phi3-small-8k-Instruct	63.83	61.17	76.67	56.33	72.67	62.00	69.25	60.17
Phi3-medium-4k-Instruct	<b>73.50</b>	<b>67.00</b>	<b>79.00</b>	<b>68.00</b>	<b>81.33</b>	72.67	<b>76.83</b>	<b>68.67</b>
Phi3.5-mini-Instruct	61.00	58.17	72.67	65.67	67.00	60.33	65.42	60.59
Phi3.5-MoE-Instruct	71.50	56.83	71.67	67.33	73.67	<b>77.67</b>	72.09	64.67
Llama2-7B-Chat	19.17	21.83	20.67	24.33	14.33	27.33	18.34	23.83
Llama2-13B-Chat	47.50	36.33	52.67	38.33	46.67	44.33	48.59	38.83
Llama2-70B-Chat	52.67	61.33	60.33	60.67	53.67	63.33	54.84	61.67
Llama3.1-8B-Instruct	63.50	55.00	72.00	67.67	65.00	62.33	66.00	60.00
Llama3.1-70B-Instruct	<b>82.17</b>	<b>76.67</b>	<b>83.67</b>	<b>83.33</b>	<b>83.67</b>	<b>78.00</b>	<b>82.92</b>	<b>78.67</b>
Qwen-7B-Chat	45.83	40.50	54.00	50.33	46.67	43.67	48.08	43.75
Qwen-14B-Chat	63.17	57.00	69.33	63.33	67.67	59.33	65.84	59.17
Qwen-72B-Chat	71.33	69.00	81.33	77.33	73.33	68.67	74.33	71.00
Qwen2-7B-Instruct	70.17	69.83	73.00	74.33	75.67	75.00	72.25	72.25
Qwen2-72B-Instruct	79.00	74.83	81.67	77.67	80.33	75.00	80.00	75.58
Qwen2.5-7B-Instruct	71.67	68.67	73.33	70.33	66.67	65.67	70.84	68.34
Qwen2.5-14B-Instruct	79.00	73.33	80.33	79.33	84.00	77.33	80.58	75.83
Qwen2.5-32B-Instruct	81.83	75.67	<b>88.33</b>	<b>86.33</b> <sup>†</sup>	<b>85.33</b>	80.00	<b>84.33</b>	79.42
Qwen2.5-72B-Instruct	<b>83.33</b>	<b>77.83</b>	85.67	83.00	80.67	<b>82.33</b>	83.25	<b>80.25</b>
GPT-3.5-Turbo 1106	62.67	63.33	78.33	78.33	61.67	68.00	66.34	68.25
GPT-4o mini 2024-07-18	81.67	75.67	84.67	84.00	81.00	80.67	82.25	79.00
GPT-4o 2024-05-13	<b>86.33</b> <sup>†</sup>	<b>81.50</b> <sup>†</sup>	<b>89.33</b> <sup>†</sup>	<b>84.33</b>	<b>87.67</b> <sup>†</sup>	<b>82.67</b> <sup>†</sup>	<b>87.42</b> <sup>†</sup>	<b>82.50</b> <sup>†</sup>

Table 2: Evaluation Results for MOTIVEBENCH across 7 popular model families in 3 domains, including Base and CoT prompting. The best results in each series are highlighted in **Bold**, with the best overall results marked by <sup>†</sup>.

tween MOTIVEBENCH and existing benchmarks, we leverage LiveBench (White et al., 2024), which typically assess general capabilities like coding, mathematics, reasoning, data analysis, language comprehension, and instruction following.

Figure 6 shows the Pearson correlation coefficients (Cohen et al., 2009) of rankings across different ability dimensions for several popular LLMs. It is evident that Motive is distinct from existing evaluation dimensions, with an average correlation coefficient of 0.8175. This suggests that by introducing the Motive dimension, we can explore patterns or relationships in human capabilities that traditional evaluation metrics do not observe.

**Insight 2: Differences Between GPT-4o and Human in Motivation and Behavior.** In Table 2, we observe that GPT-4o is the most advanced model in our benchmark. Therefore, we are curious to investigate the situations in which this model fails

to demonstrate human-like intelligence. For questions answered incorrectly by GPT-4o, we examine its reasoning and thought processes. We have summarized the following findings:

**1) Over-Rationalization, Lacking Emotional Insight.** GPT-4o often relies on logical reasoning, neglecting broader practical experience or emotional context, leading to reasoning that may be disconnected from real-world complexities.

**2) Weak Logical Precision, Prone to General Assumptions.** GPT-4o’s reasoning can be overly simplistic, often based on general assumptions or external knowledge, without fully addressing the specific context or details of a situation.

**3) Overly Idealistic, Ignoring Complex Realities.** It tends to assume people follow social norms or moral codes, ignoring more complex or real-world challenges that could affect behavior.

**4) Lack of Awareness of Behavioral Impact.** It may prioritize actions that seem easy or plausible

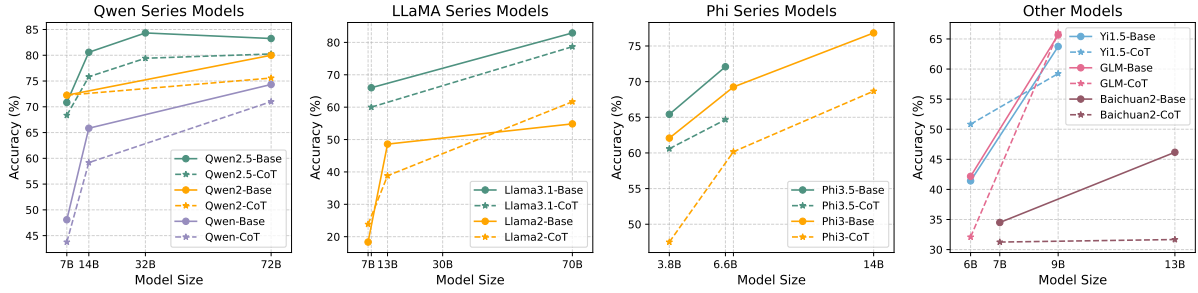


Figure 5: Scaling of various model families across their different versions and sizes in motivational reasoning ability.

Needs Hierarchy	GPT series		LLaMA series		Qwen series		Phi series		Yi series		AVG.
	4o	4o mini	3.1 70B	3.1 8B	2.5 72B	2.5 32B	3.5 MoE	3 Medium	1.5 34B	1.5 9B	
<b>L1: Physiological</b>	95.83	80.00	94.31	74.55	90.92	98.33	75.69	85.83	74.86	72.78	84.31
<b>L2: Safety</b>	90.49	83.85	85.61	69.24	86.71	87.49	83.71	78.50	69.88	71.05	80.65
<b>L3: Love &amp; Belonging</b>	84.81	72.43	77.80	60.78	85.74	80.94	67.35	73.75	68.11	57.83	72.95
<b>L4: Esteem</b>	89.08	79.10	80.90	60.69	78.75	84.65	60.95	70.73	72.12	58.00	73.50
<b>L5: Self-actualization</b>	86.40	80.06	86.76	69.31	86.67	81.90	77.17	83.04	75.48	69.52	79.63

Table 3: Hierarchy of needs-oriented evaluation results for different model families and their strongest models.

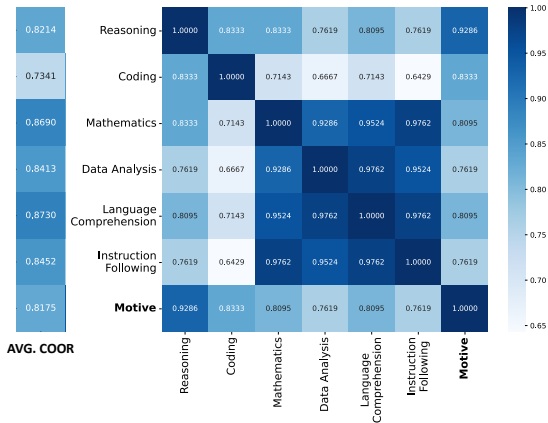


Figure 6: The correlation coefficients between the evaluation results of six general capabilities and the Motive.

but fail to consider their actual effectiveness or long-term impact in real-life scenarios.

Figure 7 presents an example of GPT-4o demonstrating weak logical coherence and overly idealistic. More cases can be found in the Appendix F.

**Insight 3: Limitations of LLMs in Complete Data Annotation.** Recently, many researchers have been curious about whether "large language models replace human annotators" or "replace human participants in social science experiments". By using our MOTIVEBENCH as the lens, we find that relying solely on LLMs for data annotation, particularly in the field of human social behaviors, presents several challenges:

**1) Logical or Factual Errors.** LLMs may generate inaccurate or misleading questions due to limited understanding of psychological theories, re-

sulting in the content of the questions or options lacking logical consistency. Therefore, we introduce a manual correction step to ensure the reliability of the reasoning questions.

**2) Limited Understanding of Human Dynamics.** The tasks we consider often involve nuanced psychological and sociological dimensions, which LLMs may struggle to accurately capture due to the complexity of human thought processes.

**3) Annotator-specific Bias.** Using a single LLM introduces biases that can distort the objective interpretation of psychological data, which is crucial for accurate research. To minimize such biases, our framework employs various models at different stages of the pipeline, and resorts to human labeling to revise questions as the final step.

**Insight 4: Trends of Motivational Intelligence in LLMs of Varying Sizes** Our experiments with LLMs of varying sizes reveal distinct trends in motivational intelligence. Smaller models (6B-13B) struggle with complex reasoning, long-chain inferences, and emotional understanding in the informative contexts we provide, often producing superficial responses. Medium-sized models (13B-34B) show improved performance, are able to capture motivational cues, provided better contextual understanding and somewhat human-like thinking patterns. In addition, large models (70B+), particularly GPT-4o and Qwen2.5-72B, excel in complex reasoning and nuanced motivational modeling. They integrate scenario and character's profile information effectively, understanding motivations,

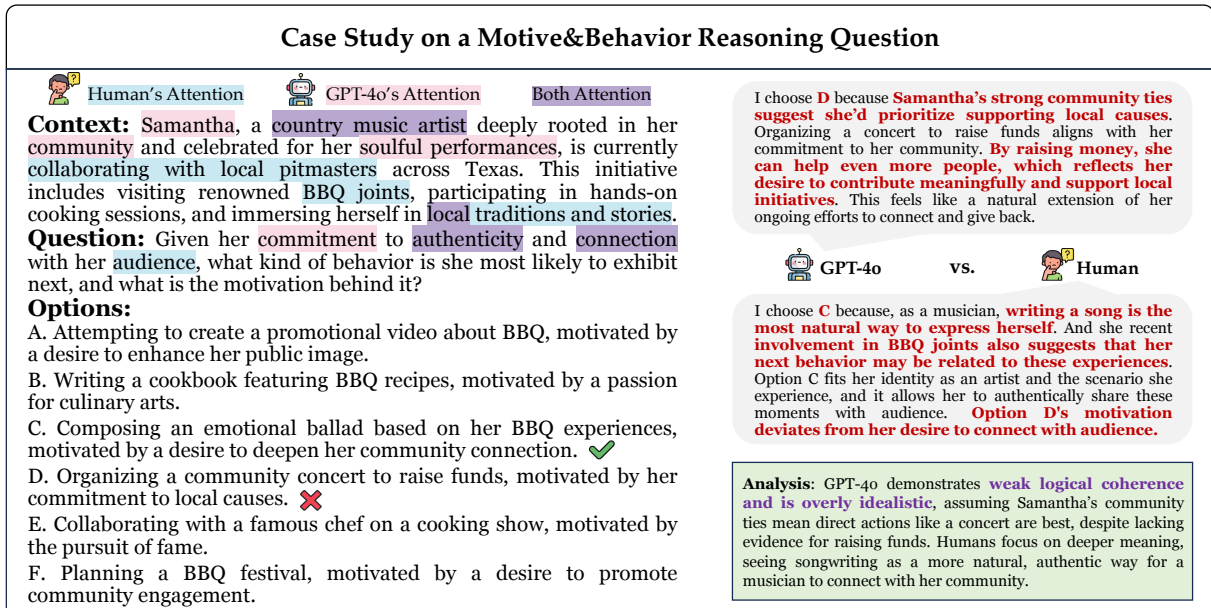


Figure 7: A case study on motivation and behavior reasoning: we analyze GPT-4o’s attention within the question and its explicitly generated reasoning process to uncover key differences between its approach and human cognition.

emotions, and logical relationships, thus showing behavior patterns closer to human intelligence.

However, current LLMs still differ from humans in some aspects of motivational reasoning (as shown in Insight 2), making them a major challenges in fully open-ended, dynamic scenarios, such as real-time strategy games or unpredictable social interactions. This represents the frontier of agent simulation and decision-making systems.

## 5 Related Work

Currently, the intellectual capabilities of LLMs have reached an unexpected level, but they also display a certain unreliability—sometimes becoming confused, while at other times demonstrating abilities that far exceed human in specific evaluations. But how far are they from human-like cognition? An increasing number of studies are now approaching this issue from Theory of Mind (ToM) perspective. Sap et al. (2022) find that GPT-3 lacks social intelligence by using SOCIALIQA (Sap et al., 2019) and ToMi (Le et al., 2019). Ullman (2023) demonstrate that small variations that preserve the principles of ToM can drastically alter the outcomes. Similar results can also be found from Shapira et al. (2023). Strachan et al. (2024) observe that GPT-4 excelled at identifying indirect requests, false beliefs, and misdirection, but struggled with detecting faux pas, and it still lags behind humans in overall ToM performance (Chen et al., 2024). In addition, some studies have conducted experiments from a

behavioral perspective. Xiao et al. (2023) reveal that current LLMs struggle to align their behaviors with assigned characters. Furthermore, Zhou et al. (2023) and Wang et al. (2024a) propose interactive and sandbox benchmarks, showing GPT-4 excels in conversational scenarios but struggles to exhibit social commonsense reasoning and deal with social tasks (Lee et al., 2024).

Different from existing research, we aim to explore a new dimension for evaluating—Motive. This dimension examines the alignment between LLMs and human behavior in dynamic, unstructured environments. MOTIVEBENCH incorporates rich scenarios, profiles, motivations, and behavioral data, offering a comprehensive assessment.

## 6 Conclusions

We introduce MOTIVEBENCH, the first systematic benchmark to evaluate the human-like motivational and behavioral reasoning ability of LLMs with detailed, realistic situations. Our results reveal that some advanced LLMs like GPT-4o generally demonstrate strong performance. However, from the perspective of needs hierarchy, most LLMs struggle with understanding "love & belonging" needs. Further in-depth analysis indicates that even the most advanced LLMs still exhibit deviations from human-like reasoning in motivation. By introducing MOTIVEBENCH, we aim to provide insights from a new dimension, enabling future models to exhibit more human-like cognition processes.



## 525 Limitations

### 526 Fully Automated Questions Generation

527 In the current pipeline for generating questions in  
528 MOTIVEBENCH, we still rely on human effort to  
529 manually check and revise the quality of questions.  
530 This approach poses challenges for automatically  
531 refreshing the benchmark to avoid data contamina-  
532 tion for future LLM releases and for scaling to a  
533 larger set of test questions.

534 To address these challenges, it is necessary to  
535 train a revision model using our existing manually  
536 corrected data, with the goal of fully automating  
537 the entire question generation pipeline. This would  
538 enable the benchmark to dynamically update it-  
539 self, thereby maximizing its value in the rapidly  
540 evolving era of LLMs. To achieve this, a potential  
541 approach involves leveraging advanced techniques,  
542 focusing particularly on fine-tuning existing pre-  
543 trained models with our manually corrected dataset.  
544 This process will help the model learn the intri-  
545 cate patterns and nuances required for high-quality  
546 question generation. Additionally, we plan to incor-  
547 porate continuous learning mechanisms, allowing  
548 the model to adapt to new data and evolving trends.  
549 By doing so, we aim to enhance not only the ac-  
550 curacy and relevance of the generated questions  
551 but also ensure that the benchmark remains aligned  
552 with the latest developments in the field of large  
553 language models. The result will be a more dy-  
554 namic, scalable, and efficient system capable of  
555 keeping pace with advancements in AI technology.

### 556 From Situational QA to Realistic Simulations

557 MOTIVEBENCH fundamentally employs the "situa-  
558 tional question-answering" paradigm, where LLMs  
559 are prompted to answer questions about the next  
560 immediate step in various scenarios. However, this  
561 paradigm still deviates from real-world human so-  
562 cial activities, where individuals take a sequence of  
563 actions, form longer life stories, and behave sponta-  
564 neously without being asked, "what will you do?".

565 To overcome this limitation, we can consider us-  
566 ing the paradigm of LLMs performing role-playing  
567 in a simulation sandbox (Wang et al., 2024b).  
568 Given an initial scenario, LLMs act as specific  
569 characters with preset personas to engage in daily  
570 activities or achieve goal-oriented tasks. By analyz-  
571 ing the behavior trajectories elicited in the sandbox,  
572 we can assess LLMs' motivational reasoning and  
573 proactive action-taking capabilities in a more com-  
574 prehensive manner.

## References

- 576 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed 576  
577 Awadallah, Ammar Ahmad Awan, Nguyen Bach, 577  
578 Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat 578  
579 Behl, et al. 2024. Phi-3 technical report: A highly ca- 579  
580 pable language model locally on your phone. *arXiv 580*  
581 *preprint arXiv:2404.14219*. 581
- 582 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama 582  
583 Ahmad, Ilge Akkaya, Florencia Leoni Aleman, 583  
584 Diogo Almeida, Janko Altschmidt, Sam Altman, 584  
585 Shyamal Anadkat, et al. 2023. Gpt-4 technical report. 585  
586 *arXiv preprint arXiv:2303.08774*. 586
- 587 Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. 587  
588 2023. Using large language models to simulate mul- 588  
589 tiple humans and replicate human subject studies. 589  
590 In *International Conference on Machine Learning*, 590  
591 pages 337–371. PMLR. 591
- 592 Lisa P Argyle, Ethan C Busby, Nancy Fulda, Joshua R 592  
593 Gubler, Christopher Rytting, and David Wingate. 593  
594 2023. Out of one, many: Using language mod- 594  
595 els to simulate human samples. *Political Analysis*, 595  
596 31(3):337–351. 596
- 597 Layla El Asri, Jing He, and Kaheer Suleman. 2016. 597  
598 A sequence-to-sequence model for user simula- 598  
599 tion in spoken dialogue systems. *arXiv preprint 599*  
600 *arXiv:1607.00070*. 600
- 601 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, 601  
602 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei 602  
603 Huang, et al. 2023. Qwen technical report. *arXiv 603*  
604 *preprint arXiv:2309.16609*. 604
- 605 Albert Bandura. 2001. Social cognitive theory: An 605  
606 agentic perspective. *Annual review of psychology*, 606  
607 52(1):1–26. 607
- 608 Joyce Berg, John Dickhaut, and Kevin McCabe. 1995. 608  
609 Trust, reciprocity, and social history. *Games and 609*  
610 *economic behavior*, 10(1):122–142. 610
- 611 Daniil A Boiko, Robert MacKnight, and Gabe Gomes. 611  
612 2023. Emergent autonomous scientific research ca- 612  
613 pabilities of large language models. *arXiv preprint 613*  
614 *arXiv:2304.05332*. 614
- 615 Barbara Brehm. 2014. *Psychology of health and fitness*. 615  
616 FA Davis. 616
- 617 Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen, 617  
618 Guanqun Bi, Gongyao Jiang, Yaru Cao, Mengting 618  
619 Hu, Yunghwei Lai, Zexuan Xiong, et al. 2024. 619  
620 Tombench: Benchmarking theory of mind in large 620  
621 language models. *arXiv preprint arXiv:2402.15052*. 621
- 622 Israel Cohen, Yiteng Huang, Jingdong Chen, Jacob Ben- 622  
623 esty, Jacob Benesty, Jingdong Chen, Yiteng Huang, 623  
624 and Israel Cohen. 2009. Pearson correlation coeffi- 624  
625 cient. *Noise reduction in speech processing*, pages 625  
626 1–4. 626

627	Edward L Deci and Richard M Ryan. 2012. Self-determination theory. <i>Handbook of theories of social psychology</i> , 1(20):416–436.	680
628		681
629		682
630	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	683
631		684
632		685
633		686
634		687
635	Alice H Eagly and Wendy Wood. 2012. Social role theory. <i>Handbook of theories of social psychology</i> , 2:458–476.	688
636		689
637		690
638	Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. Scaling synthetic data creation with 1,000,000,000 personas. <i>arXiv preprint arXiv:2406.20094</i> .	691
639		692
640		693
641		694
642	Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, et al. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. <i>arXiv preprint arXiv:2406.12793</i> .	695
643		696
644		697
645		698
646		699
647	Sandra Graham. 1991. A review of attribution theory in achievement contexts. <i>Educational Psychology Review</i> , 3:5–39.	700
648		701
649		702
650	Martin Hagger and Nikos Chatzisarantis. 2005. <i>The social psychology of exercise and sport</i> . McGraw-Hill Education (UK).	703
651		704
652		705
653	Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. <i>arXiv preprint arXiv:2308.00352</i> .	706
654		707
655		708
656		709
657		710
658	Yupeng Hou, Jiacheng Li, Zhankui He, An Yan, Xiushi Chen, and Julian McAuley. 2024. Bridging language and items for retrieval and recommendation. <i>arXiv preprint arXiv:2403.03952</i> .	711
659		712
660		713
661		714
662	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. <i>ACM Transactions on Information Systems</i> .	715
663		716
664		717
665		718
666		719
667		720
668	Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. 2024. Qwen2.5-coder technical report. <i>arXiv preprint arXiv:2409.12186</i> .	721
669		722
670		723
671		724
672	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. <i>arXiv preprint arXiv:2410.21276</i> .	725
673		726
674		727
675		728
676		729
677	Henry J Bindel Jr. 1991. Becoming an effective classroom manager: A resource for teachers, by bob f. steere.	730
678		731
679		732
	Yeonghun Kang and Jihan Kim. 2023. Chatmof: An autonomous ai system for predicting and generating metal-organic frameworks. <i>arXiv preprint arXiv:2308.01423</i> .	733
		734
	Alan E Kazdin, American Psychological Association, et al. 2000. <i>Encyclopedia of psychology</i> , volume 8. American Psychological Association Washington, DC.	735
		736
	Simon Keizer, Milica Gasic, Filip Jurcicek, François Mairesse, Blaise Thomson, Kai Yu, and Steve Young. 2010. Parameter estimation for agenda-based user simulation. In <i>Proceedings of the SIGDIAL 2010 Conference</i> , pages 116–123.	737
		738
	Florian Kreyszig, Iñigo Casanueva, Pawel Budzianowski, and Milica Gasic. 2018. Neural user simulation for corpus-based policy optimisation for spoken dialogue systems. <i>arXiv preprint arXiv:1805.06966</i> .	739
		740
	Matthew Le, Y-Lan Boureau, and Maximilian Nickel. 2019. Revisiting the evaluation of theory of mind through question answering. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 5872–5877.	741
		742
	Sangmin Lee, Minzhi Li, Bolin Lai, Wenqi Jia, Fiona Ryan, Xu Cao, Ozgur Kara, Bikram Boote, Weiyan Shi, Diyi Yang, et al. 2024. Towards social ai: A survey on understanding social interactions. <i>arXiv preprint arXiv:2409.15316</i> .	743
		744
	Bethanie Maples, Merve Cerit, Aditya Vishwanath, and Roy Pea. 2024. Loneliness and suicide mitigation for students using gpt3-enabled chatbots. <i>npj mental health research</i> , 3(1):4.	745
		746
	AH Maslow. 1943. A theory of human motivation. <i>Psychological Review google schola</i> , 2:21–28.	747
		748
	Nikhil Mehta, Milagro Teruel, Patricio Figueroa Sanz, Xin Deng, Ahmed Hassan Awadallah, and Julia Kiseleva. 2023. Improving grounded language understanding in a collaborative environment by interacting with agents through help feedback. <i>arXiv preprint arXiv:2304.10750</i> .	749
		750
	AI Meta. 2024. Introducing llama 3.1: Our most capable models to date. <i>Meta AI Blog</i> .	751
		752
	Rémi Radel, Dusan Pjavec, Karen Davranche, Fabienne d’Arripe Longueville, Serge S Colson, Thomas Lapole, and Mathieu Gruet. 2016. Does intrinsic motivation enhance motor cortex excitability? <i>Psychophysiology</i> , 53(11):1732–1738.	753
		754
	Hannah Rashkin, Antoine Bosselut, Maarten Sap, Kevin Knight, and Yejin Choi. 2018a. Modeling naive psychology of characters in simple commonsense stories. <i>arXiv preprint arXiv:1805.06533</i> .	755
		756



841	Alex Young, Bei Chen, Chao Li, Chengen Huang,	from others). Satisfying these needs boosts	890
842	Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng	an individual’s sense of self-worth and social	891
843	Zhu, Jianqun Chen, Jing Chang, et al. 2024. Yi:	standing.	892
844	Open foundation models by 01. ai. <i>arXiv preprint</i>		
845	<i>arXiv:2403.04652</i> .		
846	Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and	• <b>Self-Actualization Needs:</b> The highest level	893
847	Minlie Huang. 2023. Large language models are not	of need, self-actualization refers to the indi-	894
848	robust multiple choice selectors. In <i>The Twelfth Inter-</i>	vidual’s pursuit of fulfilling their potential and	895
849	<i>national Conference on Learning Representations</i> .	achieving their ideal self. Self-actualization	896
850	Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang,	is typically manifested through creativity, per-	897
851	Haofei Yu, Zhengyang Qi, Louis-Philippe Morency,	sonal growth, and realizing one’s own value.	898
852	Yonatan Bisk, Daniel Fried, Graham Neubig, et al.		
853	2023. Sotopia: Interactive evaluation for social	Furthermore, the Reiss Motivation Profile	899
854	intelligence in language agents. <i>arXiv preprint</i>	(RMP) (Reiss, 2004), proposed by psychologist	900
855	<i>arXiv:2310.11667</i> .	Steven Reiss, is also a theoretical model that aims	901
856	<b>A Details of the Hierarchy of Needs</b>	to understand an individual’s motivation by an-	902
857	Maslow’s hierarchy of needs (Maslow, 1943) is a	alyzing their preferences across 16 fundamental	903
858	motivational theory proposed by American psychol-	needs. Each person responds differently to these	904
859	ogist Abraham Maslow in 1943, which explains the	16 needs, and this variation determines their be-	905
860	prioritization and fulfillment of human needs. He	havior, decision-making, and lifestyle. The core	906
861	categorized human needs into five levels, ranging	assumption of the RMP model is that individual be-	907
862	from basic survival needs to higher psychological	haviors are driven by these needs, and the intensity	908
863	needs, ultimately culminating in self-actualization.	of these needs shapes one’s behavioral patterns.	909
864	Below is a detailed description of each level of	Needs like <b>Curiosity</b> , <b>Idealism</b> , and <b>Honor</b>	910
865	need, along with corresponding example scenarios	drive individuals to seek knowledge, moral val-	911
866	from our MOTIVEBENCH as shown in Figure 8.	ues, and social recognition, reflecting a pursuit of	912
867		cognitive and social fulfillment. <b>Independence</b>	913
868	• <b>Physiological Needs:</b> Physiological needs are	and <b>Power</b> highlight desires for autonomy and con-	914
869	the most fundamental survival needs for hu-	trol, with those seeking independence preferring	915
870	mans, including air, water, food, warmth, and	autonomy, while those seeking power focus on in-	916
871	sleep. These needs form the base of the hier-	fluencing others. Once cognitive and control needs	917
872	archy, and only once they are met can individ-	are met, emotional needs like <b>Acceptance</b> , <b>Family</b> ,	918
873	uals engage in other activities.	and <b>Romance</b> emphasize social connections, while	919
874		<b>Social Contact</b> and <b>Status</b> stress the importance	920
875	• <b>Safety Needs:</b> Once physiological needs are	of recognition and belonging. <b>Order</b> and <b>Saving</b>	921
876	met, humans seek security, which includes	reflect desires for stability, and <b>Tranquility</b> and	922
877	physical safety, financial stability, health, and	<b>Physical Activity</b> focus on inner peace and health.	923
878	environmental consistency. This need reflects	Finally, <b>Vengeance</b> and <b>Eating</b> address responses	924
879	an individual’s desire for order, protection,	to injustice and the enjoyment of life.	925
880	and predictability in the future.		
881		<b>B Question Format Example</b>	926
882	• <b>Love and Belonging Needs:</b> Also known as	Figure 9 shows an example of three types of tasks	927
883	social needs, people seek to build interper-	in the same scenario, from Persona-Hub.	928
884	sonal relationships, gain friendships, experi-		
885	ence love, and feel a sense of belonging to a	<b>C Evaluated LLMs</b>	929
886	social group. This need involves the desire to	We evaluate 29 popular LLMs across a range of	930
887	integrate into society, be accepted, and inter-	parameter sizes, including several models from the	931
888	act with others.	GPT series (Hurst et al., 2024) (GPT-4o 2024-05-	932
889		13, GPT-4o mini 2024-07-18, and GPT-3.5-Turbo	933
	• <b>Esteem Needs:</b> Esteem needs are divided	1106), the LLaMA series (Meta, 2024; Touvron	934
	into intrinsic esteem (self-confidence, self-	et al., 2023) (LLaMA 3.1 and LLaMA 2), the Qwen	935
	respect, and independence) and extrinsic es-	series (Hui et al., 2024; Yang et al., 2024; Bai et al.,	936
	teem (recognition, appreciation, and status		

<p><b>Level 5: Self-Actualization Needs</b></p>	<p><b>Example Scenario:</b> Samantha, a graduate student majoring in analytical psychology, is deeply engaged in her thesis research. Recently, she attended an international conference and participated in a challenging debate with Dr. Martin, a respected philosopher, on topics related to free will and determinism. Afterward, she continued to reflect on the discussion, thinking about its potential connections to her academic work. Additionally, she had received constructive feedback from her advisor on ways to develop her research further.</p>
<p><b>Level 4: Esteem Needs</b></p>	<p><b>Example Scenario:</b> Suresh, an avid local football fan from Gujarat, operates a well-known blog focused on local sports events. He has supported his favorite team for many years and frequently uses his platform to highlight their achievements. Recently, he was present at a critical match between two rival teams, during which his favorite team lost due to a disputed penalty call. In the aftermath of the match, despite facing backlash and threats from upset fans of the rival team, Suresh dedicated several hours to writing a detailed post that critiqued the referee's decision and its implications for the game.</p>
<p><b>Level 3: Love &amp; Belonging Needs</b></p>	<p><b>Example Scenario:</b> As Harry dashed through the downpour, he stumbled upon a quaint little pub that seemed like an oasis amidst the stormy evening, its warm glow inviting him to take shelter. Recently moved to a new city, Harry felt increasingly isolated. Inside the pub, he was drawn to the lively atmosphere and began trying different craft beers. Nervous about attending a wedding alone, he found a friendly bartender and chatted about shared interests like classic rock music. After an hour of engaging conversation, Harry invited the bartender to join him at the wedding.</p>
<p><b>Level 2: Safety Needs</b></p>	<p><b>Example Scenario:</b> Owen, an IT specialist, is meticulous and focuses on creating a secure and hazard-free home office environment, taking great care in managing every aspect of his workspace to avoid potential risks. Recently, he invested in a range of safety measures, including fireproof storage for sensitive documents and high-quality surge protection for his electronics. He often discusses office safety practices with his peers and follows several technology newsletters.</p>
<p><b>Level 1: Physiological Needs</b></p>	<p><b>Example Scenario:</b> Yara, a college student balancing two part-time jobs and maintaining a demanding academic schedule, browses online reviews one evening in search of an affordable brunch spot. She comes across a quaint diner known for its sizable breakfast platters and quick service. Yara decides to go for brunch the next morning. She arrives at the diner, orders a large breakfast platter, and devours her meal efficiently. After the meal, she shows notable satisfaction.</p>

Figure 8: Example scenarios for Maslow's hierarchy of needs in MOTIVEBENCH.

**An example from Persona-Hub**

**1. Motivational Reasoning Question**  
**Context:** Dr. Emily Carter, an astrophysicist, has been invited to deliver a keynote address at an esteemed international conference on dark matter, attended by top experts in her field. The conference is known for its rigorous standards and high expectations for speakers. In the days before her speech, she spends long evenings preparing her presentation, often sacrificing rest and personal time. Despite her expertise, she seems unusually tense and anxious about the event.  
**Question:** What is the most likely motivation behind her meticulous preparation for the speech?  
**Options:**  
A. She wants to impress her colleagues and friends.  
B. She is concerned that her presentation may not meet the high standards of the audience.  
C. She hopes her presentation will lead to future research funding.  
D. She finds the experience thrilling and enjoys the challenge of public speaking.  
E. She wishes to build professional connections with other scientists.  
F. She feels driven to share her work to advance scientific knowledge.  
**Correct Answer: B**

**2. Behavioral Reasoning Question**  
**Context:** Dr. Emily Carter is a respected astrophysicist with a strong commitment to her work and her family. She has received an offer for a prestigious position as a visiting professor at a top university overseas, which would enhance her career and provide financial benefits. Despite her professional success, she is very close to her family, especially her younger brother, who requires her support due to a disability. In addition to her work, she volunteers in her community to help families in similar situations.  
**Question:** Faced with this opportunity, driven by her desire to balance her career and family responsibilities, what behavior is Dr. Carter most likely to exhibit in response to this offer?  
**Options:**  
A. Accept the position immediately, viewing it as a unique chance to advance her career without hesitation.  
B. Reject the offer outright, choosing to dedicate herself fully to her family and the needs of her brother.  
C. Propose a modified working arrangement that allows her to fulfill both her professional duties and family obligations.  
D. Take a temporary leave from her current job to focus on supporting her family during this transition.  
E. Agree to the position but plan to take frequent extended breaks to care for her brother.  
F. Delay her decision for several months, considering how this change might affect her personal life.  
**Correct Answer: C**

**3. Motive&Behavior Reasoning Question**  
**Context:** Dr. Emily Carter, a renowned astrophysicist, has just returned from a year-long sabbatical where she conducted groundbreaking research on dark matter. Upon returning, she faces a significant backlog of urgent emails and requests from colleagues and students, including a tight deadline for a grant application and a critical experiment that requires her oversight. Furthermore, her younger brother, diagnosed with a chronic illness during her absence, needs ongoing support.  
**Question:** Given this scenario, what kind of behavior is Dr. Carter most likely to exhibit next, and what might motivate that behavior?  
**Options:**  
A. Prioritize her research to recover lost ground, motivated by the desire to reclaim her status in the academic community.  
B. Request an extended leave to fully support her brother, motivated by her commitment to family obligations.  
C. Strategically delegate tasks to her team to manage her workload, motivated by the need to harmonize her professional duties with personal responsibilities.  
D. Concentrate exclusively on her career, neglecting her brother's needs, motivated by an unyielding ambition for professional achievement.  
E. Organize meetings to share her research insights and delegate responsibilities, motivated by a desire to reintegrate into her professional environment effectively.  
F. Work excessive hours to complete all tasks independently, motivated by a strong sense of personal obligation.  
**Correct Answer: C**

Figure 9: An example question from Persona-Hub.

2023) (Qwen 2.5, Qwen 2, and Qwen), the Phi series (Abdin et al., 2024) (Phi 3.5 and Phi 3), the GLM series (GLM et al., 2024) (ChatGLM 3 and GLM 4), as well as other models like Baichuan 2 (Yang et al., 2023) and Yi 1.5 (Young et al., 2024). These models span a wide spectrum of architectures and parameter configurations, offering a comprehensive evaluation of current LLM performance across various tasks and benchmarks. All of our experiments are conducted on a machine with four A100 80GB GPUs.

**D Detailed Results of three tasks** 948

Tables 4 and 5 present the performance of different models on motivation reasoning, behavior reasoning, and motivation-behavior reasoning tasks. We observe that larger models tend to perform better than smaller models. 949 950 951 952 953

**E Experiment Prompts** 954

Table 6 presents the detailed prompts used for model evaluation in MOTIVEBENCH, including the base prompt and the CoT prompt. Tables 7, 8, and 9 present the prompts for the questioner, reviewer, and modifier in our multi-agent question generation and modification framework. Specifically, we empirically set the maximum number of modification rounds to 5. If no issues are identified within this threshold across the aspects managed by the three reviewers, the modification process is stopped. 955 956 957 958 959 960 961 962 963 964

**F Case Study** 965

In this section, we analyze the differences between GPT-4o and human reasoning in motivational and behavioral tasks through detailed case studies. Specifically, for questions where GPT-4o selects incorrect answers, we prompt the model to explain its reasoning behind the chosen option and analyze why the correct answer is appropriate. By comparing GPT-4o’s reasoning process with that of humans, we uncover key gaps in understanding. Tables 10, 11, 12, 13, 14 provide illustrative examples from different perspectives, showcasing GPT-4o’s limitations in accurately capturing human-like reasoning, particularly in interpreting motivations, emotional nuances, and contextual behaviors. These cases highlight the discrepancies between GPT-4o’s approach and the more contextually sensitive reasoning exhibited by humans. 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982

MRQ: Motivational Reasoning Question  
 BRQ: Behavioral Reasoning Question  
 MBQ: Motive&Behavior Reasoning Question

Base Method	Vitual Profiles Persona-Hub			Real-life Scenarios					
	MRQ	BRQ	MBQ	Amazon			Yelp		
LLMs	MRQ	BRQ	MBQ	MRQ	BRQ	MBQ	MRQ	BRQ	MBQ
Baichuan2-7B-Chat	<b>81.83</b>	68.83	57.67	69.67	68.33	74.67	74.00	61.67	68.00
Baichuan2-13B-Chat	80.83	<b>77.17</b>	<b>67.33</b>	<b>77.67</b>	<b>83.00</b>	<b>78.33</b>	<b>80.33</b>	<b>75.33</b>	<b>79.00</b>
ChatGLM3-6B	84.67	70.67	63.83	72.33	75.00	73.67	76.00	69.67	74.00
GLM4-9B-Chat	<b>93.67</b>	<b>82.17</b>	<b>77.00</b>	<b>92.67</b>	<b>89.67</b>	<b>89.00</b>	<b>93.33</b>	<b>82.67</b>	<b>84.67</b>
Yi1.5-6B-Chat	83.33	71.50	56.33	84.00	83.67	69.67	83.00	73.67	71.33
Yi1.5-9B-Chat	92.17	<b>84.17</b>	77.67	<b>91.33</b>	85.00	86.67	88.33	86.00	84.67
Yi1.5-34b-Chat	<b>96.33</b>	82.33	<b>81.00</b>	90.67	<b>89.33</b>	<b>91.33</b>	<b>95.33</b>	<b>91.00</b>	<b>91.00</b>
Phi3-mini-4k-Instruct	89.83	81.83	78.50	85.33	88.33	89.00	86.33	80.00	87.00
Phi3-small-8k-Instruct	94.50	84.33	80.33	92.33	90.00	92.00	90.33	87.67	89.00
Phi3-medium-4k-Instruct	96.33	<b>89.67</b>	<b>84.50</b>	91.33	<b>94.00</b>	<b>93.67</b>	<b>97.00</b>	<b>91.67</b>	90.67
Phi3.5-mini-Instruct	90.67	82.83	79.83	88.67	89.00	87.33	89.00	84.00	88.00
Phi3.5-MoE-Instruct	<b>96.50</b>	88.17	82.67	<b>93.00</b>	87.00	89.33	93.33	86.00	<b>91.00</b>
Llama2-7B-Chat	66.83	47.67	48.00	61.67	47.67	64.33	62.00	45.33	51.67
Llama2-13B-Chat	84.50	75.00	70.17	80.00	78.67	80.00	81.00	75.00	75.00
Llama2-70B-Chat	92.17	79.33	69.33	88.33	84.00	79.33	86.67	79.33	76.33
Llama3.1-8B-Instruct	95.50	<b>94.17</b>	77.00	<b>91.33</b>	90.33	88.00	89.33	81.33	87.33
Llama3.1-70B-Instruct	<b>98.67</b>	92.00	<b>90.17</b>	90.67	<b>94.67</b>	<b>95.67</b>	<b>94.67</b>	<b>93.67</b>	<b>93.33</b>
Qwen-7B-Chat	82.17	78.17	66.33	78.00	82.67	79.67	79.00	72.00	78.33
Qwen-14B-Chat	92.00	86.50	78.17	89.33	88.67	86.00	90.33	82.33	87.00
Qwen-72B-Chat	95.50	87.17	84.17	93.67	94.00	92.33	91.67	87.67	88.33
Qwen2-7B-Instruct	92.17	88.00	84.00	89.00	92.00	89.33	93.33	87.33	91.67
Qwen2-72B-Instruct	97.50	90.33	88.17	93.33	93.33	94.33	93.33	90.00	92.67
Qwen2.5-7B-Instruct	93.00	85.50	86.67	90.67	88.00	90.67	84.00	85.67	91.67
Qwen2.5-14B-Instruct	98.33	91.33	87.00	92.33	92.67	93.00	<b>95.00</b>	93.00	<b>94.00</b>
Qwen2.5-32B-Instruct	98.83	91.67	90.00	<b>97.33</b>	<b>94.33</b>	<b>96.67</b>	<b>95.00</b>	<b>95.33</b>	92.67
Qwen2.5-72B-Instruct	<b>99.00</b>	<b>92.17</b>	<b>90.17</b>	94.67	94.00	96.33	94.00	90.33	<b>94.00</b>
GPT-3.5-Turbo 1106	93.33	86.00	77.50	93.00	92.33	91.00	93.33	82.67	84.00
GPT-4o mini 2024-07-18	97.50	91.83	90.00	94.33	93.67	95.67	95.33	91.33	91.33
GPT-4o 2024-05-13	<b>98.83</b>	<b>94.50</b>	<b>90.83</b>	<b>95.00</b>	<b>97.00</b>	<b>97.33</b>	<b>96.33</b>	<b>93.67</b>	<b>93.67</b>

Table 4: The experimental results of the vanilla prompt-based method on the three types of tasks.

MRQ: Motivational Reasoning Question  
 BRQ: Behavioral Reasoning Question  
 MBQ: Motive&Behavior Reasoning Question

CoT Method	Vitual Profiles Persona-Hub			Real-life Scenarios					
	MRQ	BRQ	MBQ	Amazon			Yelp		
LLMs	MRQ	BRQ	MBQ	MRQ	BRQ	MBQ	MRQ	BRQ	MBQ
Baichuan2-7B-Chat	<b>80.17</b>	64.50	59.50	<b>68.00</b>	64.67	73.33	<b>68.33</b>	58.67	62.00
Baichuan2-13B-Chat	61.00	<b>75.50</b>	<b>64.00</b>	59.67	<b>71.33</b>	<b>74.33</b>	61.33	<b>66.67</b>	<b>69.00</b>
ChatGLM3-6B	75.33	66.50	59.17	66.00	68.00	68.33	67.33	67.00	71.67
GLM4-9B-Chat	<b>92.17</b>	<b>82.33</b>	<b>77.17</b>	<b>92.00</b>	<b>91.00</b>	<b>89.67</b>	<b>92.33</b>	<b>81.67</b>	<b>89.00</b>
Yi1.5-6B-Chat	86.83	75.00	71.83	80.00	83.33	83.67	81.67	77.67	82.00
Yi1.5-9B-Chat	89.83	82.33	72.67	87.33	85.00	<b>86.00</b>	86.00	86.00	83.67
Yi1.5-34b-Chat	<b>93.50</b>	<b>83.50</b>	<b>80.33</b>	<b>89.33</b>	<b>91.67</b>	<b>86.00</b>	<b>88.33</b>	<b>88.33</b>	<b>87.33</b>
Phi3-mini-4k-Instruct	84.83	76.67	68.17	82.33	74.33	79.00	82.00	73.00	81.00
Phi3-small-8k-Instruct	92.00	82.33	<b>80.33</b>	85.67	74.67	87.00	87.00	79.00	88.33
Phi3-medium-4k-Instruct	<b>94.50</b>	<b>86.33</b>	80.17	<b>92.00</b>	86.00	87.33	94.33	85.33	89.67
Phi3.5-mini-Instruct	86.17	81.67	79.17	84.33	87.67	<b>88.00</b>	86.67	81.67	84.33
Phi3.5-MoE-Instruct	93.00	83.33	71.33	<b>92.00</b>	<b>90.67</b>	82.33	<b>96.00</b>	<b>86.67</b>	<b>91.33</b>
Llama2-7B-Chat	67.50	57.83	46.17	61.00	60.00	62.67	66.00	58.33	62.67
Llama2-13B-Chat	81.33	70.67	59.17	70.33	75.67	71.67	78.67	68.00	72.67
Llama2-70B-Chat	89.67	85.50	77.67	85.33	86.67	83.67	89.33	82.67	85.33
Llama3.1-8B-Instruct	90.33	81.33	73.17	88.00	88.00	86.33	88.67	83.67	85.33
Llama3.1-70B-Instruct	<b>98.17</b>	<b>90.33</b>	<b>84.33</b>	<b>92.00</b>	<b>94.67</b>	<b>94.00</b>	<b>94.67</b>	<b>88.33</b>	<b>91.00</b>
Qwen-7B-Chat	81.17	73.83	65.83	78.33	79.33	79.33	81.33	68.67	76.67
Qwen-14B-Chat	88.17	82.50	75.83	89.67	85.33	82.67	89.33	80.67	81.33
Qwen-72B-Chat	93.83	86.33	84.00	91.00	91.67	92.33	89.33	88.00	87.33
Qwen2-7B-Instruct	91.33	86.67	86.33	90.33	92.67	88.67	93.00	86.33	93.00
Qwen2-72B-Instruct	96.50	88.67	86.50	91.33	92.67	91.67	93.67	90.33	87.00
Qwen2.5-7B-Instruct	92.33	86.50	84.50	88.67	92.67	87.00	86.67	84.33	90.00
Qwen2.5-14B-Instruct	96.17	88.00	84.67	90.67	93.33	93.67	93.33	91.33	90.67
Qwen2.5-32B-Instruct	96.00	<b>90.33</b>	86.17	<b>95.33</b>	94.00	<b>95.67</b>	94.67	<b>92.33</b>	90.67
Qwen2.5-72B-Instruct	<b>98.00</b>	88.83	<b>87.50</b>	94.00	<b>94.33</b>	94.33	<b>95.33</b>	91.67	<b>93.67</b>
GPT-3.5-Turbo 1106	92.50	84.17	81.00	93.00	93.00	89.67	93.00	83.33	87.33
GPT-4o mini 2024-07-18	<b>96.67</b>	88.67	86.33	<b>93.33</b>	95.00	94.67	94.33	91.00	<b>93.00</b>
GPT-4o 2024-05-13	96.17	<b>92.17</b>	<b>88.00</b>	92.33	<b>95.67</b>	<b>96.33</b>	<b>96.00</b>	<b>93.00</b>	92.67

Table 5: The experimental results of the CoT prompt-based method on the three types of tasks.



<b>Base Prompt for Evaluation</b>
<p>The following is a {Question_Type}. You should: {Type_Interpretation}. Carefully read the given question, fully immerse yourself in the role of the character described, and reason based on the information provided. Your answer should rely strictly on the given details.</p> <p>Note:</p> <ol style="list-style-type: none"> <li>1. Based on the content of the given question, please infer the most likely answer.</li> <li>2. You must select one answer from the given options: "A, B, C, D, E, F" as the most likely choice. Even if the question does not provide sufficient information to determine the correct answer, you should randomly choose one option as your output.</li> <li>3. The result can only return <b>**one character without any other explanation**</b>.</li> </ol> <p>Question: {Question_Content}</p> <p>Options: {Options}</p>
<b>CoT Prompt for Evaluation</b>
<p>The following is a {Question_Type}. You should: {Type_Interpretation}. Carefully read the given question, fully immerse yourself in the role of the character described, and reason based on the information provided. Your answer should rely strictly on the given details.</p> <p>Note:</p> <ol style="list-style-type: none"> <li>1. Based on the content of the given question, please think step by step and infer the most likely answer.</li> <li>2. You must select one answer from the given options: "A, B, C, D, E, F" as the most likely choice. Even if the question does not provide sufficient information to determine the correct answer, you should randomly choose one option as your output.</li> <li>3. Please first think through the question step by step, analyze the reasoning process for the possible answers, and finally output the most likely answer's letter. <b>**The last line of your reply should only contain one character of your final choice.**</b></li> </ol> <p>Question: {Question_Content}</p> <p>Options: {Options}</p>
<b>Illustration</b>
<p>{Question_Type}: {Type_Interpretation}</p> <ol style="list-style-type: none"> <li>1. <b>Motivational Reasoning Question:</b> Based on the given scenario and the character's profile, determine the most likely motivation behind the character's behavior.</li> <li>2. <b>Behavioral Reasoning Question:</b> Based on the given scenario and the character's profile, determine the most likely behavior the character would take next, given the motivation.</li> <li>3. <b>Motive&amp;Behavior Reasoning Question:</b> Based on the given scenario and the character's profile, determine the most likely motivation the character would develop next and the corresponding behavior would take.</li> </ol> <p>{Options}</p> <p>To minimize the potential bias caused by the order of options, we will randomize the order of options six times and calculate the average result from these six experiments. Specifically, the order sequences used will be: [1, 2, 3, 4, 5, 6], [6, 5, 4, 3, 2, 1], [3, 1, 6, 5, 4, 2], [2, 3, 5, 6, 1, 4], [5, 4, 1, 2, 6, 3], and [4, 6, 2, 1, 3, 5], ensuring that each option appears in every possible position across the six sequences.</p>

Table 6: Prompts for evaluation.

### Prompts of Questioner

Consider the four elements of scenario, profile, motivation, and behavior. In a given scenario, a character with a specific profile will perform a certain behavior based on a certain motivation. You are a professional psychologist and sociologist, skilled at creating challenging reasoning questions based on given scenarios to test participants' motivation and behavior reasoning abilities.

**\*\*Please create three questions based on the given scenario:\*\***

1. **Motivational Reasoning Question:** Given a complex scenario, a specific profile, and a given behavior, infer the most likely motivation behind the character's behavior. The question should not contain any direct description related to the predicted motivation.
2. **Behavioral Reasoning Question:** Given a complex scenario, a specific profile, and a given motivation, infer the most likely behavior the character will perform based on that motivation. The question should not contain any direct description related to the predicted behavior.
3. **Motive&Behavior Reasoning Question:** This is a more advanced test. The question should only include the complex scenario and the character's profile. Using only the complex scenario and specific profile, infer the most likely motivation the character will have and the corresponding behavior they will perform.

To summarize, all three questions are based on the same story scenario and character profile setup. The motivation reasoning question requires the addition of a behavior in the question stem and asks the participant to infer the motivation for that behavior. The behavior reasoning question requires the addition of a motivation in the question stem and asks the participant to infer the behavior that may result from that motivation. The motivation and behavior reasoning question does not need any additional information and requires the participant to infer both the motivation and behavior of the character based on the given scenario and profile.

**\*\*Note:\*\***

1. You will be provided with a simple scenario description. Please rewrite this scenario by correcting any logical inconsistencies, and add relevant details to make the scenario, profile, motivation, and behavior more vivid and complex.
2. Choose the most appropriate motivation and behavior to create the questions. However, ensure that the motivation and behavior are only related to real human needs, not to any POIs or products in the text.
3. The three questions are independent of each other and should be answered separately, meaning that each question should only rely on its own stem and not contain any information from the others. Therefore, please ensure that each question has enough rich and complex scenario and profile information to support correct reasoning.
4. Each question should have only one correct answer, along with five distractors. The distractors must be related to certain parts of the information in the question. Please analyze why each option is correct or incorrect.
5. The question stem must include irrelevant or redundant information that creates distractions and challenges. This is necessary to ensure each question is challenging. The correct answer must not be explicitly stated in the question.

Table 7: Detailed prompts of questioner in the multi-agent framework.

### Prompts of Reviewers

You are a strict and discerning psychologist and sociologist, capable of precisely identifying issues in the given behavior and motivation reasoning questions and offering improvement suggestions.

I will provide you with three behavior and motivation reasoning questions. Please evaluate them based on the following aspects:

1. **Reasonableness of the Question Information and Type**: Specifically, all three questions should contain a concrete scenario and character profile. The motivation reasoning question should include additional behavioral information about the character. The behavior reasoning question should include additional motivational information about the character. The motivation and behavior reasoning question should not contain any direct clues about the motivation or behavior.
2. **Logical Consistency and Reasonableness of the Four-Tuple**: Assess whether, in the given scenario, a character with a specific profile would logically perform the stated behavior based on the provided motivation.
3. **Sufficiency of Information to Derive the Correct Answer**: Examine whether the information provided in each question is enough to infer the correct answer. If not, suggest modifications to the scenario or character profile to make the information clearer or more comprehensive.
4. **Challenge and Difficulty of the Question**: Evaluate whether the question presents an appropriate level of difficulty and challenge for the respondent.
5. **Correct Answer Must Not Be Explicitly Stated**: Ensure that the correct answer does not appear explicitly in the question information and can only be deduced through reasoning steps.
6. **Clarity and Plausibility of Distractor Options**: Evaluate whether the incorrect options are misleading and whether they correspond to distracting information within the question. If they do not, suggest adding the relevant distracting information or modifying the options.
7. **Adequate Distractors and Redundant Information**: Ensure that each question includes enough irrelevant or redundant information to make the question challenging, but without disrupting the logic needed to deduce the correct answer.
8. **Objectivity and Neutrality of the Question**: Ensure that the question is presented in a neutral and objective manner, with no implicit suggestion of the correct answer.

Please provide specific modification suggestions for the question set and give your feedback to the question author in a reasonable tone. Summarize your evaluation into a single paragraph of suggestions.

(All the aspects listed above are of concern, and each reviewer will be asked to focus on different aspects.)

Table 8: Detailed prompts of reviewers in the multi-agent framework.

### Prompts of Modifier

Consider the four elements of scenario, character profile, motivation, and behavior. In the given scenario, a character with a specific profile will perform a certain behavior based on a particular motivation. You are a professional psychologist and sociologist, skilled in refining motivation and behavior reasoning test questions and providing relevant suggestions for improvement.

**\*\*The specific types of the three questions are as follows:\*\***

1. **Motivational Reasoning Question:** Based on a complex scenario, a specific character profile, and a given behavior, deduce the most likely motivation behind the character's action. The question should not include any description related to the predicted motivation.
2. **Behavioral Reasoning Question:** Based on a complex scenario, a specific character profile, and a given motivation, deduce the most likely behavior the character will perform based on that motivation. The question should not include any description related to the predicted behavior.
3. **Motive&Behavior Reasoning Question:** The question should only include a complex scenario and character profile. This is a more difficult question type, where the respondent must deduce the most likely behavior and corresponding motivation of the character based solely on the scenario and character profile.

In summary, all three questions are based on the same story scenario and character profile settings. For the motivation reasoning question, an additional behavior is given, and the task is to deduce the motivation behind that behavior. In the behavior reasoning question, an additional motivation is given, and the task is to deduce the behavior that would most likely result from that motivation. The motivation and behavior reasoning question, however, does not provide any additional information, requiring the respondent to deduce both the motivation and behavior from the scenario and character profile. It is crucial that the story scenario and character profile in the question are rich enough to support reasoning and lead to the correct answer.

**\*\*Specific Requirements:\*\***

1. Carefully consider each suggestion based on the given questions and selectively make reasonable changes to the questions.
2. Do not delete the distracting information related to the incorrect answers, as this is necessary to ensure the questions remain challenging.
3. The three questions are independent of each other and are to be answered separately. Respondents should only reason based on the question provided, without seeing any other information. Therefore, ensure that each question has sufficiently rich and complex scenario and character profile information.
4. After making revisions, analyze each option to determine why it is correct or incorrect. If there are any issues, modify the question again to ensure the uniqueness of the correct answer.

Table 9: Detailed prompts of modifier in the multi-agent framework.

### Case 1 - Motive&Behavior Reasoning Question

**Context:** Samantha, a grateful accident survivor, was involved in a severe car crash six months ago. She has since undergone multiple surgeries and intense physical therapy. As a freelance writer before the accident, Samantha now spends much of her free time reading and writing poetry, which she shares occasionally on her personal blog. She has been sharing her recovery journey on social media and feels strongly about using her experience to make a positive impact. Recently, she has gained a lot of traction and connected with many individuals through her posts, deepening her sense of responsibility to those who are still struggling.

**Question:** What kind of behavior is Samantha most likely to exhibit next, and what is the motivation behind it?

Options:

A. Motivation: Cultivating a personal brand; Behavior: Posting artistic photos of her daily life and updates about her writing process.

B. Motivation: Highlighting the importance of emotional resilience; Behavior: Hosting online webinars focused on mental health strategies.

C. Motivation: Encouraging community support; Behavior: Organizing small group meetups for accident survivors to share their experiences.

D. Motivation: Seeking validation from peers; Behavior: Posting emotionally charged poetry on social media to gain likes and shares.

E. Motivation: Gaining recognition for her journey; Behavior: Collaborating with influencers to promote her story.

F. Motivation: Finding solace through expression; Behavior: Writing a memoir to reflect on her healing process.

**Correct Answer: C**

**GPT-4o's Answer: B**

**Analysis:** In this scenario, the appropriate action should be to offer emotional support, not discuss mental health strategies. GPT-4o's choice (hosting a webinar on mental health strategies) focuses on professional methods, which doesn't align with Samantha's current situation. She seeks to inspire others through her personal experiences, not teach strategies. GPT-4o's reasoning is too theoretical, lacking the empathy and life experience humans use to understand motivations. Humans, considering Samantha's struggles, would focus on actions that resonate with her personal healing process, such as sharing her story to help others.

Table 10: Case 1 on a Motive&Behavior Reasoning Question.

### Case 2 - Motive&Behavior Reasoning Question

**Context:** James, a seasoned stockbroker specializing in tech and software stocks, has recently noticed a growing interest among younger investors in sustainable and socially responsible investments. Despite his initial skepticism, he recognizes the potential effects of this trend on his career. Additionally, he faces increasing competition from newer, digitally-savvy brokers capitalizing on this shift. Furthermore, he has come across various articles detailing the increasing demand for sustainable investments.

**Question:** Given these circumstances, what kind of behavior is James most likely to exhibit next, and what could be the motivation behind it?

**Options:**

- A. He will develop a marketing strategy aimed at promoting sustainable tech stocks, driven by the desire to connect with a younger audience interested in socially responsible investments.
- B. He will partner with a fintech firm specializing in sustainable investments, motivated by the necessity to broaden his service offerings and enhance client retention.
- C. He will begin writing articles for finance magazines, driven by the ambition to share his insights on the importance of sustainable investing among his peers.
- D. He will initiate a webinar series focusing on sustainable investment trends, motivated by the goal of showcasing his expertise and engaging with potential clients.
- E. He will host social events for potential investors, driven by the intention to foster relationships and promote discussions around sustainable investing.
- F. He will create a newsletter highlighting sustainable investment options, motivated by the aim of educating clients about emerging trends in the market.

**Correct Answer: A**

**GPT-4o's Answer: D**

**Analysis:** When faced with new market trends, humans typically prioritize directly addressing market demands and customer interests. For example, the growing interest of young investors in sustainable investments leads to a marketing strategy tailored to this group, which aligns with real market needs. GPT-4o tends to suggest that James might showcase his expertise through a webinar, but this motivation focuses more on "self-promotion" rather than directly responding to market demands or attracting a specific group, failing to address the competitive pressures and market changes.

Table 11: Case 2 on a Motive&Behavior Reasoning Question.

### Case 3 - Motive&Behavior Reasoning Question

**Context:** Yara, a new mother with a newborn baby girl who has a history of allergies, recently dined at a café that provided detailed ingredient lists and used allergen-safe cooking methods. She was satisfied with the café's attention to allergen management and its ability to cater to her needs. Yara is also known to actively participate in community groups focused on managing allergies in children.

**Question:** In this scenario, what is Yara most likely to do next, and what is her primary motivation?

**Options:**

- A. Sharing her experience with others, motivated by her commitment to helping the allergy community.
- B. Thanking the café staff, motivated by appreciation for their allergen-safe practices.
- C. Researching other restaurants, motivated by a desire for variety in dining options.
- D. Leaving a negative review elsewhere, motivated by frustration over previous dining challenges.
- E. Avoiding dining out altogether, motivated by concerns about public allergens.
- F. Offering advice to another parent in the café, motivated by her interest in parenting discussions.

**Correct Answer: B**

**GPT-4o's Answer: A**

**Analysis:** GPT-4o's analysis overlooks emotion-driven behavior by focusing on Yara's rational and altruistic motives, assuming she would share her experience to help others. This perspective ignores the possibility of a direct emotional response, such as expressing gratitude for the cafe service. Furthermore, GPT-4o overinterprets Yara's background in community management, predicting that her actions would be more focused on helping others or sharing experiences, rather than simply thanking the staff. In contrast, humans are more likely to recognize that, despite Yara's involvement in the community, her immediate interaction with the cafe staff would be influenced by her emotional response, such as gratitude, fitting the context of the situation.

Table 12: Case 3 on a Motive&Behavior Reasoning Question.

#### Case 4 - Motive&Behavior Reasoning Question

**Context:** Samantha owns a travel agency that specializes in personalized service and unique travel experiences. Recently, she's been thinking about ways to make her business more environmentally friendly and believes that adopting sustainable practices could also boost her agency's reputation. During a meeting with her accountant, Mark, they reviewed various financial strategies to implement her ideas. Although the agency already has basic recycling and uses digital communication to reduce waste, Samantha is determined to make a bigger impact in the competitive travel market.

**Question:** What kind of behavior is Samantha most likely to exhibit next, and what is the motivation behind it?

**Options:**

- A. Announcing a new program to contribute most profits to local environmental projects, motivated by a desire to build the agency's reputation for community involvement.
- B. Rushing to install solar panels on all properties without detailed cost planning, motivated by an urgent need to show visible commitment to sustainability.
- C. Delaying new projects until further discussions with stakeholders, motivated by caution about potential financial risks.
- D. Expanding the recycling program to engage customers in eco-friendly actions, motivated by a focus on community-based solutions.
- E. Launching a promotional campaign about the agency's past sustainable practices, motivated by the desire to draw media attention.
- F. Organizing workshops for employees on sustainable practices, motivated by a goal to enhance internal awareness.

**Correct Answer: D**

**GPT-4o's Answer: F**

**Analysis:** GPT-4o overlooked the emphasis on "enhancing market competitiveness" in the question and focused excessively on the superficial logic of "sustainability." However, the purpose of sustainability is to enhance market competitiveness, and merely raising internal employees' awareness does not contribute to improving market competitiveness.

Table 13: Case 4 on a Motive&Behavior Reasoning Question.



### Case 5 - Motive&Behavior Reasoning Question

**Context:** Fiona, a young woman working as an editor for a prestigious publishing house, lives alone in a vibrant urban neighborhood known for its diverse cultures. One afternoon, after a challenging week at her job, she decides to visit Bella Vita, a charming pizzeria in a more upscale area. Bella Vita is famous for its delicious pizzas and warm Italian atmosphere, complete with nostalgic music and friendly staff. As she sits by the window, enjoying the sunlight, she finds herself laughing softly, and her exhaustion starts to fade. The cozy ambiance surrounds her, bringing her feelings of comfort and joy. Fiona highly values her personal time, often enjoying these quiet moments for reflection and renewal, while also cherishing fond memories of family gatherings at similar Italian restaurants.

**Question:** As she listens to the familiar tunes and observes families enjoying meals together, based on what motivation is she most likely to exhibit what behavior next?

**Options:**

- A. Reach out to a friend to share her experience, motivated by her desire for emotional connection.
- B. Jot down her thoughts about the atmosphere, driven by her need for self-expression.
- C. Plan to revisit the restaurant with her family, inspired by her longing for shared memories.
- D. Explore other nearby restaurants, motivated by her curiosity about the local dining scene.
- E. Compliment the staff for their service, reflecting her appreciation for kindness and hospitality.
- F. Take a photograph to post online, motivated by her interest in sharing aesthetic moments with others.

**Correct Answer: C**

**GPT-4o's Answer: B**

**Analysis:** GPT-4o tends to over-rely on explicit textual details while overlooking implicit behavioral tendencies and deeper emotional motivations. For instance, it often focuses on directly stated traits in the prompt (e.g., “She is an editor, so she may prefer writing”) and limits its reasoning to surface-level information, ignoring how emotions like nostalgia might influence behavior. In contrast, humans naturally consider the emotional undertones within a situation, such as how a familial atmosphere may evoke empathy and drive planning. Additionally, GPT-4o primarily relies on explicit contextual details to infer motivations, whereas humans are more sensitive to subtle emotional cues embedded in the broader scenario, allowing for a more nuanced understanding of behavior.

Table 14: Case 5 on a Motive&Behavior Reasoning Question.