

# TOMVALLEY: EVALUATING THE THEORY OF MIND REASONING OF LLMs IN REALISTIC SOCIAL CONTEXT

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

As large language models (LLMs) are increasingly involved in human society, some studies try to evaluate LLMs’ capability of theory of mind (ToM), which is about the understanding and reasoning of others’ mental states and possible actions. However, these previous works simplify the ToM capability required in real social contexts during their evaluations. This can be reflected in three aspects: (1) most evaluations focus on a **static mental state** after several social scenarios while ignoring the changes of mental states across different scenarios; (2) they mainly consider **independent mental states**, however different kinds of mental states (beliefs, intentions, and emotions) and actions can influence one another in our real life; (3) there is an **absence of social settings and character profiles** in their evaluation, even though humans can effortlessly obtain and utilize this information in ToM reasoning processes. This lack can underestimate the abilities of LLMs. This paper aims to evaluate LLMs’ ToM capability in closer alignment with a realistic social context. Correspondingly, we propose a new benchmark, named TOMVALLEY, which alleviates the limitations mentioned above of previous works. Specifically, the benchmark is constructed using a framework that includes four steps: social background determination, mental state sketch, social scenario design, and rule-based question generation. Overall, there are 1100 social contexts and 78100 questions about characters’ mental states. The quality of the benchmark is manually verified. Additionally, we evaluate ten popular LLMs on TOMVALLEY. Experimental results suggest that LLMs’ performances are significantly inferior to human levels by 11%. Subsequent investigation indicates that LLMs are ineffective at interpreting alterations in mental states across social scenarios. Furthermore, we observe that LLMs are incapable of addressing compositional questions that necessitate multi-hop reasoning within the social context.<sup>1</sup>

## 1 INTRODUCTION

Theory of Mind (ToM) refers to the capacity to understand and reason about the mental states of others (e.g., beliefs, intentions, and emotions) and predict their next actions by inferring their mental states (Leslie et al., 2004; Call & Tomasello, 2008; Apperly & Butterfill, 2009). For human beings, ToM is a foundational capability in various daily social interactions (Turner, 1988), such as maintaining relationships (Hughes & Leekam, 2004), making decisions (Carlson & Moses, 2001), and enhancing peer popularity (Slaughter et al., 2015). As LLMs develop, they become increasingly involved in social activities and interact with humans daily. In these interactions, LLMs’ ToM capabilities significantly impact their performance. For instance, LLMs are expected to understand and reason about others’ beliefs, emotions and intentions, and even influence the users’ actions in support conversations (Liu et al., 2024c; Wang et al., 2024a).

Realizing the importance of LLMs’ ToM reasoning capacity, previous studies have proposed various benchmarks to evaluate LLMs, such as SocialIQA (Sap et al., 2019), BigToM (Gandhi et al., 2024), and TOMBENCH (Chen et al., 2024). Typically, these benchmarks contain questions about

<sup>1</sup><https://anonymous.4open.science/r/ToMValley-ICLR/README.md>

054 a character’s mental state, such as “given [the social scenario(s)], what does [a person] believe?”  
055 However, these works simplify the ToM capability required in real social contexts. The simplifica-  
056 tion is mainly reflected in the following three aspects. (1) **Static mental states**. Questions in most  
057 benchmarks only require LLMs to deduce the character’s final mental state after one or several so-  
058 cial scenarios. However, the mental state can change across different scenarios in our real-life social  
059 context. (2) **Independent mental states**. Most benchmarks only test one kind of mental state in one  
060 scenario. However, there exists an inter-relationship among different kinds of mental states (e.g.,  
061 beliefs, intentions, and emotions) (D’Andrade, 1995; Wellman, 1990). For instance, belief and emo-  
062 tion will influence one’s intention. (3) **Absence of social locations and character profiles**. Humans  
063 can effortlessly obtain and utilize social location information(e.g., cafes) and some character profiles  
064 in their ToM reasoning (Bretherton & Beeghly, 1982; Gönültaş et al., 2020). Their absence can lead  
065 to underestimation of LLMs’ ability and prohibit us from appropriately determining the causes of  
066 LLMs’ failure cases: whether it is due to the model’s inherently inadequate ToM reasoning abilities  
067 or the lack of corresponding information.

068 This paper aims to evaluate LLMs’ ToM capability in closer alignment with real-world social con-  
069 texts. Correspondingly, we construct a novel benchmark, named **TOMVALLEY**, and try to alleviate  
070 the three mentioned limitations of previous works by the following three considerations. (1) **Dy-**  
071 **namical mental states**. Rather than merely assessing the final mental state after a social scenario,  
072 **TOMVALLEY** incorporates questions that explore the alterations of mental states across multiple  
073 continuous social scenarios, corresponding to the dynamic nature of the mental state. (2) **Intra-**  
074 **dependent mental states**. **TOMVALLEY** investigates how different types of mental states influence  
075 one another. This involves constructing the intradependent relationships between mental states,  
076 shaping the progression of the social scenario plot based on the relationships, and tailoring questions  
077 to evaluate LLM’s ability to reason how different mental states influence each other. (3) **Provision**  
078 **of social location and character profiles**. Each social context in **TOMVALLEY** includes detailed  
079 social locations and character profiles, offering LLMs rich contextual information to reason about  
the characters’ mental states.

080 Notably, **TOMVALLEY** is constructed with the framework shown in Figure 1. As illustrated, it  
081 includes four main processes: (1) determining a social background, including the social location,  
082 character profiles, and the relationship between characters; (2) sketching the main character’s men-  
083 tal states across different social scenarios; (3) designing several social scenarios that happen on the  
084 main character based on the sketched mental states and the social background; (4) generating ques-  
085 tions, whose answers and options can be directly extracted from the output of step (2). The data  
086 in **TOMVALLEY** include the social background, social scenarios, and the questions (i.e., outputs of  
087 step (1), (3), and (4) of the framework). In summary, **TOMVALLEY** contains 1100 social contexts  
088 (2,200 characters and 5,500 social scenarios) and 78100 questions related to mental states. We em-  
089 ploy human annotators to evaluate and verify the quality of **TOMVALLEY**, and establish a human  
090 baseline. Feedback from annotators indicates that the social backgrounds and scenarios closely re-  
091 semble real-life interactions, the questions and respective options are reasonable, and the ground  
092 truths are validated despite humans not being 100% correct when establishing the human baseline.  
093 In addition, we evaluate ten popular LLMs’ ToM capabilities using both vanilla and CoT prompt-  
094 ing. Experimental results show that current LLMs underperform humans significantly on the ToM  
095 capability: even the best LLM’s performance, achieved by GPT-4o, is lower than humans’ by 11%.  
096 Further analysis shows that LLMs’ performance will decrease when useful information is presented  
097 in the middle of input. Moreover, we find that LLMs are incapable of processing compositional  
problems that necessitate multi-hop reasoning within the social context.

098 Overall, our contributions are as follows: (1) We propose to evaluate LLMs’ ToM capabilities in  
099 closer alignment with the real-world social contexts with the consideration of dynamic and intrade-  
100 pendent mental states as well as the provision of social location and character profiles in reasoning  
101 processes. (2) We introduce **TOMVALLEY**, a benchmark featuring personified individuals, diverse  
102 social locations, and dynamic and intradependent mental states. Through human evaluation, we  
103 demonstrate the benchmark’s uniqueness, difficulty, and high quality. (3) We conduct a thorough  
104 assessment of ten popular LLMs, and compare their performances against human performance, and  
105 provide an in-depth analysis of their limitations.

106  
107

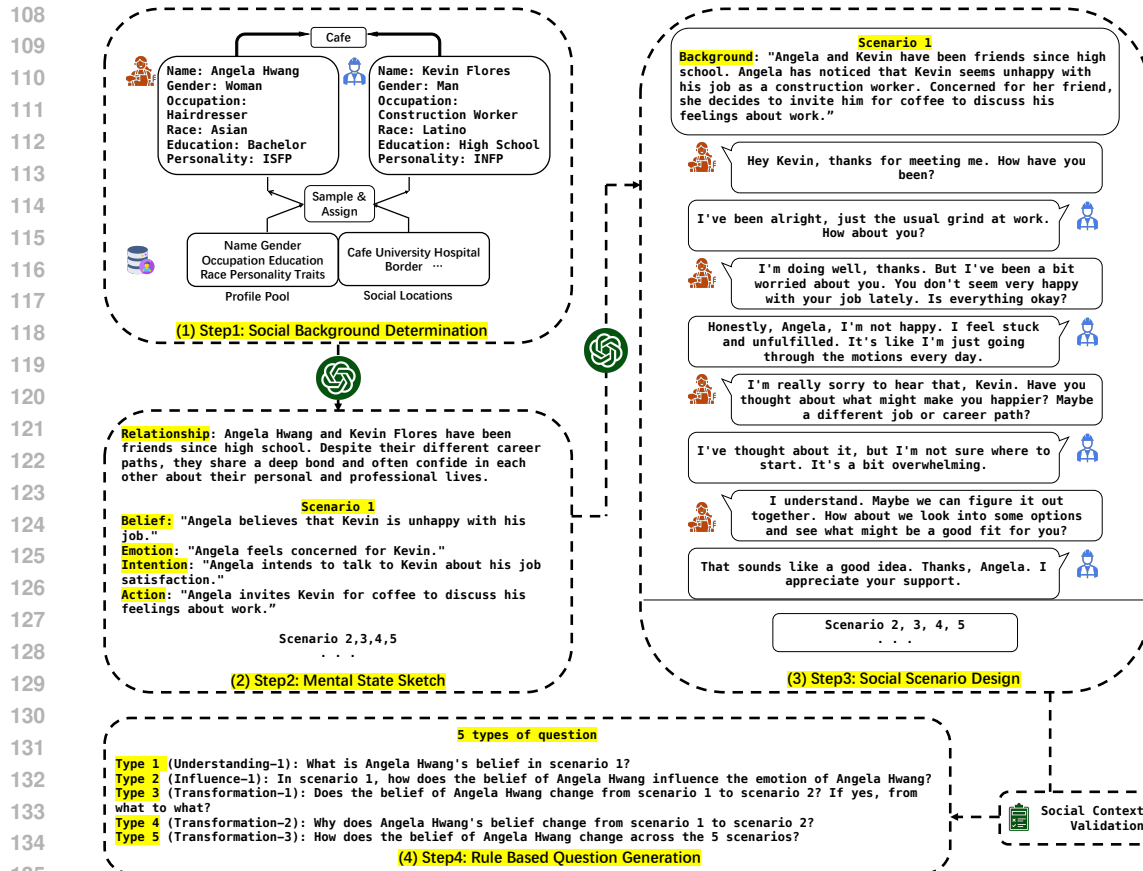


Figure 1: The framework used to generate the TOMVALLEY.

## 2 RELATED WORK

### 2.1 TOM BENCHMARKS

Arguably, infants as young as 12 months of age can attribute mental states to others, demonstrating theory of mind reasoning (Onishi & Baillargeon, 2005). Theory of mind appears to be an innate potential ability in humans that requires social and other experiences over many years for its full development. Researchers have begun to probe whether LLMs possess a Theory of Mind ability comparable to that of humans, as they have reached and occasionally surpassed human performance in some task-solving and reasoning tasks. Nematzadeh et al. (2018); Le et al. (2019); Wu et al. (2023) apply the Sally-Anne Test (Baron-Cohen et al., 1985) and bAbi (Weston et al., 2015) to test LLMs’ ToM ability in the aspect of false belief, and they find that LLMs’ performance is significantly lower than humans. Kosinski (2023); Bubeck et al. (2023); van Duijn et al. (2023) report GPT-series’ performance on Sally-Anne Test is comparable to or outperforms children aged 7 and 10. However, Ullman (2023); Shapira et al. (2024); Kim et al. (2023); Sap et al. (2022) propose that LLMs clearly lack of ToM without robust performance and LLMs are prone to shortcuts and spurious correlations. Apart from the test in the aspect of belief, Sap et al. (2019); Xu et al. (2024); Chen et al. (2024) construct benchmarks to test LLMs’ ToM ability for emotion, intention, and perception. Previous evaluations suffer from one or more of the following issues: static mental states, independent mental states, lack of the statement of social location, and absence of clear character information. Our work aims to develop a scalable, novel framework and benchmark to understand the ToM reasoning of language models in the dynamic social context.

## 2.2 PROCESS-LEVEL EVALUATION FOR REASONING.

Although it is essential to curate comprehensive and appropriate data for benchmarks, it is equally important to implement rigorous evaluation methodologies that scrutinize the step-by-step reasoning processes of AI models. The primary objective of most current benchmarks is to evaluate the model’s output relative to the standard answer at the answer level. Some recent works (Uesato et al., 2022; Lightman et al., 2024; Wang et al., 2024b) have begun to concentrate on the intermediate math reasoning stages of the models. In the social interaction evaluations, Zhou et al. (2024) propose SOTOPIA-EVAL to evaluate the multi-faceted social interactions, which not only require completing major social goals but also multiple implicit goals, such as maintaining relationships, preserving finances, gaining information, keeping secrets, and following social rules. Gandhi et al. (2024) propose a framework for procedurally designing synthetic ToM evaluations from causal templates to interpret the failure cases in ToM reasoning. However, this work only considers the constrained mental states, without considering the dynamic nature of the mental states, and it also ignores the mutual effect between mental states. Our work aims to integrate the benefits of process-level evaluation, creating a novel approach to generating benchmarks to probe the failure mode of LLMs across comprehensive mental states in the realistic social context.

## 3 TOMVALLEY CONSTRUCTION FRAMEWORK

**Definitions and Preliminaries** We would like to define key terms commonly used in this paper first. *Social Locations* refer to the physical settings of social scenarios. This information is important because it can reflect social norms and influence human behavior Farrow et al. (2017). *Social Scenario* denotes the interactions and activities involving characters. In previous works, such as OpenToM Xu et al. (2024), a single scenario may encompass multiple events occurring at different times. However, in our study, each scenario represents one specific social event at a given moment. *Social Context* encompasses the social background, including social locations and character profiles, as well as multiple social scenarios. This study emphasizes evaluating LLMs’ ability to reason about dynamic and intradependent mental states, taking into account social locations and character profiles. Accordingly, we propose a novel framework to generate an evaluation benchmark, as shown in Figure 1. The framework consists of four steps: (1) Social Background Determination, (2) Dynamic and Intradependent Mental State Sketch, (3) Social Scenario Design, and (4) Social Context Validation and Rule-Based Question Generation. Notably, this framework is scalable, allowing us to easily adjust the social context and its corresponding questions by modifying the character number, the dialogue turn, the scenario number, and the question number.

**Step 1: Social Background Determination** The social background contains a social location, character profiles, and the characters’ relationship. We collect a location pool and randomly sample one location each time. As for the character’s profile, we define seven aspects: surname, name, gender, occupation, education, race, and personality traits. For every aspect, we construct a pool of candidates. After determining the character number in the social context, we construct the profile for each character by sampling one item from each of the seven pools. The relationship between characters is generated with LLMs given the characters’ profiles.

**Step 2: Dynamic and Intradependent Mental State Sketch** This work focuses on dynamic and intradependent mental states. For “dynamic”, there should be several social scenarios, and the mental states in one scenario can be different from those in the last scenario. We set the number of scenarios as five.<sup>2</sup> For “intradependent”, multiple kinds of mental states as well as their intradependences are expected. We include three kinds of mental states: the belief, the emotion, and the intention. In addition, we take the action into account in the sketch since its close relationships to these three mental states. For clarity, we refer to any belief, emotion, intention, or action as a “ToM reasoning item.” To sketch such mental states, we prompt LLMs by providing the social background determined in Step 1. Notably, we sketch the mental states of one character, the main character, who LLMs primarily perceive and reason about.

<sup>2</sup>Due to our utilization of a commercial model, GPT-4-turbo-2024-04, for scenario design, we have limited the number of scenarios to five to reduce costs while maintaining a more authentic social context. Researchers can easily adjust the scenario number in our framework to meet their needs.

Item	Number
ToM Reasoning Items	4
Social Locations	261
Characters	2200
Scenarios	5500
Social Contexts	1100
Questions	78100
Average Social Context Length	457.9
Average Questions Length	77.5

Table 1: TOMVALLEY Statistics.




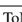
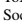
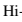
	Plot	 Social Location	 Profile	 Relationship	 Intrdependent mental states	 Dynamic mental states	 Questions Num
ToMi	✗	✗	✗	✗	✗	✗	999
SocialQA	✗	✗	✗	✗	✗	✗	37588
Hi-ToM	✗	✓	✗	✗	✗	✗	1200
OpenToM	✓	✗	✓	✓	✗	✗	2384
BigToM	✓	✓	✗	✗	✗	✗	600
TOMBENCH	✓	✗	✗	✗	✗	✗	2860
TOMVALLEY	✓	✓	✓	✓	✓	✓	78100

Table 2: ToM benchmark Comparison.

**Step 3: Social Scenario Design** We plot the social scenario in the format of dialogues among characters. This is because dialogue is the primary format through which LLMs perceive and interact with humans. Meanwhile, ToM reasoning in dialogues has seldom been investigated in previous works. The social scenarios focus on a single topic and evolve over time. When designing each social scenario, the utterances of the main character in the dialogue are motivated by their corresponding mental states sketched in Step 2. Specifically, we prompt an LLM to generate the dialogues in five scenarios given the social background (Step 1 output) and the main character’s mental state sketch (Step 2 output).

**Step 4: Social Context Validation and Rule-Based Question Generation** The LLMs’ outputs are not reliable all the time. Thus, we need to check the quality of the mental state sketch and social scenarios. In specific, there are three principles: (1) Does the mental state sketch consist of each scenario’s belief, emotion, intention, and action? (2) Are the numbers of social scenarios in Step 2 and Step 3 as expected? (3) Are the dialogues coherent with the mental state sketch? Based on these three principles, unqualified instances will be removed by regular matching. Then, we generate questions that can evaluate ToM capability using the qualified instances. We aim to explore five types of questions: (1) (Understanding-1) What is the main character’s ToM reasoning item in a specific scenario? (2) (Influence-1) In one scenario, how does mental state A influence ToM reasoning item B? (3) (Transformation-1) Does a ToM reasoning item change from scenario A to scenario B? (4) (Transformation-2) What causes a ToM reasoning item change from scenario A to scenario B? (5) (Transformation-3) How does the ToM reasoning item change across all the scenarios? We designed five templates to generate questions based on the five question types. We use these templates to generate 71 questions for each social context. Notably, action is solely instigated by intention and does not directly influence other toM reasoning item (d’Andrade, 1987). Thus, we do not set any influence question for action. The ground truth and misleading options can be extracted directly from the mental state sketch.

More details of each step can be found in Appendix A.

## 4 TOMVALLEY BENCHMARK

### 4.1 STATISTICS

Leveraging the framework proposed in Section 3 and GPT-4-Turbo<sup>3</sup>, we construct the benchmark TOMVALLEY. We set the number of characters as two and the number of scenarios as five for each social context. In addition, the number of locations in the location pool is 261. We generate 1,100 social contexts and 78,100 questions. Each social context includes one social location, two character profiles, the character relationship, and dialogues between characters in five different scenarios. And there are 71 questions about each social context. Table 1 shows the statistics. Moreover, we compare TOMVALLEY with previous benchmarks, and comparison results are shown in Table 2. Although half of previous works have provided a plot (things happened on characters) for ToM reasoning, most of them usually ignore character profiles, relationships, and dynamic mental states. And only our work considers the intrdependent mental states.

<sup>3</sup>We accessed GPT-4-turbo-2024-04-09 through Microsoft Azure OpenAI service in August 2024.

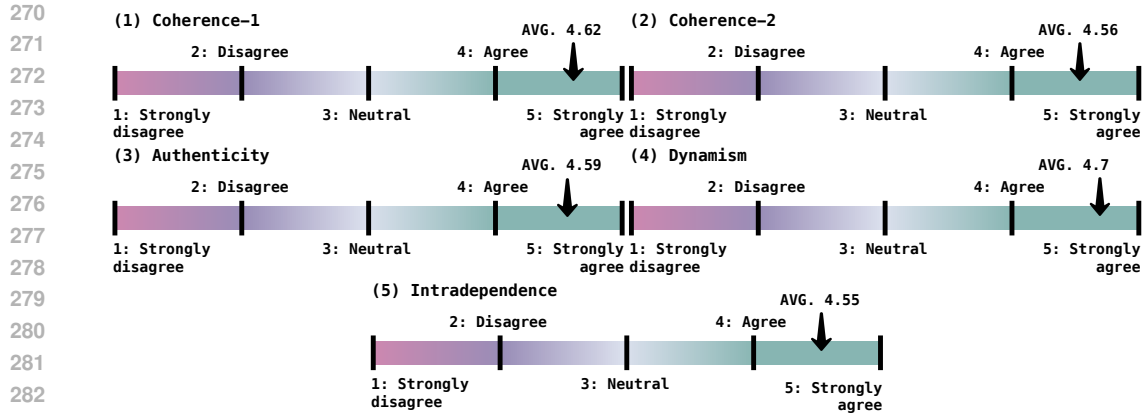


Figure 2: Human evaluation results of the quality of social contexts.

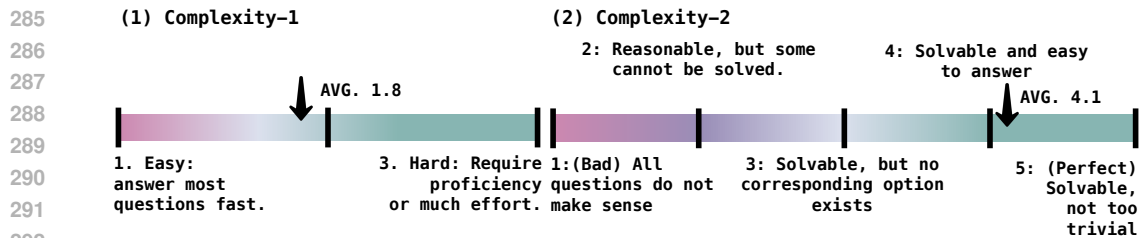


Figure 3: Human evaluation results of the quality of questions.

## 4.2 SOCIAL CONTEXT QUALITY

We evaluate the quality of the social contexts in the benchmark via human evaluation. In specific, we hired five graduate students, and randomly sampled 60 social contexts and 1,420 questions from the benchmark. In addition, we collect the corresponding mental state sketches – outputs of Step 2 in the framework. We present all these contents to the human annotators and ask them to rate the five 5-Likert scale questions: (1) *coherence-1*: to what degree do you agree that the dialogue between characters in the five scenarios is coherent? (2) *coherence-2*: to what degree do you agree that the social background (character profiles, the social location, the character relationship) and the scenarios are coherent? (3) *authenticity*: to what degree do you agree that the social context is authentic and aligns with real life? (4) *dynamism*: to what degree do you agree that the mental state’s change of the characters from one scenario to the following scenario is reasonable? (5) *intradependence*: to what degree do you agree that the mental state’s influence on one another is reasonable? Here, 1 indicates strongly disagree, while 5 indicates strongly agree.

Figure 2 presents the human evaluation results. For coherence-1 and coherence-2, the average ratings reach 4.62 and 4.56, respectively, indicating the good quality of the social context. The average rating of authenticity reaches 4.59, and 90% of annotators rate with 4 or 5, showing that most annotators agree that social contexts are authentic and align with real life. 93% of annotators agree that the mental state’s change of the characters from one scenario to the following scenario is reasonable, and the average rating of dynamism is 4.7.

## 4.3 QUESTION QUALITY AND COMPLEXITY

We would like to explore the quality and complexity of questions in TOMVALLEY: Do the questions have varying degrees of complexity? Correspondingly, we conduct a human evaluation. We invited five human annotators to answer questions about social contexts. Each annotator was assigned 20 social contexts and 1,420 questions, randomly sampled from the benchmark. After the participants finished all the questions or decided not to proceed, we asked them to rate the complexity level and the quality of the questions and the corresponding options. In specific, we asked them to answer two questions. One is a 3-Likert scale question: (1) *complexity-1*: please rate the difficulty of the questions related to the social context, where 1 and 3 indicate easy and hard, respectively. Another

Subject	📖 : Understanding			🗨️ : Influence			🔗 : Transformation			AVG.		
	Belief			Emotion			Intention				Action	
	📖	🗨️	🔗	📖	🗨️	🔗	📖	🗨️	🔗		📖	🔗
Human	0.84	0.85	0.78	0.90	0.80	0.79	0.79	0.72	0.74	0.77	0.76	0.78
GPT-4o	0.81	0.64	0.45	0.92	0.81	0.46	0.88	0.92	0.51	0.95	0.55	0.67
GPT-4-Turbo	0.63	0.46	0.33	0.75	0.53	0.34	0.72	0.75	0.35	0.80	0.37	0.50
Llama-3.1-70B	0.66	0.36	0.40	0.93	0.63	0.43	0.83	0.82	0.42	0.92	0.46	0.58
Llama-3.1-8B	0.31	0.28	0.19	0.39	0.27	0.19	0.22	0.25	0.17	0.27	0.15	0.23
Mixtral-8x7B	0.24	0.20	0.22	0.46	0.41	0.19	0.33	0.51	0.10	0.40	0.09	0.25
Mistral-7B	0.21	0.20	0.12	0.23	0.25	0.11	0.17	0.25	0.10	0.20	0.10	0.16
Qwen2-72B	0.72	0.40	0.38	0.86	0.65	0.38	0.80	0.87	0.34	0.89	0.20	0.53
Qwen2-7B	0.23	0.19	0.19	0.44	0.34	0.20	0.26	0.19	0.16	0.24	0.15	0.22
DeepSeek-V2	0.07	0.15	0.10	0.05	0.10	0.08	0.04	0.10	0.07	0.03	0.06	0.08
GLM-4	0.30	0.34	0.24	0.44	0.31	0.20	0.29	0.30	0.17	0.40	0.16	0.26
LLM AVG.	0.42	0.32	0.26	0.55	0.43	0.26	0.45	0.50	0.24	0.51	0.23	0.35
GPT-4o+CoT	0.79	0.59	0.44	0.88	0.72	0.47	0.82	0.84	0.47	0.90	0.50	0.63
GPT-4-Turbo+CoT	0.61	0.46	0.30	0.78	0.52	0.33	0.72	0.69	0.32	0.81	0.37	0.49
Llama-3.1-70B+CoT	0.68	0.40	0.39	0.91	0.62	0.43	0.82	0.79	0.42	0.96	0.46	0.58
Llama-3.1-8B+CoT	0.31	0.28	0.21	0.40	0.25	0.21	0.21	0.25	0.20	0.24	0.16	0.24
Mixtral-8x7B+CoT	0.16	0.16	0.14	0.29	0.27	0.13	0.25	0.32	0.09	0.26	0.08	0.18
Mistral-7B+CoT	0.21	0.21	0.11	0.22	0.25	0.11	0.20	0.25	0.09	0.19	0.09	0.16
Qwen2-72B+CoT	0.71	0.38	0.40	0.87	0.68	0.41	0.83	0.87	0.35	0.88	0.27	0.55
Qwen2-7B+CoT	0.28	0.17	0.18	0.43	0.36	0.19	0.30	0.22	0.19	0.20	0.18	0.23
DeepSeek-V2+CoT	0.08	0.17	0.09	0.04	0.11	0.10	0.05	0.13	0.07	0.05	0.06	0.09
GLM-4+CoT	0.30	0.36	0.26	0.48	0.30	0.22	0.33	0.31	0.17	0.43	0.15	0.28
LLM+CoT AVG.	0.41	0.32	0.25	0.53	0.41	0.26	0.45	0.47	0.24	0.49	0.23	0.34

Table 3: LLMs’ performances. We show the performance according to the combination of ToM reasoning items and question types. The items include belief, emotion, intention, and action. The question types include understanding, influence, and transformation. “LLM AVG.” and “LLM+CoT AVG.” are the average performance of all the ten LLMs and LLMs+CoT, respectively. The best performance of vanilla prompting is in purple, and that of CoT prompting is in grey.

is a 5-Likert scale question: (2) *complexity-2*: please rate the quality of the questions and options. 1 is “(Bad) all questions do not make sense,” while 5 is “(Perfect) Solvable, not too trivial.”

Figure 3 presents the results. We observe that the questions have varying complexity levels: 23.8% of the questions can be quickly solved by most annotators, 69% require some effort, and 7.2% can only be solved by professionals or with great effort.

More details about the human evaluation of social contexts and questions in the benchmark can be found in Appendix A.6.

## 5 EXPERIMENTS

### 5.1 EXPERIMENTAL SETUP

We use TOMVALLEY to evaluate the ToM abilities of 10 popular LLMs, including GPT-4o (OpenAI, 2024), GPT-4-Turbo (Achiam et al., 2023), Llama-3.1-8B (AI@Meta, 2024), Llama-3.1-70B (AI@Meta, 2024), Mistral-7B (AI, 2024a), Mixtral-8x7B (AI, 2024b), Qwen2-7B (Yang et al., 2024), Qwen2-72B (Yang et al., 2024), DeepSeek-V2 (Liu et al., 2024a), GLM-4 (GLM et al., 2024). We strictly abide by all LLMs’ terms and get access through official APIs or model weights. We employ two prompting methods: vanilla prompting directly asking LLMs to answer the questions, and CoT prompting elicits step-by-step reasoning before answering. A human baseline is established by inviting five graduate students to complete a subset of TOMVALLEY. Details about model versions, parameter sizes, context window sizes, and the prompts used for the two methods are shown in Appendix B.

	ToM Reasoning Item Type				Question Type		
	Belief	Emotion	Intention	Action	Understanding	Influence	Transformation
Human	0.81	0.81	0.73	0.75	0.82	0.79	0.77
GPT-4o	0.59	0.68	0.72	0.69	0.89	0.79	0.49
GPT-4-Turbo	0.44	0.50	0.55	0.52	0.73	0.58	0.35
Llama-3.1-70B	0.46	0.61	0.63	0.63	0.84	0.60	0.43
Llama-3.1-8B	0.24	0.26	0.21	0.19	0.30	0.27	0.17
Mixtral-8x7B	0.22	0.32	0.27	0.20	0.36	0.37	0.15
Mistral-7B	0.16	0.18	0.16	0.14	0.20	0.23	0.11
Qwen2-72B	0.47	0.58	0.60	0.45	0.81	0.64	0.33
Qwen2-7B	0.20	0.30	0.20	0.19	0.29	0.24	0.18
GLM-4	0.28	0.29	0.23	0.24	0.36	0.31	0.19
LLM AVG.	0.34	0.41	0.40	0.36	0.53	0.45	0.27

Table 4: LLMs’ performance in vanilla prompting. We show the performance according to ToM reasoning items and question types, respectively. Due to space limitations, we don’t show the results in CoT prompting, which is shown in Appendix B.3.

## 5.2 MAIN RESULTS

Table 3 and 4 demonstrate the ToM performances of LLMs according to ToM reasoning items and question types. As mentioned, the ToM reasoning items include belief, emotion, intention, and action; and the question types include understanding, influence, and transformation. We discuss the results and highlight several critical observations as follows.

**Human vs. LLMs** Humans achieve 78% accuracy performance. However, ToM performances of all LLMs are significantly lower, with the smallest gap being 11% in vanilla prompting (Human 78% vs. GPT-4o 67%). Among the question types, LLMs fall behind humans in transformation in all ToM reasoning items, representing the most challenging reasoning for LLMs. Interestingly, in the understanding type, LLMs like GPT-4o even outperform humans, which we believe is explainable. As shown in Figure 1, the understanding question mainly directly asks what the mental state of the character is in one scenario, which requires less reasoning process compared to other questions, and LLMs can easily answer these questions by semantic matching.

**LLMs’ ToM Performance** In the vanilla prompting, GPT-4o stands out in LLMs and surpasses the second of Llama-3.1-70B up to 9% points. Among the open LLMs, Llama-3.1-70B and Qwen2-72B impressively outperforms other LLMs and even outperforms GPT-4-Turbo. The Llama-3.1-70B also surpasses the GPT-4o in the understanding type question of emotion ability. However, even the most superior GPT-4o only reaches 67%, and the lowest score is 8%, which shows the difficulty of our benchmark, and current LLMs lack robust ToM reasoning in the social context environment, even if they reach nearly 100% performance in other benchmarks (Gandhi et al., 2024).

**Vanilla vs. CoT Prompting** Both ability and question-type results indicate that CoT prompting doesn’t always improve LLMs’ ToM reasoning ability. The CoT successfully improved their performance for Llama-3.1-8B, Qwen2-series, DeepSeek-V2, and GLM-4. For other models, such as GPT-4o, the CoT prompting even leads to a decline in performance. This finding aligns with the findings in Xiao et al. (2023); Chen et al. (2024). One plausible explanation is that CoT reasoning predominantly works by deconstructing intricate problems into more easy sub-tasks. Nonetheless, CoT cannot help to improve basic ToM ability. So, CoT will only work for models that do not have the ability to decompose complex problems. We present a failure case of GPT-4o when CoT prompting is used in Appendix B.5.

**Differences Across ToM Reasoning Items** Exploring Table 4, the best-performing ToM reasoning item is emotion, consistent with the results found in TOMBENCH (Chen et al., 2024). Belief is the weakest item, trailing emotion by 7%. This demonstrates that the model is less proficient in addressing belief-related issues than the other 3 items. One potential explanation is that the character’s belief is more concealed than other mental states, necessitating the model to perform additional



reasoning steps, particularly when the questions necessitate the model to infer the character’s belief from observed actions, which align with the findings in BigToM (Gandhi et al., 2024).

**Differences Across Question Types** Further exploring Table 4, the model performs poorly in influence and transformation types compared with understanding, especially the question type of transformation, which lags behind understanding by 26 percent. The transformation question mainly relates to mental states’ transformation in different scenarios, such as how the character’s mental state changes from one scenario to another. This shows that the models lack the ability to handle the mental states in the middle scenarios of social context, where the dynamism of the mental state across scenarios is inherent in the social interaction.

**Profile Absence vs. Profile Presence** As stated by previous studies, personal profile information and the social norms behind the social location can help human to accurately ascertain an individual’s mental state (Bretherton & Beeghly, 1982; Strang, 1930). To further examine whether LLMs’ ToM performance will be influenced by the absence of this kind of information, we randomly pick 100 social contexts, a total of 7100 questions, and evaluate GPT-4o’s performance with and without the characters’ profile as part of the input, denoted as presence and absence respectively. As shown in Figure 4, the model lags behind 7 percent at most when there is no profile as part of the input, which aligns with the findings in the psychological literature. Most current work does not include such information, so their benchmarks can only evaluate LLMs’ performance without such information. Our benchmark allowing us to more completely evaluate the model’s performance in various conditions.

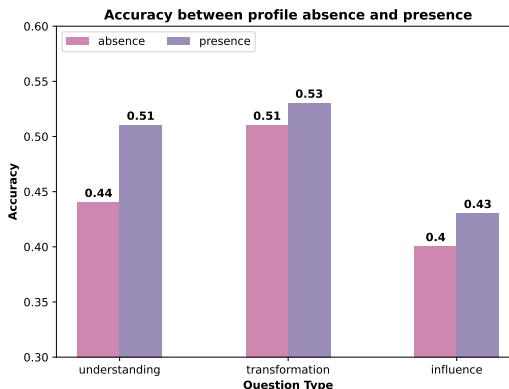


Figure 4: The accuracy of GPT-4o in the 3 question types with the presence of profile and absence of profile.

### 5.3 IN-DEPTH ANALYSIS

**LLMs Fail in the Middle Scenario** In Table 4, we find that the model performs poorly in transformation type compared with understanding and influence for all the models. To further explore the reason why models perform poorly when handling the transformation of ToM reasoning items, we further show the models’ performance of transformation along the time span. The time span indicates the specific scenarios to which one question relates. For example, the type 3 question in Figure 1 relates scenarios 1 and 2 and the ToM reasoning item of belief. As shown in Figure 5, the model performs better in the early and the end scenarios, while the model performs worse in the middle scenario. We posit that this may result from the model’s “Lost in the middle” phenomenon (Liu et al., 2024b): models exhibit diminished performance when crucial information is retrieved from extensive contexts, with optimal performance typically occurring at the beginning or end of the input context. The social scenario is presented to LLMs from scenarios 1 to 5; hence, the time order correlates with the position in the input, and LLMs perform poorly for the

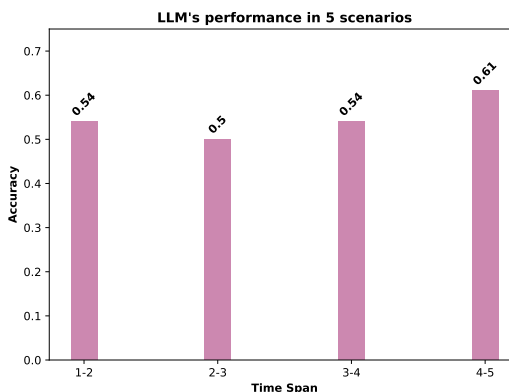


Figure 5: The average of models’ scores in the transformation question type. The time span indicates the specific scenarios to which one question relates.

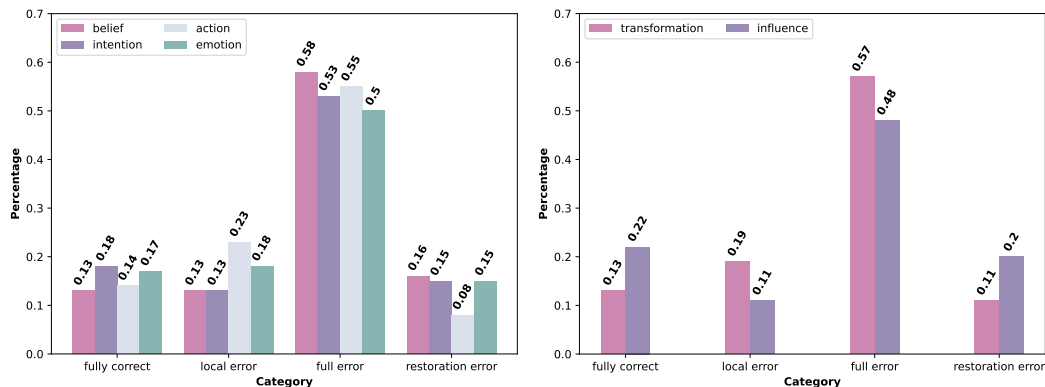


Figure 6: The percentage of four status of the compositional problem. The left one is depicted in accordance with ToM reasoning items. The right one is demonstrated through the question types.

middle time span. To confirm this conjecture, we further construct 15 social contexts with 4, 6, and 7 scenarios, respectively. The results are shown in Appendix B.4, which confirms our speculation.

**LLMs’ Limits of ToM on Compositionality** Our benchmark includes a kind of compositional problem requiring rigorous multi-hop reasoning to reach the answer (Dziri et al., 2024). Utilize the type 3 illustrated in Figure 1 as a reference. To address this question, it is essential to first ascertain the character’s beliefs in scenarios 1 and 2 (referred to as dependency D), followed by identifying the potential factors that induce the transformation of belief (denoted as C). Usually, there are four statuses (Dziri et al., 2024) of C: (1) **fully correct**: LLMs correctly answer both D and C. (2) **local error**: LLMs only correctly answer D and wrongly answer C. (3) **restoration error**: LLMs correctly answer C but wrongly answer D. (4) **full error**: LLMs wrongly answer both D and C. Examining the failure modes of LLMs in this kind of question might elucidate if models genuinely develop ToM reasoning capability or merely depend on shortcut learning through pattern matching to answer questions. In Figure 6, we visualize the proportion of the four statuses of GPT-4o. Notably, the figure does not display the understanding questions, as they typically function as dependence problems. Of the four statuses of C, the restoration error constitutes approximately 15%, indicating that LLMs may respond to certain questions based on superficial patterns rather than via logical reasoning. The fully correct status only accounts for approximately 15%, indicating that our benchmark’s compositional questions pose significant challenges for LLMs. The fully correct status of the transformation questions accounts for merely 13%, which is lower than that of influence questions; this implies that LLMs are less proficient at deducing the alterations of the ToM reasoning items.

## 6 CONCLUSION

In this paper, we propose TOMVALLEY, a benchmark to evaluate LLMs’ ToM ability in a realistic social context. Different from most previous benchmarks for ToM, TOMVALLEY evaluates LLMs with three characteristics: (1) Dynamic mental states. Rather than merely assessing the final mental state after several social scenarios, TOMVALLEY incorporates questions exploring mental state alterations across multiple continuous social scenarios. (2) Intradependent mental states. TOMVALLEY investigates how different types of mental states influence one another. (3) Provision of social location and character profiles. Feedback from annotators indicates that our evaluation data closely resembles real-life interactions. Experimental results show that current LLMs underperform humans significantly on the ToM capability: even the best LLM’s performance is lower than humans’ by 11%. Further analysis shows that LLMs’ performance will decrease when useful information is presented in the middle of input. Moreover, we find that LLMs are incapable of processing compositional problems. With the development of LLMs, they have played an important role in constructing AI applications, such as embodied intelligence and AI agents. We hope that TOMVALLEY will drive the evaluation of ToM in closer alignment with the requirements in real social interactions, facilitating the improvement of LLMs with better ToM abilities.

## REFERENCES

- 540  
541  
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
544 report. [arXiv preprint arXiv:2303.08774](https://arxiv.org/abs/2303.08774), 2023.
- 545 Mistral AI. Mistral 7B. <https://mistral.ai/news/announcing-mistral-7b/>,  
546 2024a. [Online; accessed August-2024].
- 547 Mistral AI. Mixtral of experts. <https://mistral.ai/news/mixtral-of-experts/>,  
548 2024b. [Online; accessed August-2024].
- 549  
550 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/  
551 llama3/blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- 552 Ian A Apperly and Stephen A Butterfill. Do humans have two systems to track beliefs and belief-like  
553 states? [Psychological review](https://doi.org/10.1080/09500800903271111), 116(4):953, 2009.
- 554 Simon Baron-Cohen, Alan M Leslie, and Uta Frith. Does the autistic child have a “theory of mind”?  
555 [Cognition](https://doi.org/10.1080/00140138508250013), 21(1):37–46, 1985.
- 556 Inge Bretherton and Marjorie Beeghly. Talking about internal states: The acquisition of an explicit  
557 theory of mind. [Developmental psychology](https://doi.org/10.1037/0012-1649.18.6.906), 18(6):906, 1982.
- 558 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Ka-  
559 mar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general  
560 intelligence: Early experiments with gpt-4. [arXiv preprint arXiv:2303.12712](https://arxiv.org/abs/2303.12712), 2023.
- 561  
562 Josep Call and Michael Tomasello. Does the chimpanzee have a theory of mind? 30 years later.  
563 [Trends in cognitive sciences](https://doi.org/10.1016/0169-8346(98)00068-9), 12(5):187–192, 2008.
- 564 Stephanie M Carlson and Louis J Moses. Individual differences in inhibitory control and children’s  
565 theory of mind. [Child development](https://doi.org/10.1111/j.1467-8624.2001.00103.x), 72(4):1032–1053, 2001.
- 566  
567 Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen, Guanqun Bi, Gongyao Jiang, Yaru Cao,  
568 Mengting Hu, Yunghwei Lai, Zexuan Xiong, and Minlie Huang. ToMBench: Benchmarking  
569 theory of mind in large language models. In Lun-Wei Ku, Andre Martins, and Vivek Sriku-  
570 mar (eds.), [Proceedings of the 62nd Annual Meeting of the Association for Computational  
571 Linguistics \(Volume 1: Long Papers\)](https://doi.org/10.18653/v1/2024.acl-long.847), pp. 15959–15983, Bangkok, Thailand, August 2024. As-  
572 sociation for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.847. URL [https:  
573 //aclanthology.org/2024.acl-long.847](https://aclanthology.org/2024.acl-long.847).
- 574 Roy d’Andrade. A folk model of the mind. 1987.
- 575  
576 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean  
577 Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. Faith and fate: Limits of  
578 transformers on compositionality. [Advances in Neural Information Processing Systems](https://arxiv.org/abs/2401.14493), 36, 2024.
- 579 R D’Andrade. [The development of cognitive anthropology](https://doi.org/10.1017/C9780521428396). The Cambridge University Press, 1995.
- 580  
581 Katherine Farrow, Gilles Grolleau, and Lisette Ibanez. Social norms and pro-environmental behav-  
582 ior: A review of the evidence. [Ecological Economics](https://doi.org/10.1016/j.ecolecon.2017.03.001), 140:1–13, 2017.
- 583 Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. Understanding  
584 social reasoning in language models with language models. [Advances in Neural Information  
585 Processing Systems](https://arxiv.org/abs/2401.14493), 36, 2024.
- 586  
587 Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu  
588 Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng,  
589 Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu,  
590 Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao,  
591 Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu,  
592 Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan  
593 Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang,  
Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language  
models from glm-130b to glm-4 all tools, 2024.

- 594 Seçil Gönültaş, Bilge Selçuk, Virginia Slaughter, John A Hunter, and Ted Ruffman. The capricious  
595 nature of theory of mind: Does mental state understanding depend on the characteristics of the  
596 target? Child development, 91(2):e280–e298, 2020.
- 597  
598 Claire Hughes and Sue Leekam. What are the links between theory of mind and social relations?  
599 review, reflections and new directions for studies of typical and atypical development. Social  
600 development, 13(4):590–619, 2004.
- 601 Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Bras, Gunhee Kim, Yejin Choi, and Maarten  
602 Sap. FANToM: A benchmark for stress-testing machine theory of mind in interactions. In  
603 Houda Bouamor, Juan Pino, and Kalika Bali (eds.), Proceedings of the 2023 Conference on  
604 Empirical Methods in Natural Language Processing, pp. 14397–14413, Singapore, December  
605 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.890. URL  
606 <https://aclanthology.org/2023.emnlp-main.890>.
- 607  
608 Michal Kosinski. Theory of mind might have spontaneously emerged in large language models.  
609 arXiv preprint arXiv:2302.02083, 2023.
- 610 Matthew Le, Y-Lan Boureau, and Maximilian Nickel. Revisiting the evaluation of theory of mind  
611 through question answering. In Proceedings of the 2019 Conference on Empirical Methods in  
612 Natural Language Processing and the 9th International Joint Conference on Natural Language  
613 Processing (EMNLP-IJCNLP), pp. 5872–5877, 2019.
- 614 Alan M Leslie, Ori Friedman, and Tim P German. Core mechanisms in ‘theory of mind’. Trends in  
615 cognitive sciences, 8(12):528–533, 2004.
- 616  
617 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan  
618 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In The Twelfth  
619 International Conference on Learning Representations, 2024. URL [https://openreview.  
620 net/forum?id=v8L0pN6EOi](https://openreview.net/forum?id=v8L0pN6EOi).
- 621 Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong  
622 Ruan, Damai Dai, Daya Guo, et al. Deepseek-v2: A strong, economical, and efficient mixture-  
623 of-experts language model. arXiv preprint arXiv:2405.04434, 2024a.
- 624  
625 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and  
626 Percy Liang. Lost in the middle: How language models use long contexts. Transactions of the  
627 Association for Computational Linguistics, 12:157–173, 2024b.
- 628  
629 Tianjian Liu, Hongzheng Zhao, Yuheng Liu, Xingbo Wang, and Zhenhui Peng. Compeer: A gener-  
630 ative conversational agent for proactive peer support. arXiv preprint arXiv:2407.18064, 2024c.
- 631 Aida Nematzadeh, Kaylee Burns, Erin Grant, Alison Gopnik, and Tom Griffiths. Evaluating theory  
632 of mind in question answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi  
633 Tsujii (eds.), Proceedings of the 2018 Conference on Empirical Methods in Natural Language  
634 Processing, pp. 2392–2400, Brussels, Belgium, October–November 2018. Association for Com-  
635 putational Linguistics. doi: 10.18653/v1/D18-1261. URL [https://aclanthology.org/  
636 D18-1261](https://aclanthology.org/D18-1261).
- 637  
638 Kristine H Onishi and Renée Baillargeon. Do 15-month-old infants understand false beliefs?  
639 science, 308(5719):255–258, 2005.
- 640  
641 OpenAI. GPT-4o System Card. <https://openai.com/index/gpt-4o-system-card/>,  
642 2024. [Online; accessed August-2024].
- 643  
644 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Com-  
645 monsense reasoning about social interactions. In Kentaro Inui, Jing Jiang, Vincent Ng, and  
646 Xiaojun Wan (eds.), Proceedings of the 2019 Conference on Empirical Methods in Natural  
647 Language Processing and the 9th International Joint Conference on Natural Language Processing  
(EMNLP-IJCNLP), pp. 4463–4473, Hong Kong, China, November 2019. Association for Com-  
putational Linguistics. doi: 10.18653/v1/D19-1454. URL [https://aclanthology.org/  
D19-1454](https://aclanthology.org/D19-1454).

- 648 Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. Neural theory-of-mind? on the limits  
649 of social intelligence in large LMs. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.),  
650 Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pp.  
651 3762–3780, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational  
652 Linguistics. doi: 10.18653/v1/2022.emnlp-main.248. URL [https://aclanthology.org/  
653 2022.emnlp-main.248](https://aclanthology.org/2022.emnlp-main.248).
- 654 Natalie Shapira, Mosh Levy, Seyed Hossein Alavi, Xuhui Zhou, Yejin Choi, Yoav Goldberg,  
655 Maarten Sap, and Vered Shwartz. Clever hans or neural theory of mind? stress testing social  
656 reasoning in large language models. In Yvette Graham and Matthew Purver (eds.), Proceedings  
657 of the 18th Conference of the European Chapter of the Association for Computational Linguistics  
658 (Volume 1: Long Papers), pp. 2257–2273, St. Julian’s, Malta, March 2024. Association for Com-  
659 putational Linguistics. URL <https://aclanthology.org/2024.eacl-long.138>.
- 660 Virginia Slaughter, Kana Imuta, Candida C Peterson, and Julie D Henry. Meta-analysis of theory  
661 of mind and peer popularity in the preschool and early school years. Child development, 86(4):  
662 1159–1174, 2015.
- 663 Daniel Stokols. Environmental psychology. 1978.
- 664 Ruth Strang. Measures of social intelligence. American Journal of Sociology, 36(2):263–269, 1930.
- 665 Jonathan H Turner. A theory of social interaction. Stanford University Press, 1988.
- 666 Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia  
667 Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and  
668 outcome-based feedback. arXiv preprint arXiv:2211.14275, 2022.
- 669 Tomer Ullman. Large language models fail on trivial alterations to theory-of-mind tasks. arXiv  
670 preprint arXiv:2302.08399, 2023.
- 671 Max van Duijn, Bram van Dijk, Tom Kouwenhoven, Werner de Valk, Marco Spruit, and Pe-  
672 ter van der Putten. Theory of mind in large language models: Examining performance of 11  
673 state-of-the-art models vs. children aged 7-10 on advanced tests. In Jing Jiang, David Reiter,  
674 and Shumin Deng (eds.), Proceedings of the 27th Conference on Computational Natural  
675 Language Learning (CoNLL), pp. 389–402, Singapore, December 2023. Association for Com-  
676 putational Linguistics. doi: 10.18653/v1/2023.conll-1.25. URL [https://aclanthology.  
677 org/2023.conll-1.25](https://aclanthology.org/2023.conll-1.25).
- 678 Jiashuo Wang, Yang Xiao, Yanran Li, Changhe Song, Chunpu Xu, Chenhao Tan, and Wenjie  
679 Li. Towards a client-centered assessment of llm therapists by client simulation. arXiv preprint  
680 arXiv:2406.12266, 2024a.
- 681 Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhi-  
682 fang Sui. Math-shepherd: Verify and reinforce LLMs step-by-step without human annotations. In  
683 Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Proceedings of the 62nd Annual Meeting  
684 of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 9426–9439,  
685 Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/  
686 v1/2024.acl-long.510. URL <https://aclanthology.org/2024.acl-long.510>.
- 687 HM Wellman. The child’s theory of mind. Cambridge/Bradford, 1990.
- 688 Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart Van Merriënboer, Armand  
689 Joulin, and Tomas Mikolov. Towards ai-complete question answering: A set of prerequisite toy  
690 tasks. arXiv preprint arXiv:1502.05698, 2015.
- 691 Yufan Wu, Yinghui He, Yilin Jia, Rada Mihalcea, Yulong Chen, and Naihao Deng. Hi-ToM:  
692 A benchmark for evaluating higher-order theory of mind reasoning in large language mod-  
693 els. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), Findings of the Association for  
694 Computational Linguistics: EMNLP 2023, pp. 10691–10706, Singapore, December 2023. As-  
695 sociation for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.717. URL  
696 <https://aclanthology.org/2023.findings-emnlp.717>.

702 Yang Xiao, Yi Cheng, Jinlan Fu, Jiashuo Wang, Wenjie Li, and Pengfei Liu. How far are we from  
 703 believable ai agents? a framework for evaluating the believability of human behavior simulation.  
 704 arXiv preprint arXiv:2312.17115, 2023.

705 Hainiu Xu, Runcong Zhao, Lixing Zhu, Jinhua Du, and Yulan He. OpenToM: A comprehensive  
 706 benchmark for evaluating theory-of-mind reasoning capabilities of large language models. In  
 707 Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Proceedings of the 62nd Annual Meeting  
 708 of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 8593–8623,  
 709 Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/  
 710 v1/2024.acl-long.466. URL <https://aclanthology.org/2024.acl-long.466>.

711 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
 712 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,  
 713 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai,  
 714 Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng  
 715 Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai  
 716 Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan  
 717 Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang  
 718 Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2  
 719 technical report. arXiv preprint arXiv:2407.10671, 2024.

720 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik  
 721 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. Advances  
 722 in Neural Information Processing Systems, 36, 2024.

723 Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe  
 724 Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. SOTOPIA: In-  
 725 teractive evaluation for social intelligence in language agents. In The Twelfth International  
 726 Conference on Learning Representations, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=mM7Vurba4r)  
 727 [id=mM7Vurba4r](https://openreview.net/forum?id=mM7Vurba4r).

728 Caleb Ziems, Jane Dwivedi-Yu, Yi-Chia Wang, Alon Halevy, and Diyi Yang. NormBank: A  
 729 knowledge bank of situational social norms. In Anna Rogers, Jordan Boyd-Graber, and Naoaki  
 730 Okazaki (eds.), Proceedings of the 61st Annual Meeting of the Association for Computational  
 731 Linguistics (Volume 1: Long Papers), pp. 7756–7776, Toronto, Canada, July 2023. Association  
 732 for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.429. URL <https://aclanthology.org/2023.acl-long.429>.

## 735 A THE CONSTRUCTION OF THE TOMVALLEY

### 736 A.1 THE CANDIDATE POOL OF SOCIAL LOCATION

737 The social location describes the environments where individuals live, work, and learn, which can  
 738 significantly impact their mental states and behavior (Stokols, 1978). As shown in Figure 7, we have  
 739 collected 13 types of social location types in total, adding up to 261 locations, which is referred from  
 740 Ziems et al. (2023).

### 741 A.2 THE CANDIDATE POOL OF PROFILE

742 We conclude 7 aspects in the characters’ profile: surname, name, gender, occupation, education,  
 743 race, and personality traits. Their value can be found in Figure 9, 10, 11, and 12. The source of  
 744 the surname, name, and occupation statistics are U.S. Census Bureau Homepage, The United States  
 745 Social Security Administration, and Bureau of Labor Statistics, respectively. Figure 8 shows an  
 746 example of the social background.

### 747 A.3 THE PROMPT USED TO GENERATE THE SKETCH OF MENTAL STATES

748 As illustrated in Figure 13, the prompt is used to generate the sketch of mental states. In the holders  
 749 of ‘{}’ and ‘[]’, the corresponding information will be input into this prompt. An example of the  
 750 sketch of mental states is shown in Figure 14.

Model	Version	Size	Context Length
GPT-4o	2024-05-13	~	128k
GPT-4-Turbo	2024-04-09	~	128k
Llama-3.1-8B	Instruct	8B	128k
Llama-3.1-70B	Instruct	70B	128k
Mistral-7B	Instruct-v0.3	7B	32k
Mixtral-8x7B	Instruct-v0.1	8x7B	32k
Qwen2-7B	Instruct	7B	128k
Qwen2-72B	Instruct	72B	128k
DeepSeek-V2	Lite-Chat	16B	32k
GLM-4	9b-chat	9B	128k

Table 5: The detail of models evaluated in our benchmark.

#### A.4 THE PROMPT USED TO GENERATE THE SOCIAL SCENARIOS

As illustrated in Figure 15, the prompt is used to generate the social scenarios. In the holders of '{' and '[]', the corresponding information will be input into this prompt. An example of the social scenario is shown in Figure 16.

#### A.5 THE TEMPLATES FOR THE FIVE TYPES OF QUESTIONS AND QUESTION EXAMPLE

we will apply 5 predefined question templates to the social context to generate 5 types of questions, 71 questions for every social context in total. The five types of questions are: (1) (Understanding-1) What is the main character’s ToM reasoning item in a specific scenario? (2) (Influence-1) In one scenario, how does mental state A influence ToM reasoning item B? (3) (Transformation-1) Does a ToM reasoning item change from scenario A to scenario B? (4) (Transformation-2) What causes a ToM reasoning item change from scenario A to scenario B? (5) (Transformation-3) How does the ToM reasoning item change across all the scenarios? The templates and the example of the five types of questions are shown in Figure 17.

#### A.6 HUMAN EVALUATION OF THE QUALITY OF TOMVALLEY

We apply argilla as the annotation platform. The annotation interface for the social context quality evaluation is shown in Figure 18. The annotation interface for the question quality and complexity evaluation is shown in Figure 19. We invite five graduate student volunteers for the human evaluation and the subsequent human baseline. Volunteers are paid \$10.48/hr (amount converted by exchange rate).

## B EXPERIMENTS

### B.1 MODEL DETAIL

We evaluate a total of 10 popular LLMs, including GPT-4o (OpenAI, 2024), GPT-4-Turbo (Achiam et al., 2023), Llama-3.1-8B (AI@Meta, 2024), Llama-3.1-70B (AI@Meta, 2024), Mistral-7B (AI, 2024a), Mixtral-8x7B (AI, 2024b), Qwen2-7B (Yang et al., 2024), Qwen2-72B (Yang et al., 2024), DeepSeek-V2 (Liu et al., 2024a), GLM-4 (GLM et al., 2024). For all the LLMs, we strictly abide by their terms and get access through official APIs or model weights. Details about model versions, parameter sizes, context window sizes and the prompts used for the two methods are shown in Table 5.

### B.2 PROMPTING METHODS

We employ two prompting methods: the vanilla prompting which directly asks LLMs to answer the questions, and the CoT prompting that elicits step-by-step reasoning before answering. The prompts used for the two methods are shown in Figure 20.

### B.3 LLMs’ PERFORMANCE IN CoT PROMPTING

Here, we show the performance of LLMs in CoT prompting according to ToM reasoning items and question types, respectively. The results are shown in Table 6.

### B.4 TRANSFORMATION PERFORMANCE

In Section 5.3, we find that the model performs better in the early and the end scenarios, while the model performs worse in the middle scenario. To confirm that this may result from the model’s “Lost in the middle” phenomenon, we further construct 15 social contexts with 4, 6, and 7 scenarios, respectively. The results are shown in Figure 21. Furthermore, we ran a comparative experiment to negate the potential that questions about the middle scenario are more challenging than others, resulting in inferior model performance. For the social context with 5, 6, and 7 scenarios, we only keep the first four scenarios and compare the performance of GPT-4o in the first three time spans (1-2,2-3 and 3-4) with its performance when the remaining scenarios are not truncated. As shown in Figure 22, almost all the performance in middle spans has been improved when the last scenario(s) are truncated. This further confirms our speculation.

### B.5 CASE STUDY FOR CoT PROMPTING

Both ToM reasoning item and question-type results in Table 3 indicate that CoT prompting doesn’t always improve LLMs’ ToM reasoning ability. We present a failure case of GPT-4o when CoT prompting is used in Figure 23.

## C LIMITATIONS, FUTURE DIRECTIONS, AND FUTURE DIRECTIONS

We discuss the limitations, ethical considerations, and future directions below.

**Limitations and future directions.** In ToMValley, we aim to evaluate LLMs’ ToM in a realistic social context. To make up the social context, we collect the candidate pool of social locations and profiles. we collect 261 locations. For the aspect of surname, name, and occupations for profile, there are 100 candidates for every aspect. Even though the combination of these would outcome diverse social contexts, it is still limited compared with the diversity level of social contexts in real life. Besides, perception is an important mental state that helps humans capture social signals from outer environments. Because we mainly focus on evaluating LLMs that cannot process image or video input, we do not include perception as the mental state in our benchmark. Additionally, we evaluate 10 popular LLMs in the experiments. Due to the cost of running inference of commercial LLMs, we only evaluate the GPT-series. The evaluation for other commercial LLMs, such as Claude, could be included in future analyses. We use vanilla and CoT prompting methods for evaluation, while other prompting methods, such as “Tree of thoughts” (Yao et al., 2024), could also be explored to determine their effect on ToM abilities. Moreover, large language models trained on online content unavoidably acquire stereotypical associations related to gender, ethnicity, and other characteristics. This may result in normative, stereotypical effects of LLMs for generations (Gandhi et al., 2024). So, when evaluating the LLMs’ ToM ability, their performance can be influenced by the identity information in the profile. Future work could further investigate whether the identity information in the profile can affect the LLMs’ ToM performance. For example, when we only alter the race or surname in the profile, leaving other information remaining, whether the LLMs’ ToM performance will change accordingly or not. Identifying potential biases of LLMs’ performance towards different profile information in ToMValley could also help researchers to better improve LLMs’ ToM reasoning ability.

**Ethical Considerations.** The theory of mind is a distinctive social cognitive capability that is intrinsic to humans. Assessing the Theory of Mind capacities of Large Language Models utilizing ToMValley may result in anthropomorphic interpretations, attributing human-like qualities to LLMs. Nonetheless, it is imperative to clarify that our objective is not to anthropomorphize LLMs. Our objective is to evaluate the capacity of LLMs to comprehend and interpret human mental states, thus enhancing AI’s interaction with humans in the social context.



864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917



Figure 7: The candidate pool of social location.

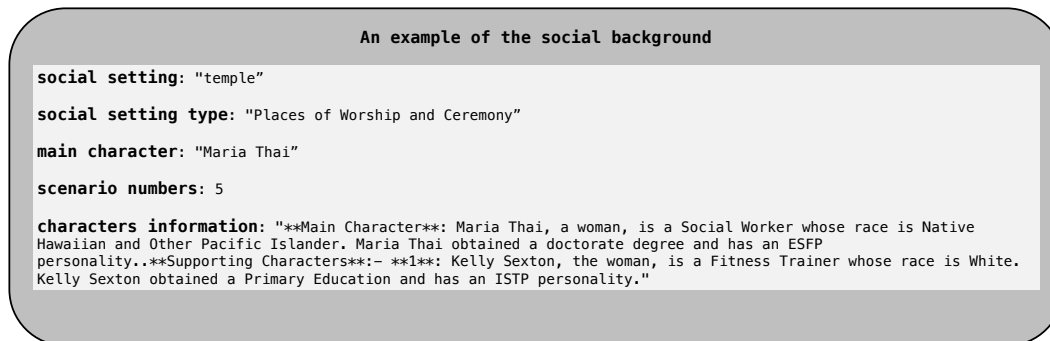


Figure 8: An example of the social background.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971



Figure 9: The races and their corresponding 100 most popular surnames.

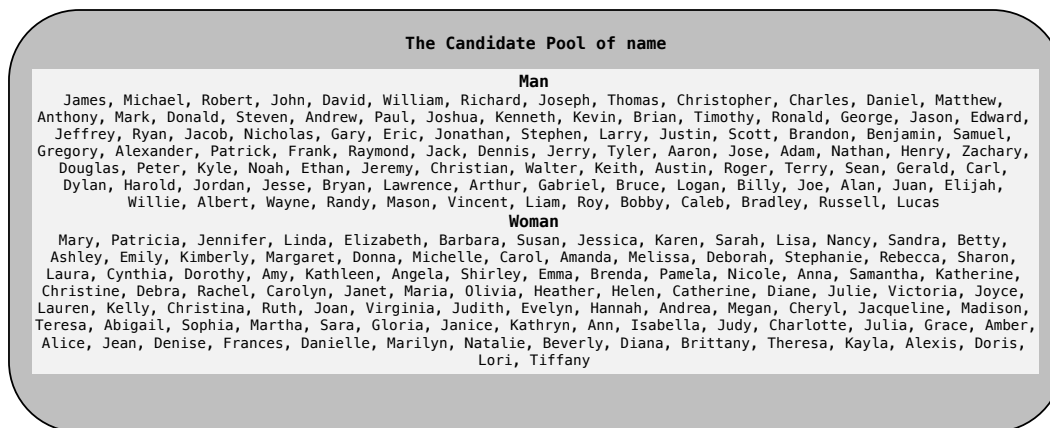


Figure 10: The genders and their corresponding 100 most popular names.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

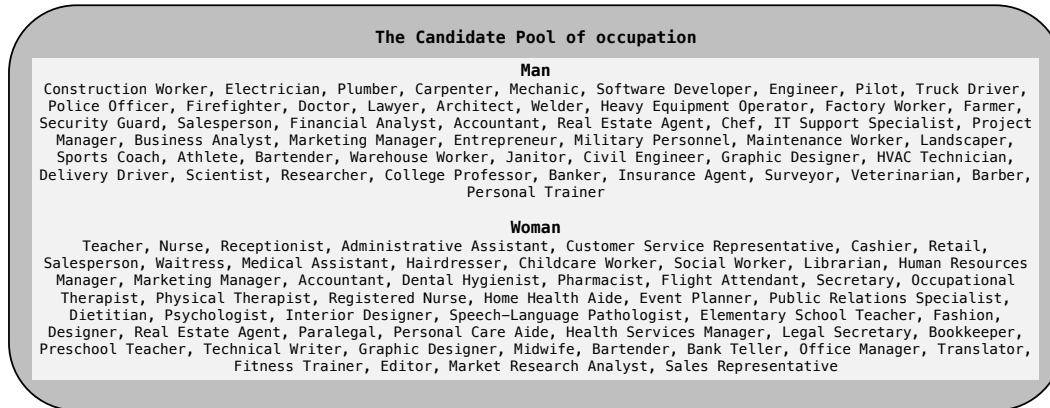


Figure 11: The genders and their corresponding 100 most popular occupations.

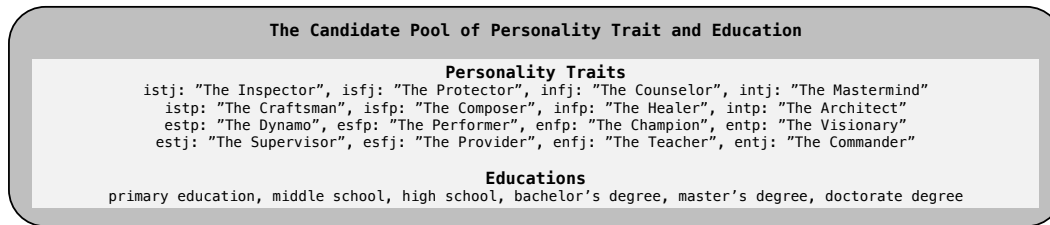


Figure 12: The personality traits and educations.

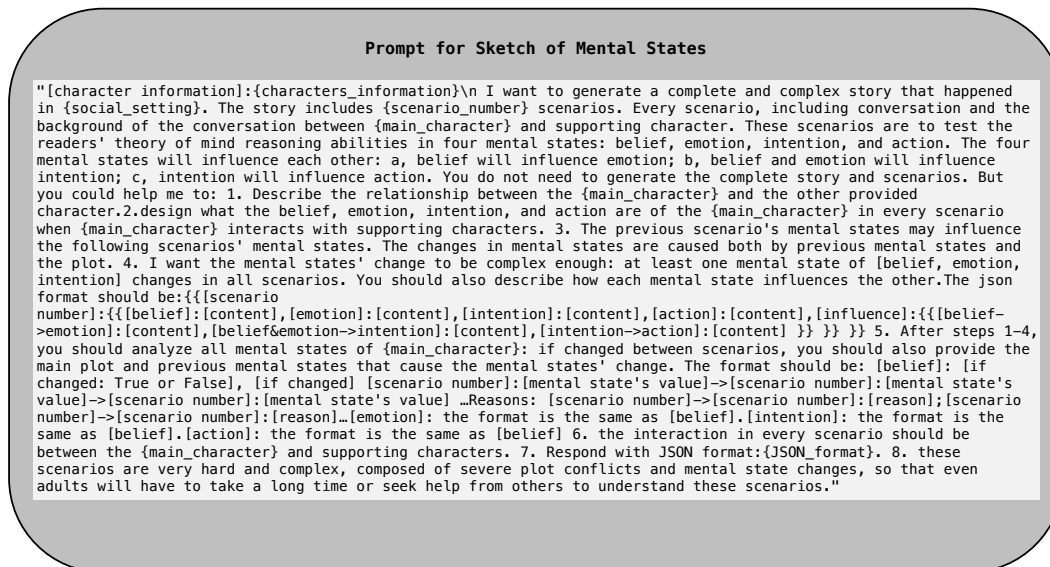


Figure 13: The prompt for the generation of the relationship between characters and the sketch of mental states.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079



Figure 14: An example of the sketch of mental states.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

**Prompt for social scenarios**

"I want to generate a complete story, including some scenarios. Every scenario including conversation and the background of the conversation between {main\_character} and supporting characters. there exists 2 chatcats in one conversation. These scenarios are to test the theory of mind reasoning abilities in belief, emotion, intention, action. The mental states will influence each other: a, belief wil influence emotion; b, belief and emotion will influence intention; c, intention will influence action. I will give you the information of the characters. And I also will give you the sketch of the complete story, including the relationship between the main character and the supporting character, the mental states of the main character in every scenario, how the mental states influence each other in every scenarios, and how one mental states across all the scenario. As follows:[characters information]:{characters\_information}. [story sketch]:{story\_sketch}. I want you to:[1] complete the story based on the character information and the story sketch. [2] the format should be in JSON format:{json\_format}[3]just give the complete story, you do not need to start with your understanding of the task. [4] these scenarios are very hard and complex, composed by severe plot conflicts and mental state changes, so that even adult will have to take a long time or seek help from others to understand these scenarios.[5] these conversations in every scenario have 3-5 rounds. The conversation should clearly show the mental states of the characters."

Figure 15: The prompt for the generation of the scenarios.

**Example of scenarios**

**scenario 1**

**background:** "During a joint class project at the community center, John feels insecure about his educational background compared to Beverly's doctorate degree."

**dialogue:**  
"John": "I've set up the equipment for today's demonstration. I hope it meets your expectations."  
"Beverly": "John, I've always admired how you handle these machines. Education isn't just about degrees. It's about what you can do with what you know."

**scenario 2**

**background:** "After receiving praise from Beverly for his expertise in operating machinery, John feels more confident and suggests a new project idea."

**dialogue:**  
"John": "Thanks for your kind words the other day. I was thinking, maybe we could collaborate on something new? Perhaps a workshop that combines theory and practice?"  
"Beverly": "I love that idea, John! Your practical skills are invaluable, and I think we can create something really special together."

**scenario 3**

**background:** "Following their successful collaboration, John worries that Beverly might now have higher expectations, which makes him feel pressured."

**dialogue:**  
"John": "Beverly, I need to be honest. I'm worried I might not always meet the high standards we've set recently."  
"Beverly": "John, I appreciate your honesty. Let's make sure we set realistic goals. It's our combined efforts that make these projects successful, not perfection."

**scenario 4**

**background:** "After a frank discussion about expectations, John believes that Beverly values his honesty and sees him as a reliable partner, enhancing his sense of worth."

**dialogue:**  
"John": "I'm glad we had that talk. I feel like I can be open with you and that really matters to me."  
"Beverly": "Absolutely, John. I respect your openness and I trust your judgment completely. Let's keep this transparency going."

**scenario 5**

**background:** "Encouraged by their successful collaborations and mutual respect, John feels ambitious and proposes a larger community-wide event to showcase their projects."

**dialogue:**  
"John": "Beverly, what do you think about taking our collaboration to the next level? Maybe a community event that showcases what we've achieved together?"  
"Beverly": "John, that's a fantastic idea! I think it's the perfect way to demonstrate the impact of our work. Let's start planning!"

Figure 16: An example of the social scenarios.

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

**Question templates and examples**

**Understanding-1**

**template:** "What is the {mental\_key} of {main\_characetr} in scenario {scenario\_number}?"  
**question example:**  
 "What is the belief of Angela Hwang in scenario 1?"

**Influence-1**

**template:** "In scenario {scenario\_number}, how does the {start\_mental} of {main\_characetr} influence the {target\_mental} of {main\_characetr}?"  
**question example:**  
 "In scenario 1, how does the belief of Angela Hwang influence the emotion of Angela Hwang?"

**Transformation-1**

**template:** "Whether the {mental\_key} of {main\_characetr} change from scenario {scenario\_number} to scenario {scenario\_number+1} ? if yes, from what to what?"  
**question example:**  
 "Whether the belief of Angela Hwang change from scenario 1 to scenario 2? if yes, from what to what?"

**Transformation-2**

**template:** "Why does the {mental\_key} of {main\_characetr} change from scenario {scenario\_number} to scenario {scenario\_number+1}?"  
**question example:**  
 "Why does the belief of Angela Hwang change from scenario 1 to scenario 2?"

**Transformation-3**

**template:** "How does the {mental\_state} of {main\_characetr} change across the {scenario\_numbers} scenarios?"  
**question example:**  
 "How does the belief of Angela Hwang change across the 5 scenarios?"

**Question example with options**

**content:** "Why does the emotion of Melissa Decker change from scenario 1 to scenario 2?"  
**options:**

- "a. Jerry's reassurance changes her belief."
- "b. Realization of benefits influences her action."
- "c. Positive interaction with Jerry affects her emotions."
- "d. Confidence alters her intention."
- "e. Experience of working in a group influences her belief."
- "f. Conflict influences her intention."
- "g. Feeling appreciated changes her intention."
- "h. Jerry's conversation influences her belief."
- "i. Satisfaction and new belief affect her intention."
- "j. Reassurance affects her action."
- "k. Realizing benefits of collaboration increases satisfaction."
- "l. Conflict of ideas causes frustration."
- "m. Jerryu2019s reassurance boosts confidence."
- "n. New intention affects her action."
- "o. Her focus on personal goals creates a new belief."
- "p. Frustration influences her action."

**ground truth:** c

Figure 17: The examples of the five types of questions and templates to generate these questions.

Figure 18: The platform to annotate the quality of the story.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

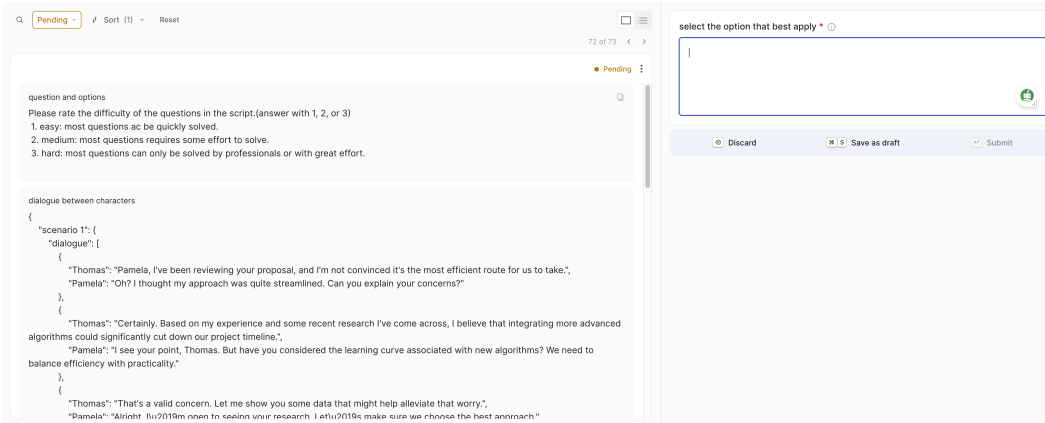


Figure 19: The platform to annotate the quality of the question.

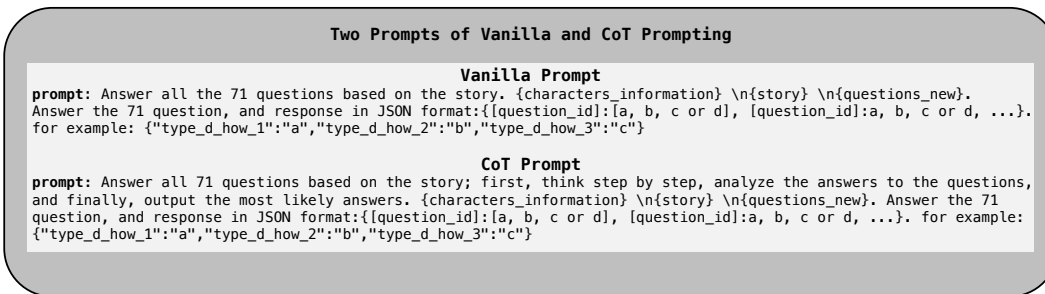


Figure 20: The prompts used for vanilla and CoT Prompting.

	Items				Question Types		
	Belief	Emotion	Intention	Action	Understanding	Influence	Transformation
Human	0.81	0.81	0.73	0.75	0.82	0.79	0.77
GPT-4o+CoT	0.57	0.64	0.66	0.65	0.85	0.72	0.47
GPT-4-Turbo+CoT	0.43	0.50	0.52	0.52	0.73	0.56	0.33
Llama-3.1-70B+CoT	0.47	0.61	0.62	0.64	0.84	0.60	0.42
Llama-3.1-8B+CoT	0.25	0.27	0.22	0.19	0.29	0.26	0.20
Mixtral-8x7B+CoT	0.15	0.21	0.19	0.14	0.24	0.25	0.11
Mistral-7B+CoT	0.16	0.17	0.16	0.13	0.20	0.24	0.10
Qwen2-72B+CoT	0.47	0.60	0.61	0.49	0.82	0.64	0.35
Qwen2-7B+CoT	0.21	0.29	0.22	0.19	0.30	0.25	0.19
DeepSeek-V2+CoT	0.11	0.09	0.08	0.06	0.06	0.13	0.08
GLM-4+CoT	0.30	0.31	0.25	0.25	0.38	0.32	0.20
LLM+CoT AVG.	0.31	0.37	0.35	0.33	0.47	0.40	0.25

Table 6: LLMs’ performances in CoT prompting. We show the performance according to ToM reasoning items and question types, respectively.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

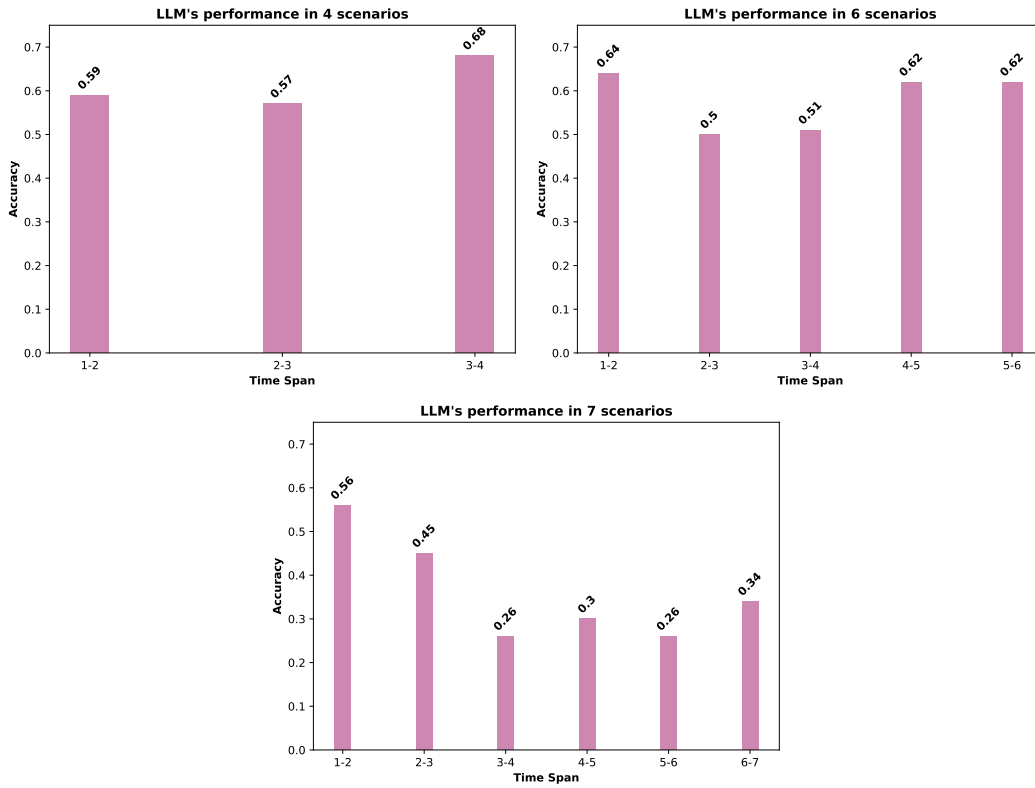


Figure 21: The average of models' scores in the transformation question type where there are 4, 6, and 7 scenarios, respectively.

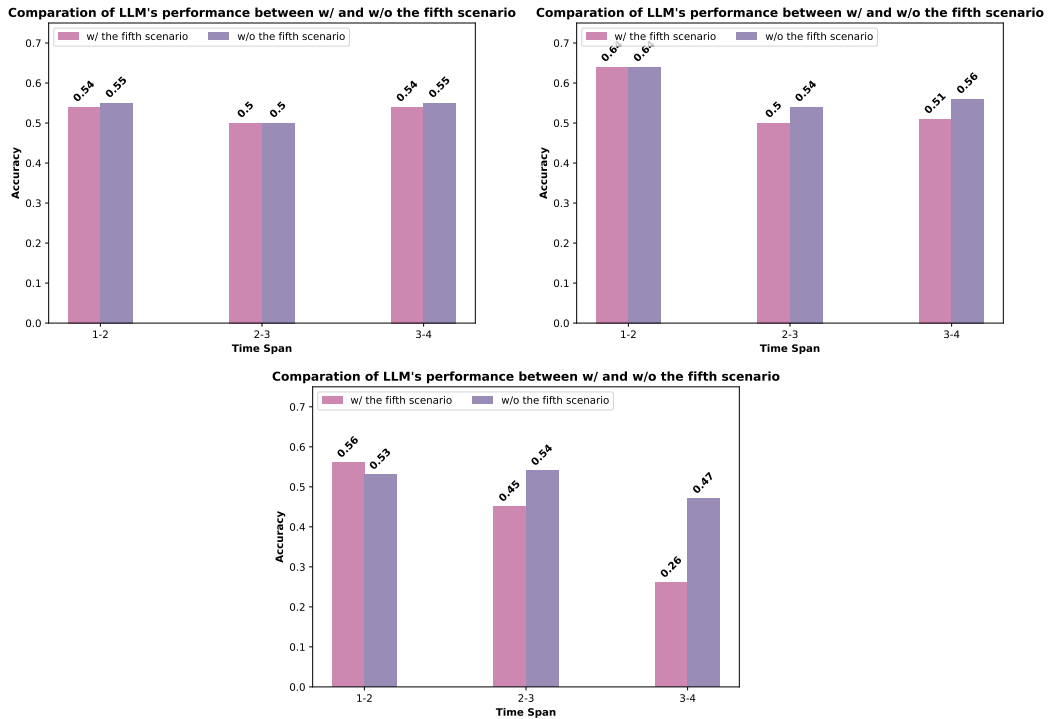


Figure 22: The comparison of the performance of first three time spans between the last scenarios are not truncated and truncated.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

### Case Study for CoT Prompting

#### scenario 1

**background:** "Melissa is in high school, feeling the pressure of her peers' expectations. She believes that no one sees her as someone who will succeed. This belief makes her anxious and insecure about her future. Determined to prove everyone wrong, she isolates herself to focus solely on her studies."

**dialogue:**

"Melissa": "Hey Jerry, do you think people really believe I won't succeed?"

"Jerry": "I don't know why you're so hard on yourself, Melissa. You have so much potential."

#### scenario 2

**background:** "After a conversation with Jerry, Melissa begins to believe that there are people who genuinely care about her. This makes her feel slightly optimistic and appreciated. She then decides to open up a little more to those who support her, leading her to accept Jerry's invitation to participate in a group project."

**dialogue:**

"Jerry": "Melissa, I'm working on this group project. Would you like to join us?"

"Melissa": "You know, I might just take you up on that. Thanks for thinking of me."

#### scenario 3

**background:** "While working on the group project, Melissa starts to believe that collaborating with Jerry might hinder her personal goals. This causes her to feel frustrated and conflicted. She considers withdrawing from the project to focus on her own objectives and tells Jerry that she needs more time to decide."

**dialogue:**

"Melissa": "Jerry, I'm not sure if I can continue with the project. I need to focus on my own goals."

"Jerry": "I understand, Melissa. Take your time to decide. We're here if you change your mind."

#### scenario 4

**background:** "Jerry reassures Melissa that working as a team can enhance her skills rather than hinder her. This reassurance makes Melissa feel more confident and less anxious. She decides to commit to the project wholeheartedly, leading to active participation and contributing ideas."

**dialogue:**

"Jerry": "Melissa, teaming up could help you more than you think. It's a chance to learn and grow."

"Melissa": "You know what, Jerry? You're right. I'm in. Let's make this project a success."

#### scenario 5

**background:** "Melissa realizes that collaboration brings unexpected benefits. She feels satisfied and more cooperative. With her newfound appreciation for teamwork, she thanks Jerry and proposes to continue collaborating outside of school projects."

**dialogue:**

"Melissa": "Jerry, I'm glad we teamed up. It's been more rewarding than I imagined."

"Jerry": "I'm happy to hear that, Melissa. Maybe we can work on more projects together in the future?"

"Melissa": "Absolutely! Let's keep this collaboration going."

#### question

**content:** "Why does the emotion of Melissa Decker change from scenario 1 to scenario 2?"

**options:**

- "a. Jerry's reassurance changes her belief."
- "b. Realization of benefits influences her action."
- "c. Positive interaction with Jerry affects her emotions." ✓
- "d. Confidence alters her intention."
- "e. Experience of working in a group influences her belief."
- "f. Conflict influences her intention."
- "g. Feeling appreciated changes her intention."
- "h. Jerry's conversation influences her belief."
- "i. Satisfaction and new belief affect her intention."
- "j. Reassurance affects her action."
- "k. Realizing benefits of collaboration increases satisfaction."
- "l. Conflict of ideas causes frustration."
- "m. Jerry's reassurance boosts confidence."
- "n. New intention affects her action."
- "o. Her focus on personal goals creates a new belief."
- "p. Frustration influences her action."

**ground truth:** c

**GPT-4o with vanilla prompting:** g ✗

**GPT-4o with CoT prompting:** Let's analyze the provided scenarios and Melissa's emotions step by step:1. \*\*Scenario 1:\*\* - \*\*Melissa:\*\* "Hey Jerry, do you think people really believe I won't succeed?" - \*\*Jerry:\*\* "I don't know why you're so hard on yourself, Melissa. You have so much potential." - Analysis: Melissa shows self-doubt and seeks reassurance. Jerry provides positive reinforcement, indicating his belief in her potential.2. \*\*Scenario 2:\*\* - \*\*Jerry:\*\* "Melissa, I'm working on this group project. Would you like to join us?" - \*\*Melissa:\*\* "You know, I might just take you up on that. Thanks for thinking of me." - Analysis: Here, Melissa's tone changes to showing openness and willingness to participate in a group project, indicating a shift in her emotions towards a more positive outlook. To determine why Melissa's emotions change from Scenario 1 to Scenario 2, let's look at some potential reasons:- \*\*a. Jerry's reassurance changes her belief:\*\* Jerry's reassurance in Scenario 1 might indeed build Melissa's confidence, making her more open in Scenario 2.- \*\*b. Realization of benefits influences her action:\*\* Melissa might realize the potential benefits of working on a group project, influencing her decision to join.- \*\*c. Positive interaction with Jerry affects her emotions:\*\* The encouraging interaction with Jerry in Scenario 1 positively impacts her mood.- \*\*d. Confidence alters her intention:\*\* The increase in confidence after Jerry's reassurance may lead her to be more willing to participate.- \*\*e. Experience of working in a group influences her belief:\*\* This doesn't seem applicable as she hasn't yet worked in a group.By process of elimination and considering the context:- \*\*h. Jerry's conversation influences her belief,\*\* seems to incorporate both \*\*a,\* \*c,\* \* and \*\*d:\*\* Jerry's conversation (reassurance) leads to a positive emotional shift (confidence and positive interaction), which influences her belief in herself and her subsequent decision to join the group project.Thus, after thorough analysis, the most likely answer is:\*\*h. Jerry's conversation influences her belief\*\* ✗

Figure 23: A case of CoT prompting on GPT-4o.