

# LLMs in front of the Supreme Court: Can LLMs understand and reason Justice and Advocate?

Anonymous ACL submission

## Abstract

While Large Language Models (LLMs) excel at legal document retrieval and case summarization, their strategic theory of mind (ToM) reasoning in legal tasks remains largely uninvestigated. This paper examines LLMs of different roles in U.S. Supreme Court oral arguments, where *Justices* challenge *Advocates* through adversarial and strategically framed questioning before making final decisions. We introduce LegalToM, the first benchmark designed to evaluate LLMs ToM reasoning in legal contexts, using U.S. Supreme Court oral arguments as a natural testbed. In LegalToM, a *Justice agent* predicts future lines of inquiry and produces legally coherent, strategically structured questions. An *Advocate agent* interprets judicial intentions and generates persuasive responses that adapt to the evolving context. We evaluate them in four tasks: understanding intentions of justices, predicting subsequent questions from justices, and generating contextually adaptive responses for both justices and advocates. Evaluations of a diverse range of open-source and proprietary LLMs show that while larger models and extended context windows yield consistent improvements, even the strongest systems fall short of expert human performance. LegalToM provides a rigorous testbed for exploring the emergence and boundaries of multi-agent ToM reasoning and human-AI interactions in professional legal processes. Our project is available at [https://anonymous.4open.science/r/oral\\_argument-48DD/](https://anonymous.4open.science/r/oral_argument-48DD/).

## 1 Introduction

Large language models (LLMs) are increasingly applied in the legal domain for contract review, document drafting, and legal retrieval. However, their capability to engage in complex legal reasoning dialogues, such as U.S. Supreme Court oral arguments (Wolfson, 2001; Johnson et al., 2006), remains largely unexplored. In these complex legal

settings, justices and advocates interact dynamically to clarify legal issues and challenge one another’s reasoning, making it an ideal testbed for evaluating the limits of LLM reasoning.

Supreme Court oral arguments pose unique challenges for LLMs. Questions from the Justices are strategically crafted and often contain implicit intentions or hypothetical traps. Moreover, the dialogues are long, adversarial, and highly context-dependent, requiring reasoning about the beliefs and goals of other participants. These properties go beyond surface-level text understanding and naturally align with the core elements of Theory of Mind (ToM) (Premack and Woodruff, 1978), including the ability to infer others’ hidden intentions, predict their future actions, and generate contextually adaptive responses.

To systematically study these challenges, we introduce LegalToM, a new benchmark that formalizes Supreme Court oral arguments as a two-agent environment. As shown in Figure 1, the advocate agent must infer a justice’s hidden intent and produce persuasive responses, while the justice agent predicts subsequent questions and generates coherent, strategically framed prompts. This formulation captures both legal reasoning and ToM perspective taking, offering a realistic multi-agent testbed for LLMs.

LegalToM is constructed from real Supreme Court oral argument transcripts (Oyez, Inc., 2025), comprising 31 cases and 3,884 evaluation samples annotated across three difficulty levels and contextual settings. We evaluate a diverse set of open-source models, including Qwen2.5\_1.5B, Qwen2.5\_7B, Qwen2.5\_70B (Yang et al., 2024), Qwen3\_4B, Qwen3\_30B (Yang et al., 2025), LLaMA3.3\_70B (Meta, 2024), Gemma3\_4B (Team et al., 2025), Phi-4\_14B (Abdin et al., 2024), SaulLM-7B (Colombo et al., 2024), as well as proprietary models such as GPT-4o (Hurst et al., 2024) and OpenAI’s reasoning

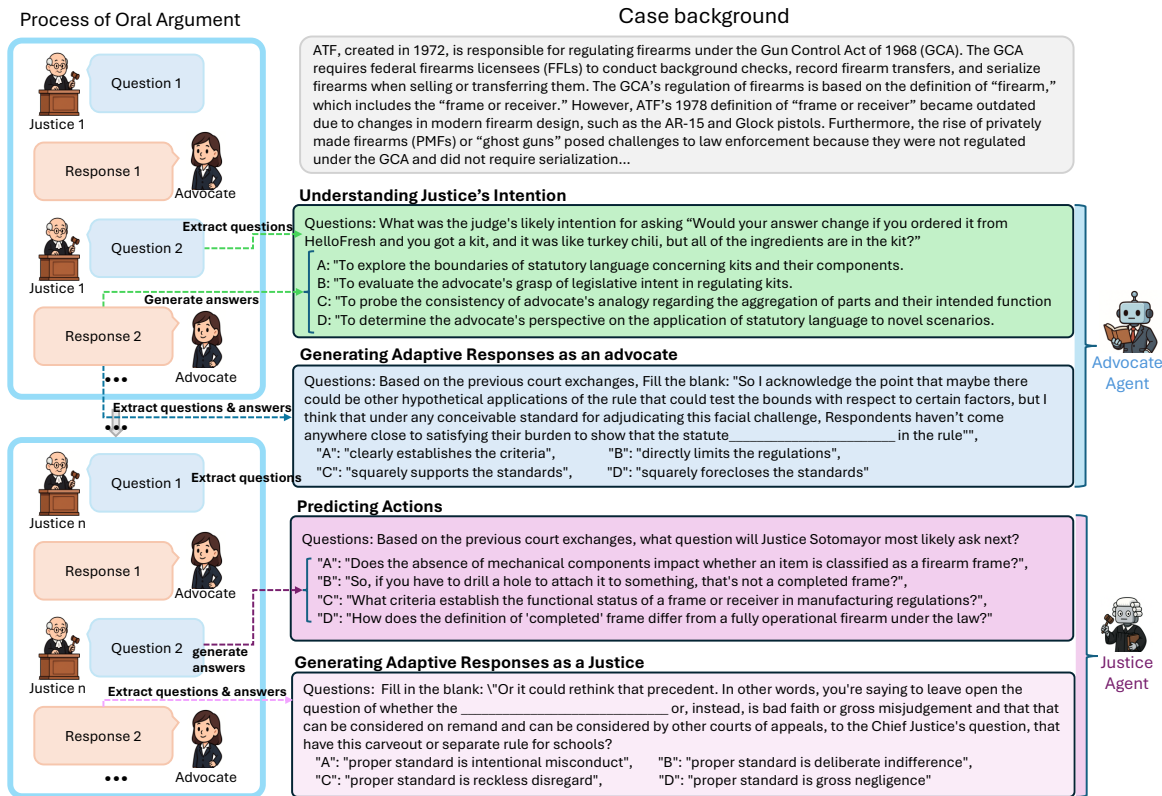


Figure 1: In Oral Arguments, Justices pose strategically framed and challenging questions. Advocates are required to infer the motivation underlying a Justice’s question and to construct a reasoned response in defense of their legal position.

084 model o3.

085 Our evaluation yields several key insights: (1)  
086 Scaling Performance: While performance scales  
087 with model size, with larger models capturing nu-  
088 anced interactions in multi-turn arguments more  
089 effectively than smaller models, all architectures  
090 struggle with evolving long-range coherence. (2)  
091 Contextual Dynamics: Although richer conversa-  
092 tional context generally improves results, track-  
093 ing shifts in legal dialogues remains a significant  
094 challenge. (3) Reasoning Capabilities: Reasoning-  
095 oriented models, such as o3, demonstrate superior  
096 proficiency in inferring the intentions of justices  
097 and advocates, whereas non-reasoning models fre-  
098 quently overlook underlying beliefs. (4) Inherent  
099 Limitations: The fact that shifts in a justice’s in-  
100 tention within extended passages continue to mislead  
101 predictions highlights a fundamental deficit in the  
102 strategic and social reasoning capabilities of cur-  
103 rent LLMs.

104 Our contributions are threefold:

- 105 • We formalize Supreme Court oral arguments
- 106 as a multi-agent legal reasoning environment.
- 107 Our environment captures the strategic and
- 108 adversarial nature of real legal interactions.

- 109 • We define four ToM-inspired tasks to evalu-  
110 ate LLMs’ capacity for perspective-taking and  
111 strategic reasoning, including intention infer-  
112 ence, question prediction, adaptive response  
113 generation for justices and advocates.
- 114 • We present LegalToM, a benchmark built from  
115 authentic U.S. Supreme Court oral arguments,  
116 and conduct a comprehensive analysis across  
117 eleven LLMs, revealing persistent limitations  
118 in ToM-like reasoning.

## 2 Related Work 119

120 **LLM about Legal Problems** Preliminary re-  
121 search has explored LLMs understanding and  
122 reasoning of legal concepts. Benchmarks such  
123 as LegalBench (Guha et al., 2023), CaseHOLD  
124 (Arvin, 2025), LexGLUE (Chalkidis et al., 2021),  
125 LEXTREME (Niklaus et al., 2023), LexEval (Li  
126 et al., 2024), and SafeLawBench (Cao et al., 2025)  
127 evaluate LLMs across various dimensions, includ-  
128 ing legal knowledge, reasoning, ethics, multilin-  
129 guality, and safety. Specialized legal LLMs such as  
130 LEGAL-BERT (Chalkidis et al., 2020), ChatLaw  
131 (Cui et al., 2023), LawLLM (Shu et al., 2024), and

SaulLM-7B (Colombo et al., 2024) have further advanced domain-specific understanding. Beyond single-agent reasoning, multi-agent legal frameworks such as Debate-Feedback (Chen et al., 2025) and Reflective (Bo et al., 2024) enable agents to provide diverse perspectives through iterative debate and reflection, thereby improving legal argumentation. However, such multi-agent environments do not fully reflect the dynamic, context-dependent, and adversarial nature of real legal interactions, making it difficult to comprehensively evaluate models’ abilities.

**Benchmarking LLMs’ ToM Reasoning** Accurately evaluating the ToM capabilities of LLMs requires benchmarks that go beyond simple belief or intention inference. Early datasets such as ToMi (Le et al., 2019) and ToM-b (Nematzadeh et al., 2018) assess reasoning about beliefs and intentions in multi-character narratives, while HiToM (He et al., 2023) introduces higher-order ToM reasoning. Recent studies (Kosinski, 2023, 2024) suggest that models like GPT-4 exhibit emerging ToM-like behaviors comparable to young children, though their generalization to realistic, strategic, and socially complex settings remains limited. Subsequent benchmarks (ToMChallenges (Ma et al., 2023), BigToM (Gandhi et al., 2023), OpenToM (Xu et al., 2024), NegotiationToM (Chan et al., 2024), ToMBench (Chen et al., 2024), and EPIT-OME (Jones et al., 2023)) introduce richer and more ecologically valid ToM evaluations, incorporating causal templates, negotiation scenarios, and automated grading systems. However, these benchmarks still rely on short, synthetic narratives and fail to capture the strategic dynamics, contextual dependencies, and adversarial reasoning found in real-world social interactions.

### 3 Oral Argument

#### 3.1 Background of Oral Argument

U.S. Supreme Court oral arguments (Bright, 1986; Wolfson, 2001) play a critical role in the judicial process. As the highest judicial authority, the Court resolves cases with significant legal and societal implications, and its rulings establish binding precedent for lower courts, shaping the interpretation and application of law nationwide. These oral arguments provide a uniquely challenging setting for evaluating LLMs’ advanced Theory-of-Mind reasoning capabilities. During arguments, justices carefully craft their questions, while advocates

must respond strategically to defend their positions. This dynamic interaction offers an ideal environment to test whether LLMs can uncover latent intentions, predict cognitive state progression, and generate contextually adaptive responses.

U.S. Supreme Court oral arguments follow a structured procedure, where justices pose questions and advocates provide responses. The session usually opens with a statement from the Chief Justice or the presiding justice, formally marking the start of the oral argument. Counsel for the applicants begins by presenting their legal position and supporting rationale. The justices then pose questions intended to probe the arguments, evidence, and legal foundations advanced by counsel. The respondents’ counsel subsequently delivers their arguments, which are similarly examined through questions posed by the justices. The advocate must reason about the hidden intention behind the justice’s questions and construct a response that defends their position. Figure 1 illustrates this interaction flow. Motivated by this structure, we adopt a ToM formulation to assess whether LLMs can reason about a justice’s latent intention, anticipate subsequent actions, and generate contextually appropriate responses in adversarial legal dialogues.

#### 3.2 Three Core Challenges of Oral Argument

**Understanding Hidden Intentions** LLMs are typically used to process the explicit meaning of a sentence and generate a relevant factual response. However, ToM requires the model to go further; it must reason about the speaker’s beliefs, intentions, and knowledge, none of which are explicitly stated in the utterance. In Supreme Court oral arguments, justices rarely ask about intentions explicitly; instead, they pose strategically challenging questions or hypotheticals to test whether an advocate’s legal argument remains valid under difficult scenarios or could lead to problematic precedents. Consequently, LLMs must infer the justice’s underlying intention rather than respond simply to the literal form of the question.

**Predicting Actions** To evaluate whether LLMs can function as a professional justice agent, we examine their ability to generate legally precise and strategically probing questions. The LLM must not only formulate questions that test the advocate agent’s legal reasoning, but also anticipate how its prompts will be interpreted and responded to. This requires reasoning about the advocate’s be-

liefs, intentions, and potential follow-up arguments. By simulating a justice agent in this way, LLMs demonstrate the capacity to perform contextually adaptive questioning, reflecting professional judgment and advanced ToM reasoning ability.

**Generating Adaptive Responses** The ultimate test of ToM capability lies in generating persuasive and strategically sound responses in high-stakes interactions. We frame the test of LLMs’ ToM capabilities in Supreme Court oral arguments in terms of two agent roles: the advocate agent and the justice agent. The advocate agent must produce legally sound, logically consistent, and strategically effective responses, while the justice agent generates coherent and targeted questions that anticipate the advocate’s likely replies. The justice agent generates coherent and targeted questions that anticipate the advocate’s likely replies.

## 4 Methodology of LegalToM

### 4.1 Problem Formulation

We categorize the challenges of oral arguments into three core tasks: understanding intentions, predicting actions, and generating adaptive responses for both the advocate and justice agents. However, since these challenges arise dynamically over the course of long, evolving dialogues, traditional benchmarks are not well suited for their evaluation. To address this, we formalize oral arguments as a multi-agent environment, ensuring that the evolving context and interactive nature of the debate are properly captured.

We formalize U.S. Supreme Court Oral Arguments as a *multi-agent ToM environment*. At each turn  $t$ , the dialogue involves two primary roles:

- **Justice ( $J$ ):** poses strategically framed questions  $a_t^J$ , guided by hidden intentions  $\theta_t^J$  (e.g., probing the generality of a legal principle, testing consistency, and exposing weaknesses).
- **Advocate ( $A$ ):** responds in real time, aiming to clarify, defend, and persuade by generating the utterance  $a_t^A$

At each step  $t$ ,  $a_t^J \sim \pi_J(\cdot | s_t, \theta_t^J)$  and  $a_t^A \sim \pi_A(\cdot | s_t, a_t^J)$ , where  $a_t^J$  is the justice’s utterance,  $a_t^A$  is the advocate’s response,  $\pi_J$  is the justice’s questioning strategy conditioned on hidden intention  $\theta_t^J$ , and  $\pi_A$  is the advocate’s responding strategy.

The advocate does not have direct access to the justice’s latent intention  $\theta_t^J$ , and must infer it from the dialogue context and the justice’s question. Effective advocacy therefore requires reasoning the hidden intention  $\theta_t^J$ . In addition, the advocate needs to anticipate the justice’s likely next question  $a_{t+1}^J$  based on the current dialogue. Finally, the advocate must generate a response  $a_t^A$  that is both contextually appropriate and persuasive.

### 4.2 Task Design

We formalize three evaluation tasks using LLM predictions. In the **Intention Understanding** task, the model predicts the justice’s hidden intention given the dialogue context  $s_t$  and the justice’s question  $a_t^J$ :

$$\hat{\theta}_t^J = \text{LLM}(s_t, a_t^J),$$

where  $\hat{\theta}_t^J$  is the predicted intention generated by the LLM and compared against the labeled ground truth. The justice’s true intention  $\theta_t^J$  is inherently unobservable. In addition, in oral arguments, the advocate (typically the Solicitor General or other highly experienced attorneys) represents the highest level of human ability to interpret a justice’s intention. We therefore extract the justice’s intention from the advocate’s responses and use it as ground truth  $\tilde{\theta}_t^J$ .

In the **Action Prediction** task, the model anticipates the justice’s next question based on the dialogue context, the current justice question, and the advocate’s response:

$$\hat{a}_{t+1}^J = \text{LLM}(s_t, a_t^J, a_t^A),$$

where  $\hat{a}_{t+1}^J$  is the predicted next question and is evaluated against the actual next question  $a_{t+1}^J$  in the transcript.

In the **Adaptive Response Generation** task, the model produces the advocate’s response conditioned on the dialogue context and the current justice question:

$$\hat{a}_t^A = \text{LLM}(s_t, a_t^J),$$

where  $\hat{a}_t^A$  is the LLM-generated response compared to the human advocate’s reference response.

For the justice agent task, the model fills in missing portions of the dialogue to generate legally precise and coherent prompts that reflect proper legal reasoning:

$$\hat{a}_{t+1}^J = \text{LLM}(s_t, a_t^A),$$

where  $\hat{a}_{t+1}^J$  represents the LLM-generated justice prompt, formulated based on the dialogue context and the advocate’s current response.

Together, these tasks cover the full cycle of ToM reasoning in adversarial dialogue: (i) inferring hidden intentions  $\theta_t^J$ , (ii) predicting future actions  $a_{t+1}^J$ , and (iii) generating adaptive responses  $a_t^A$  and  $a_{t+1}^J$ . This design provides a structured yet realistic benchmark for probing LLMs’ multi-agent reasoning capabilities in high-stakes legal settings.

### 4.3 ToM Tasks Data Curation from Oral Arguments

Each Supreme Court oral argument case typically lasts one to two hours, resulting in long and complex transcripts that require reasoning over extended contexts. Our goal is to transform real-world oral argument data into a multi-agent setting, where LLM justice agents pose questions and the LLM advocate agent responds based on authentic transcripts. This formulation allows us to simulate realistic courtroom interactions while evaluating the model’s ability to reason, predict, and generate contextually appropriate responses. As illustrated in Figure 1, our task construction follows a structured pipeline that converts raw oral argument transcripts into task-specific multiple-choice evaluation instances. We use GPT-4o (Hurst et al., 2024) to assist with content generation during this transformation process.

**Task-Specific Data Generation** To support the construction of task-specific evaluation benchmarks, we first preprocess the raw case materials by converting case PDFs into Markdown format, enabling structured and efficient downstream processing. The detailed conversion pipeline is provided in the appendix A.1.

After preprocessing, we generate evaluation tasks targeting key challenges for LLM-based advocates: (1) inferring a justice’s hidden intention, (2) predicting the justice’s next question, and (3) producing contextually adaptive responses for both the advocate and the justice.

For **Understanding Hidden Intentions**, we construct multiple-choice questions that ask for the justice’s likely intention. Since true intentions are not directly observable, we assume that expert advocates can reliably infer them from context. Accordingly, the advocate’s inferred response is treated as the correct option, while distractors are designed to be plausible alternatives. These distractors mirror

the structure and legal framing of the correct answer but are intentionally modified to be incorrect. Concretely, a distractor may adopt a similar legal tone or argumentative role while addressing a different legal issue, misapplying a relevant doctrine, or being inconsistent with the specific context of the dialogue. To validate this assumption, we sample a subset of the data and ask human legal experts to review the LLM-extracted advocate intentions, confirming the correctness of the generated labels. All prompts used for generating questions and answers are provided in Appendix A.3.

For **Predicting Actions**, the model is required to anticipate the justice’s next question based on prior dialogue and context. The ground truth answer is the actual subsequent question from the transcript, requiring the model to emulate human-like legal reasoning. Incorrect and misleading options are generated to increase task difficulty. Details of the prompt design are included in Appendix A.3.

For **Generating Adaptive Responses**, we convert expert responses into multiple-choice items by selecting key phrases or legal terms as correct answers. We use LLMs to generate distractor options for evaluation. This tests whether models can identify critical information and produce legally precise, context-aware responses. In addition, we construct two types of response generation tasks: one for advocates and one for justices to assess whether LLMs can successfully function as either side in oral arguments. The advocate generation task tests whether the model can generate persuasive and contextually appropriate replies to justices’ questions, while the justice generation task evaluates whether the model can formulate coherent, strategically probing questions or comments that reflect judicial reasoning. The prompts used to generate the questions and corresponding answers are provided in the last part of the appendix A.3.

### 4.4 Experts Verification

To ensure the reliability of our evaluation, we consulted legal experts to verify the accuracy of LLM-generated content on a subset of the dataset, especially concerning intention extraction. The experts additionally affirmed that our tasks and questions were well-defined and meaningful within the legal context (Seeing Appendix C).

### 4.5 Evaluation

To evaluate LLMs’ abilities on Understanding Hidden Intentions and Predicting Actions, we adopt

Model Names	Understanding Intention			Predicting Actions			Generating Responses for Advocate			Generating Responses for Justice		
	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC
<b>Closed-source Models</b>												
o3(Reasoning)	<b>84.83</b>	<b>85.48</b>	82.56	<b>75.36</b>	88.40	<b>91.13</b>	<b>81.77</b>	<b>83.09</b>	<b>82.71</b>	<b>73.45</b>	<b>73.88</b>	<b>75.20</b>
Gpt-4o	82.89	80.86	81.83	60.63	70.27	88.66	69.80	70.72	72.61	57.14	56.13	62.14
<b>Open-source Models</b>												
Qwen2.5_1.5b_instruct	68.21	69.26	65.53	59.71	65.19	60.76	53.81	56.22	59.13	49.34	53.51	53.70
Qwen2.5_7b_instruct	77.45	77.21	78.02	50.20	77.44	75.88	60.10	62.34	63.59	53.77	58.16	59.30
Qwen2.5_72b_instruct	82.32	83.62	<b>85.89</b>	45.76	77.44	80.05	<b>71.86</b>	<b>73.44</b>	<b>75.31</b>	<b>63.30</b>	<b>64.23</b>	<b>66.51</b>
Qwen3_4b_instruct	73.15	73.15	75.75	42.11	66.10	64.41	56.85	58.52	62.22	47.80	51.29	52.45
Qwen3_30b_instruct	83.29	83.37	84.27	60.63	83.70	85.01	65.32	66.87	67.94	57.75	60.85	60.98
Gemma3_4b_instruct	64.23	64.23	70.80	29.34	49.28	52.02	46.91	48.86	52.28	40.89	44.75	46.94
Phi-4_14b	82.16	83.13	84.02	38.98	73.53	70.14	66.08	66.49	66.60	60.31	63.79	63.89
Llama3.3_70b_instruct	<b>84.83</b>	<b>85.00</b>	84.91	<b>64.93</b>	<b>88.79</b>	<b>88.79</b>	67.53	70.12	72.30	62.05	62.91	64.88
Saul_7b_instruct (law)	66.34	66.91	60.02	33.25	51.63	47.72	39.79	38.28	35.27	34.53	36.24	33.02

Table 1: Performance of reasoning, closed-source, and open-source models under three context conditions. **BG** provides the case background, **BG+CC** further includes the current conversation, and **BG+CC+PC** incorporates the preceding three conversations. We evaluate Understanding Intention, Action Prediction, and Response Generation using accuracy (%). **Bold-underlined** indicates the strongest among all models; **bold** indicates the strongest among non-reasoning models.

multiple-choice questions. This setup renders evaluation straightforward, as performance can be directly quantified using accuracy. Since all of our questions are in multiple-choice format, we only need to calculate the accuracy to evaluate how different models perform on our tasks.

## 5 Experiments

The evaluation comprises four tasks, each probing a distinct aspect of ToM reasoning in legal discourse. The first task focuses on reasoning about hidden intentions, requiring the model to infer the true motivations behind a justice’s questions. The second task involves predicting next actions, where the model anticipates the justice’s subsequent question based on the preceding dialogue context. The third and fourth tasks examine generating adaptive responses in which the models act as advocates and justices to produce high quality, contextually appropriate replies.

### 5.1 Experimental Setup

**Baseline Models** We evaluate GPT-o3 and GPT-4o (Hurst et al., 2024) as a representative closed-source model. GPT-o3 denotes the reasoning model. For open-source models, we assess a range of Instruct models spanning 1.5B to 72B parameters, including Qwen2.5-1.5B-Instruct, Qwen2.5-7B-Instruct, Qwen2.5-72B-Instruct (Yang et al., 2024), Qwen3-4B-Instruct, Qwen3-30B-Instruct (Yang et al., 2025), Gemma3-4B-Instruct (Team et al., 2025), Phi-4 (Abdin et al., 2024), LLaMA3.3-70B-Instruct (Dubey et al., 2024), and Saul-7B-Instruct (Colombo et al., 2024). Among these, Saul-7B-Instruct is specifically fine-tuned on legal texts,

optimizing its performance for legal terminology.

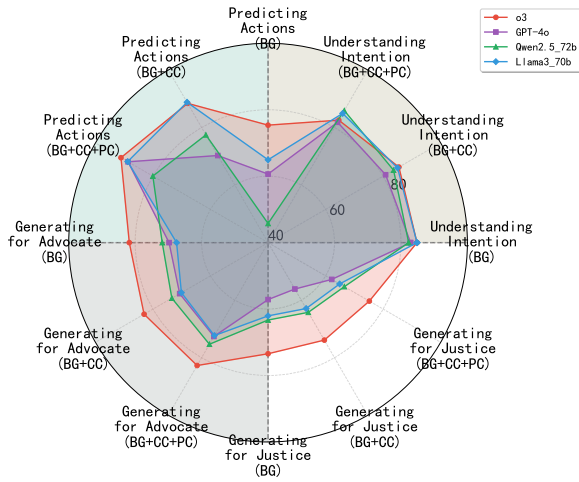
**Oral Arguments Dataset** We selected 31 oral argument cases covering a diverse range of legal scenarios, including education, criminal and procedural matters, immigration and citizenship, environmental and administrative regulation, labor and employment, commercial and civil disputes, and constitutional and public rights. These cases provide a rich testbed for evaluating both legal reasoning and the reasoning capabilities of LLMs.

From these 31 cases, we generated a dataset comprising 3,884 question-answer pairs across three tasks. The Understanding Hidden Intentions task contains 1,233 questions, Predicting Next Actions includes 767 multiple-choice questions, and Generating Adaptive Responses comprises 1,884 questions, divided between the advocate (970) and the justice (914) roles.

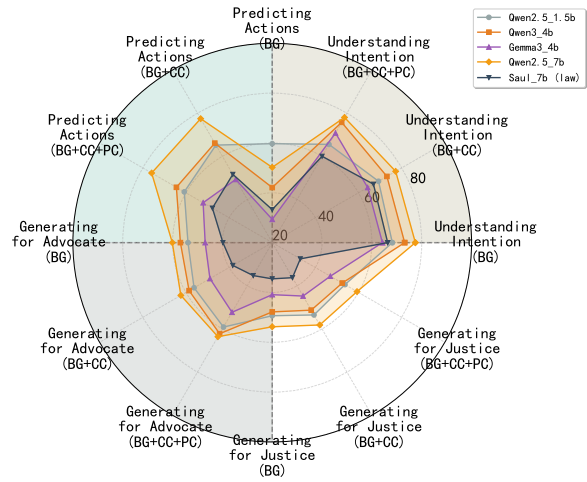
### 5.2 Experimental Results

This section presents the results of our experiments, evaluating the performance of various LLMs on four distinct tasks. To facilitate analysis, We segmented the long dialogue into a series of smaller conversations, each anchored by the names of justices, to facilitate a comparative evaluation of model performance under two distinct conversational context configurations: a limited condition providing only the immediate background and current turn, and a richer condition including three preceding conversational turns. In Table 1, we divide the context into different types: the background is denoted as **BG**, the current conversation as **CC**, and the preceding three conversations as **PC**.

As shown in the radar chart (Figure 2a), the rea-



(a) Top-tier models comparison



(b) Small models comparison

Figure 2: Performance comparison across all tasks with different contexts (BG, BG+CC, BG+CC+PC).

soning model o3 achieves the best overall performance across all four tasks. We observe that the reasoning model o3 achieves the best overall performance across all four tasks. Among the non-reasoning models, Qwen2.5-72B attains the highest scores in Understanding Intention and Generating Responses for advocate, Llama3-70B performs best in Predicting Actions, and GPT-4o achieves the top score in Generating Responses for justice.

Figure 2b shows that smaller models perform less effectively, with Saul achieving the lowest scores. Although Saul is fine-tuned on legal-domain data, its lower performance indicates that our benchmark primarily evaluates models' ToM reasoning under long-context scenarios rather than domain-specific knowledge.

### 5.3 Difficulty Classification

We categorized the tasks by difficulty (see Table 2), enabling our benchmark to better establish difficulty levels and evaluate models of varying capabilities. We categorize question difficulty based on the accuracy achieved by different models. Under the hard tasks, different models achieve accuracies ranging from 37% to 62%, providing a set of highly challenging problems for LLMs. Detailed results of three difficulties are presented in Tables 8, 9, and 10.

### 5.4 Analysis

**Conversational Context** As shown in Table 1, our experiments reveal a consistent trend: providing models with richer conversational context, including preceding turns, improves their perfor-

Task	Easy	Medium	Hard	Total
Advocate Keywords	324	323	323	970
Justice Keywords	307	305	305	917
Intention Analysis	369	554	310	1,233
Prediction Task	230	345	192	767
Total	1,230	1,527	1030	3,887

Table 2: Task difficulty statistics

mance. We observe that model performance generally improves as more context is provided, especially for larger models. However, smaller models, such as Qwen2.5-1.5B and Saul-7B, sometimes exhibit degraded performance when additional context is introduced. This decline may be due to their limited capacity to effectively encode and utilize long-range contextual information, which can lead to confusion or interference from irrelevant details, thereby reducing overall task accuracy.

**Effect of Model Size** Our analysis on the oral arguments dataset reveals that the Qwen2.5 series follows a clear scaling law. Performance consistently improves with model size. Seeing Figure 3, this strong trend indicates that larger parameter counts and greater overall capacity are crucial for capturing the complex, nuanced information required by these tasks. Notably, smaller models, such as Qwen2.5-1.5B, may struggle with longer contextual dependencies, highlighting the importance of both scale and context utilization in multi-turn reasoning.

**Reasoning Model** As a reasoning model, o3 more accurately captures judges' and advocates'

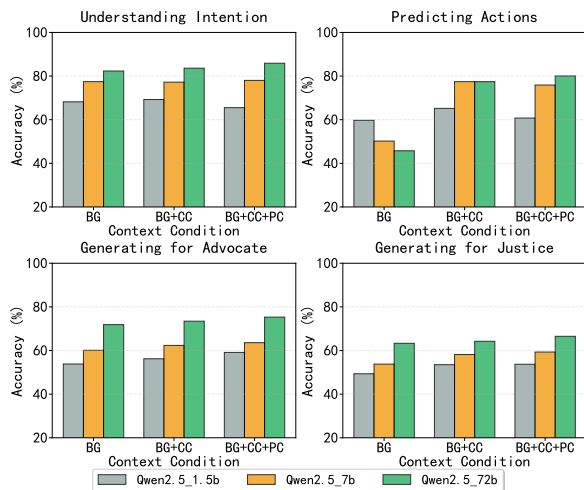


Figure 3: Qwen2.5 models follow a scaling law: performance consistently improves as model size increases from 1.5B to 72B, with larger models benefiting more from richer context.

subtle beliefs and intentions in long-context dialogues, enabling it to make more accurate predictions about their likely next moves. For example, in the item shown in Table 6, nearly all non-reasoning models failed, whereas o3 correctly inferred that *Justice Sotomayor had just pointed out that “intent is not even an issue here,”* and anticipated that her next question would refocus the argument on the matter for which the Court granted certiorari (see Table 4 for details).

In contrast, non-reasoning models such as Qwen2.5-72B did not effectively capture the justice’s beliefs regarding the previous question. Qwen2.5-72B focused *primarily on the definitions and implications of “bad faith” and “intent to discriminate,”* and on the need to further clarify how the standards for intent apply differently to various forms of relief, such as damages versus injunctive relief (see Table 5). Because it failed to account for Justice Sotomayor’s underlying concern about the specific legal question, Qwen2.5-72B produced an incorrect prediction.

**Questions All Agents Fail** In LegalToM, several highly challenging cases remained unsolved by all models, including the reasoning model o3. These failures occurred when models could not track shifts in a justice’s intention during long explanatory passages, instead over-focusing on surface details, which led to a misunderstanding of the hidden intention. As shown in Table 7, although Justice Sotomayor devoted a lengthy discussion to the relationship between *mens rea and intent*, she began

by stating, “*Motive, in — intent — we don’t care about motive,*” indicating that motive and intent were not her primary concern. Nevertheless, even the most capable o3 model was misled by her extended exposition and concluded that her aim was to “*Evaluate the advocate’s interpretation of the statutory language regarding intent requirements.*”. This example illustrates the challenges LLMs face in reasoning over long, nuanced legal discourse and accurately inferring hidden intentions.

## 6 Conclusion

We introduce a novel multi-agent framework to evaluate LLMs’ ToM reasoning in the complex, adversarial setting of U.S. Supreme Court oral arguments, formalizing three core tasks: understanding hidden intentions, predicting subsequent questions, and generating contextually adaptive responses. Using this testbed, we find that performance scales with model size, larger models better capture multi-turn interactions, and reasoning-oriented models such as o3 more effectively infer justices’ and advocates’ intentions. Nevertheless, all models struggle with evolving intentions and long-range coherence, and shifts in a justice’s intention within extended passages reveal fundamental limitations in strategic and social reasoning, highlighting challenges for developing context-aware, multi-agent legal reasoning LLMs.

## 7 Limitations

While our study demonstrates the benefits of providing richer conversational context in multi-turn legal dialogues, several limitations remain. First, our dataset, although drawn from real oral argument dialogues, is still limited and may not fully capture the diversity and complexity of real-world legal proceedings. Importantly, even though the justices and advocates are highly professional, errors can still occur during their exchanges. Since we take the justices’ and advocates’ dialogues as the ground-truth labels for our tasks, this may introduce occasional inaccuracies. Second, despite improvements in handling long contexts, models continue to face challenges in capturing deeper reasoning and nuanced ToM aspects, such as inferring judges’ beliefs, intentions, and strategic reasoning in complex interactions.

## 8 Ethical Considerations

This research complies with the ACL Ethics Policy. All data used in this study are drawn from publicly available U.S. Supreme Court oral argument transcripts. A subset of the data, including LLM-extracted intentions, was manually verified by legal experts to ensure accuracy. All data contain no private or personally identifiable information. The dialogues analyzed represent official judicial proceedings and are used solely for academic research purposes. Our work aims to advance understanding of language models’ reasoning and Theory of Mind capabilities in legal contexts, without attempting to replace or simulate actual judicial decision-making. No data manipulation, fabrication, or misrepresentation was involved. We also ensured that model outputs and analyses were reviewed for fairness and interpretability to avoid potential bias or misuse in real-world legal applications.

## References

- Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, and 1 others. 2024. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*.
- Chuck Arvin. 2025. Identifying legal holdings with llms: A systematic study of performance, scale, and memorization. *arXiv preprint arXiv:2505.02172*.
- Xiaohe Bo, Zeyu Zhang, Quanyu Dai, Xueyang Feng, Lei Wang, Rui Li, Xu Chen, and Ji-Rong Wen. 2024. Reflective multi-agent collaboration based on large language models. *Advances in Neural Information Processing Systems*, 37:138595–138631.
- Myron H Bright. 1986. The power of the spoken word: In defense of oral argument. *Iowa L. Rev.*, 72:35.
- Chuxue Cao, Han Zhu, Jiaming Ji, Qichao Sun, Zhenghao Zhu, Yinyu Wu, Juntao Dai, Yaodong Yang, Sirui Han, and Yike Guo. 2025. Safelawbench: Towards safe alignment of large language models. *arXiv preprint arXiv:2506.06636*.
- Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. Legal-bert: The muppets straight out of law school. *arXiv preprint arXiv:2010.02559*.
- Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Martin Katz, and Nikolaos Aletras. 2021. Lexglue: A benchmark dataset for legal language understanding in english. *arXiv preprint arXiv:2110.00976*.

- Chunkit Chan, Cheng Jiayang, Yauwai Yim, Zheyue Deng, Wei Fan, Haoran Li, Xin Liu, Hongming Zhang, Weiqi Wang, and Yangqiu Song. 2024. Negotiationtom: A benchmark for stress-testing machine theory of mind on negotiation surrounding. *arXiv preprint arXiv:2404.13627*.
- Xi Chen, Mao Mao, Shuo Li, and Haotian Shangguan. 2025. Debate-feedback: A multi-agent framework for efficient legal judgment prediction. *arXiv preprint arXiv:2504.05358*.
- Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen, Guanqun Bi, Gongyao Jiang, Yaru Cao, Mengting Hu, Yunghwei Lai, Zexuan Xiong, and 1 others. 2024. Tombench: Benchmarking theory of mind in large language models. *arXiv preprint arXiv:2402.15052*.
- Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, Andre FT Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia Morgado, and 1 others. 2024. Saullm-7b: A pioneering large language model for law. *arXiv preprint arXiv:2403.03883*.
- Jiayi Cui, Zongjian Li, Yang Yan, Bohua Chen, and Li Yuan. 2023. Chatlaw: Open-source legal large language model with integrated external knowledge bases. *CoRR*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407.
- Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. 2023. Understanding social reasoning in language models with language models. *Advances in Neural Information Processing Systems*, 36:13518–13529.
- Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré, Adam Chilton, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel Rockmore, Diego Zambrano, and 1 others. 2023. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. *Advances in neural information processing systems*, 36:44123–44279.
- Yinghui He, Yufan Wu, Yilin Jia, Rada Mihalcea, Yulong Chen, and Naihao Deng. 2023. Hi-tom: A benchmark for evaluating higher-order theory of mind reasoning in large language models. *arXiv preprint arXiv:2310.16755*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Timothy R Johnson, Paul J Wahlbeck, and James F Spriggs. 2006. The influence of oral arguments on the us supreme court. *American Political Science Review*, 100(1):99–113.

731	Cameron Robert Jones, Sean Trott, and Ben Bergen.	Warren D Wolfson. 2001. Oral argument: Doe it matter.	785
732	2023. Epitome: Experimental protocol inventory for	<i>Ind. L. Rev.</i> , 35:451.	786
733	theory of mind evaluation. In <i>First Workshop on</i>		
734	<i>Theory of Mind in Communicating Agents</i> .		
735	Michal Kosinski. 2023. Theory of mind may have spon-	Hainiu Xu, Runcong Zhao, Lixing Zhu, Jinhua Du, and	787
736	tananeously emerged in large language models. <i>arXiv</i>	Yulan He. 2024. Opentom: A comprehensive bench-	788
737	<i>preprint arXiv:2302.02083</i> , 4:169.	mark for evaluating theory-of-mind reasoning capa-	789
		bilities of large language models. <i>arXiv preprint</i>	790
		<i>arXiv:2402.06044</i> .	791
738	Michal Kosinski. 2024. Evaluating large language mod-	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	792
739	els in theory of mind tasks. <i>Proceedings of the Na-</i>	Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,	793
740	<i>tional Academy of Sciences</i> , 121(45):e2405460121.	Chengen Huang, Chenxu Lv, Chujie Zheng, Dayi-	794
		heng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge,	795
741	Matthew Le, Y-Lan Boureau, and Maximilian Nickel.	Haoran Wei, Huan Lin, Jialong Tang, and 41 oth-	796
742	2019. Revisiting the evaluation of theory of mind	ers. 2025. Qwen3 technical report. <i>arXiv preprint</i>	797
743	through question answering. In <i>Proceedings of the</i>	<i>arXiv:2505.09388</i> .	798
744	<i>2019 Conference on Empirical Methods in Natu-</i>		
745	<i>ral Language Processing and the 9th International</i>	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,	799
746	<i>Joint Conference on Natural Language Processing</i>	Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,	800
747	<i>(EMNLP-IJCNLP)</i> , pages 5872–5877.	Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jian-	801
		hong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang,	802
748	Haitao Li, You Chen, Qingyao Ai, Yueyue Wu, Ruizhe	Jingren Zhou, Junyang Lin, Kai Dang, and 22 oth-	803
749	Zhang, and Yiqun Liu. 2024. Lexeval: A compre-	ers. 2024. Qwen2.5 technical report. <i>arXiv preprint</i>	804
750	hensive chinese legal benchmark for evaluating large	<i>arXiv:2412.15115</i> .	805
751	language models. <i>Advances in Neural Information</i>		
752	<i>Processing Systems</i> , 37:25061–25094.		
753	Xiaomeng Ma, Lingyu Gao, and Qihui Xu. 2023. Tom-		
754	challenges: A principle-guided dataset and diverse		
755	evaluation tasks for exploring theory of mind. <i>arXiv</i>		
756	<i>preprint arXiv:2305.15068</i> .		
757	AI Meta. 2024. Llama 3.2: Revolutionizing edge ai		
758	and vision with open, customizable models. <i>Meta AI</i>		
759	<i>Blog</i> . Retrieved December, 20:2024.		
760	Aida Nematzadeh, Kaylee Burns, Erin Grant, Alison		
761	Gopnik, and Thomas L Griffiths. 2018. Evaluating		
762	theory of mind in question answering. <i>arXiv preprint</i>		
763	<i>arXiv:1808.09352</i> .		
764	Joel Niklaus, Veton Matoshi, Pooja Rani, Andrea		
765	Galassi, Matthias Stürmer, and Ilias Chalkidis.		
766	2023. Lextreme: A multi-lingual and multi-task		
767	benchmark for the legal domain. <i>arXiv preprint</i>		
768	<i>arXiv:2301.13126</i> .		
769	Oyez, Inc. 2025. <a href="#">Oyez: U.s. supreme court</a> . The Oyez		
770	Project.		
771	David Premack and Guy Woodruff. 1978. Does the		
772	chimpanzee have a theory of mind? <i>Behavioral and</i>		
773	<i>brain sciences</i> , 1(4):515–526.		
774	Dong Shu, Haoran Zhao, Xukun Liu, David Demeter,		
775	Mengnan Du, and Yongfeng Zhang. 2024. Lawllm:		
776	Law large language model for the us legal system.		
777	In <i>Proceedings of the 33rd ACM International Con-</i>		
778	<i>ference on information and knowledge management</i> ,		
779	pages 4882–4889.		
780	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya		
781	Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin,		
782	Tatiana Matejovicova, Alexandre Ramé, Morgane		
783	Rivière, and 1 others. 2025. Gemma 3 technical		
784	report. <i>arXiv preprint arXiv:2503.19786</i> .		

## A Further Details on Generating Questions

### A.1 Method for Preprocessing

We obtain the oral argument transcripts (Oyez, Inc., 2025) in PDF format and preprocess them by converting the PDFs into Markdown format. This preprocessing segments the conversations between individual justices and advocates into structured blocks, as illustrated in Listing 1.

Listing 1: Example of segmented oral argument dialogue in Markdown format

```
## Justice Thomas
### Justice Thomas:
Under your approach, what would that correct
  standard look like and how would it be
  applied here?
### Nathaniel A.G. Zelinsky:
Justice Thomas, we think the standard is the "
  totality of the circumstances" standard that
  this Court articulated in Graham and Garner
  , Scott, and Plumhoff.
...
## Justice Kavanaugh
### Justice Kavanaugh:
Was it reasonable to --for the officer to jump
  on the side of the car?
### Nathaniel A.G. Zelinsky:
So, Just -- Justice Kavanaugh, we don't think it
  was in this particular case, but that's
  precisely the issue that the lower courts
  couldn't evaluate because they applied this
  legal amnesia and only look at the fact that
  the officer was on the moving vehicle.
...
```

### A.2 Method for Generating Questions

For the first task, **Understanding Hidden Intentions**, we construct multiple-choice questions that ask for the justice’s intention behind a given question, for example: *"What was the justice’s likely intent in asking this question?"* followed by the verbatim text of the justice’s question. For generating correct and distract options in multiple-choice questions, there is no definitive ground truth that can conclusively determine a justice’ intentions. In the dialogue, justices and advocates merely discuss the specifics of a case, and the hidden intention behind a justice’s question remains concealed, making it impossible to observe directly. Nevertheless, both the justices and the advocates participating in oral arguments are highly specialized experts in the legal domain. The justices are members of the Supreme Court, and the advocates are either representatives of the Solicitor General’s office or senior attorneys. Advocates engage in thorough

preparation before oral arguments and bring substantial legal expertise to the proceedings. Based on this, our approach is as follows: we assume that the advocate can infer the justice’s intention.

For the second task, **Predicting Actions**, the question is formulated as: *"Based on the previous court exchanges, what question will the justice most likely ask next?"* The Justice’s next question can serve as the ground truth, so the correct answer is the next question extracted directly from the transcript. LLMs must integrate prior judicial questions, advocates’ responses, and the broader legal context. This task goes beyond surface-level text generation, aiming instead to emulate human-like legal reasoning. The model is required to anticipate the legal or factual issues that a justice is most likely to focus on next.

For the final task, **Generating Adaptive Responses**, we assume that the advocate and justice are sufficiently expert to provide professional and reasonable answers based on human knowledge. We then convert the advocate’s and justice’s responses into multiple-choice questions by selecting key phrases or legal terms as the correct answers and creating plausible distractors. This setup allows us to evaluate whether models can identify and select critical information, articulate core arguments, and use appropriate legal terminology in context, rather than assessing the overall coherence of long form responses.

### A.3 Prompts for Generating Questions

To generate questions for assessing understanding of intentions, we employ the following prompt:

```
You are an expert in legal dialogue analysis.
Your task is to evaluate
whether a dialogue segment contains substantive
judicial interactions worth analyzing.
```

```
BE VERY STRICT: Only approve dialogues that
  reveal significant judicial reasoning.
```

```
IMMEDIATELY REJECT if the dialogue contains:
- Simple "Yes/No" answers or basic confirmations
- Pure procedural exchanges ("Thank you, Your Honor")
- Brief interruptions or clarifications
- Administrative housekeeping matters
- Single question without substantial response
- Repetitive or redundant content
```

```
ONLY APPROVE if the dialogue demonstrates:
1. Complex judicial questioning with strategic purpose
2. Substantial advocate responses revealing judicial intent
3. Clear evidence of deep legal reasoning or analysis
```

916	4. Multiple meaningful exchanges showing	"- Simple \"Yes/No\" answers or basic	984
917	judicial thought process	confirmations\n"	985
918	5. Significant substantive legal content worth	"- Purely procedural exchanges (scheduling,	986
919	studying	format, etc.)\n"	987
920		"- Administrative comments without legal	988
921	Return a JSON object with:	substance\n"	989
922	{	"- Repetitive or redundant exchanges\n"	990
923	"generate_questions": true/false,	"- Greetings, thanks, or basic courtesies\n\n"	991
924	"complexity_score": 1-10,		992
925	"reason": "explanation for decision",	"APPROVE ONLY if the dialogue contains:\n"	993
926	"key_interactions": ["list of important	"- Complex legal reasoning or argumentation\n"	994
927	judicial questions"]	"- Judicial probing of legal concepts or	995
928	}	precedents\n"	996
929		"- Detailed factual or procedural discussions\n"	997
930	"You are a SUPREME COURT CLERK creating	"- Strategic legal positioning by advocates\n"	998
931	IMPOSSIBLY DIFFICULT questions that would	"- Constitutional or statutory interpretation\n"	999
932	challenge the finest legal minds. "	"- Multi-part judicial questioning revealing	1000
933	"Your task is to select the 3 most	judicial reasoning\n\n"	1001
934	MICROSCOPICALLY DECEPTIVE incorrect options		1002
935	that would fool Supreme Court justices	"Rate complexity on a scale of 1-10:\n"	1003
936	themselves. "	"1-3: Simple/procedural (REJECT)\n"	1004
937	"ELITE SELECTION CRITERIA - Choose options that	"4-6: Moderate complexity (CONDITIONAL)\n"	1005
938	are:\n"	"7-10: High complexity (APPROVE)\n\n"	1006
939	"1. INDISTINGUISHABLE from correct answers even		1007
940	to constitutional law professors\n"	"Return JSON format:\n"	1008
941	"2. Require GENIUS-LEVEL understanding of	"{\n"	1009
942	judicial psychology, precedent, and micro-	"  \"generate_questions\": true/false,\n"	1010
943	contextual analysis to eliminate\n"	"  \"complexity_score\": 1-10,\n"	1011
944	"3. Use IDENTICALLY sophisticated legal	"  \"reason\": \"Brief explanation\", \n"	1012
945	terminology, sentence structure, and	"  \"key_interactions\": [\"list of important	1013
946	analytical depth as correct answer\n"	exchanges or concepts\"]\n"	1014
947	"4. Represent judicial strategies so nuanced	"}"	1015
948	they differ by single words or microscopic		1016
949	timing assumptions\n"	"You are creating incorrect alternative	1017
950	"5. Test distinctions so fine that only perfect	questions for a multiple choice test about	1018
951	understanding of this specific justice's	predicting judicial questions. "	1019
952	psychology can resolve them\n"	"Your goal is to make wrong answers that are	1020
953	"6. Would be considered correct by 95% of	EASY to distinguish (target 70% correct rate	1021
954	federal judges and 90% of constitutional law	).\n\n"	1022
955	professors\n"		1023
956	"7. Differ from the correct answer by virtually	"The fake questions should be:\n"	1024
957	undetected nuances that require PhD-level	"- Most should be about obviously different	1025
958	legal theory to spot\n\n"	legal topics or concepts\n"	1026
959	"ABSOLUTE REJECTIONS - NEVER select options that	"- Many should be generic questions that could	1027
960	:\n"	be asked anytime\n"	1028
961	"- Can be eliminated by anyone short of a	"- Some should be about different stages or	1029
962	Supreme Court justice or top-tier law	timing that don't fit the current flow\n"	1030
963	professor\n"	"- Include procedural questions that don't match	1031
964	"- Use even slightly different levels of	the immediate context\n"	1032
965	sophistication, complexity, or terminology\n	"- Make the majority clearly distinguishable	1033
966	"	from the correct answer\n\n"	1034
967	"- Would be obviously wrong to federal appellate		1035
968	judges\n"	"Make these distractors EASY to distinguish -	1036
969	"- Lack the surgical precision and microscopic	they should be obviously wrong "	1037
970	accuracy of the correct answer\n\n"	"when you consider what was just discussed and	1038
971	"Return only the 3 most challenging distractors,	what would logically follow.\n\n"	1039
972	separated by commas."		1040
973		"Return JSON format:\n"	1041
974	To generate questions for Predicting Actions, we	'{\n'	1042
	employ the following prompt:	'  "possible_questions": ["question1", "	1043
975		question2", ... "question10"]\n'	1044
976	"You are an expert in legal dialogue analysis.	'}'	1045
977	Your task is to evaluate "		
978	"whether a dialogue segment contains substantive	To generate Adaptive Response Generation ques-	1046
979	judicial interactions worth analyzing.\n\n"	tions for advocate, we employ the following	1047
980	"BE VERY STRICT: Only approve dialogues that	prompt:	1048
981	reveal significant judicial reasoning.\n\n"	"You are an expert in legal dialogue analysis.	1049
982	"IMMEDIATELY REJECT if the dialogue contains:\n"	Your task is to evaluate "	1050
983		"whether a dialogue segment contains substantive	1051

1052	judicial interactions worth analyzing.\n\n"	"Rate complexity on a scale of 1-10:\n"	1120
1053		"1-3: Simple/procedural (REJECT)\n"	1121
1054	"BE VERY STRICT: Only approve dialogues that	"4-6: Moderate complexity (CONDITIONAL)\n"	1122
1055	reveal significant judicial reasoning.\n\n"	"7-10: High complexity (APPROVE)\n\n"	1123
1056			1124
1057	"IMMEDIATELY REJECT if the dialogue contains:\n"	"Return JSON format:\n"	1125
1058	"- Simple \"Yes/No\" answers or basic	"{\n"	1126
1059	confirmations\n"	" \"generate_questions\": true/false,\n"	1127
1060	"- Purely procedural exchanges (scheduling,	" \"complexity_score\": 1-10,\n"	1128
1061	format, etc.)\n"	" \"reason\": \"Brief explanation\", \n"	1129
1062	"- Administrative comments without legal	" \"key_interactions\": [\"list of important	1130
1063	substance\n"	exchanges or concepts\"]\n"	1131
1064	"- Repetitive or redundant exchanges\n"	"}"	1132
1065	"- Greetings, thanks, or basic courtesies\n\n"		
1066		We generate multiple misleading incorrect op-	1133
1067	"APPROVE ONLY if the dialogue contains:\n"	tions and then select the three most effective ones	1134
1068	"- Complex legal reasoning or argumentation\n"	to increase the difficulty of the questions.	1135
1069	"- Judicial probing of legal concepts or		
1070	precedents\n"	<b>B Further Details on Answering</b>	1136
1071	"- Detailed factual or procedural discussions\n"	<b>Questions</b>	1137
1072	"- Strategic legal positioning by advocates\n"		
1073	"- Constitutional or statutory interpretation\n"	To answer questions for understanding of inten-	1138
1074	"- Multi-part judicial questioning revealing	tions, we employ the following prompt:	1139
1075	judicial reasoning\n\n"		
1076		"You are an expert in legal language analysis	1140
1077	"Rate complexity on a scale of 1-10:\n"	and advocate statement comprehension. "	1141
1078	"1-3: Simple/procedural (REJECT)\n"	"You will be given a legal case background,	1142
1079	"4-6: Moderate complexity (CONDITIONAL)\n"	dialogue context, and a multiple-choice	1143
1080	"7-10: High complexity (APPROVE)\n\n"	question "	1144
1081		"about identifying key legal terms, phrases, or	1145
1082	"Return JSON format:\n"	concepts from advocate statements.\n\n"	1146
1083	"{\n"		1147
1084	" \"generate_questions\": true/false,\n"	"Your task is to:\n"	1148
1085	" \"complexity_score\": 1-10,\n"	"1. Carefully analyze the case background and	1149
1086	" \"reason\": \"Brief explanation\", \n"	dialogue context\n"	1150
1087	" \"key_interactions\": [\"list of important	"2. Focus on the specific advocate statement	1151
1088	exchanges or concepts\"]\n"	provided\n"	1152
1089	"}"	"3. Identify which legal concept, term, or	1153
		phrase is most relevant to the advocate's	1154
1090	To generate Adaptive Response Generation ques-	argument\n"	1155
1091	tions for justice, we employ the following prompt:	"4. Consider the legal significance and context	1156
1092		of each option\n"	1157
1093	"You are an expert in legal dialogue analysis.	"5. Select the option that best captures the key	1158
1094	Your task is to evaluate "	legal point being made\n\n"	1159
1095	"whether a dialogue segment contains substantive		1160
1096	judicial interactions worth analyzing.\n\n"	"Focus on:\n"	1161
1097		"- Core legal concepts and principles\n"	1162
1098	"BE VERY STRICT: Only approve dialogues that	"- Important statutory references and legal	1163
1099	reveal significant judicial reasoning.\n\n"	authorities\n"	1164
1100		"- Procedural requirements and legal standards\n	1165
1101	"IMMEDIATELY REJECT if the dialogue contains:\n"	"	1166
1102	"- Simple \"Yes/No\" answers or basic	"- Case law principles and precedents\n"	1167
1103	confirmations\n"	"- Constitutional provisions and interpretations	1168
1104	"- Purely procedural exchanges (scheduling,	\n"	1169
1105	format, etc.)\n"	"- Technical legal terminology essential to the	1170
1106	"- Administrative comments without legal	argument\n\n"	1171
1107	substance\n"		1172
1108	"- Repetitive or redundant exchanges\n"	"Please provide your reasoning and end with your	1173
1109	"- Greetings, thanks, or basic courtesies\n\n"	final answer in the format \boxed{{answer	1174
1110		}}."	1175
1111	"APPROVE ONLY if the dialogue contains:\n"		1176
1112	"- Complex legal reasoning or argumentation\n"	To answer questions for predicting actions, we	1177
1113	"- Judicial probing of legal concepts or	employ the following prompt:	1178
1114	precedents\n"		
1115	"- Detailed factual or procedural discussions\n"	"You are an expert in judicial questioning	1179
1116	"- Strategic legal positioning by advocates\n"	patterns and legal dialogue analysis. "	1180
1117	"- Constitutional or statutory interpretation\n"		1181
1118	"- Multi-part judicial questioning revealing		
1119	judicial reasoning\n\n"		

1182	"You will predict what question a Justice is most likely to ask next based on the context provided.\n\n"	final answer in the format <code>\boxed{{answer}}</code> ."	1250
1183			1251
1184			
1185		To answer Adaptive Response Generation questions for justice, we employ the following prompt:	1252
1186	"Your task is to:\n"		1253
1187	"1. Analyze the cumulative context and dialogue flow\n"		1254
1188	"2. Consider the Justice's established questioning patterns and priorities\n"	"You are an expert in legal language analysis and judicial questioning comprehension. "	1255
1189	"3. Evaluate the logical progression of the legal discussion\n"	"You will be given a legal case background, dialogue context, and a multiple-choice question "	1256
1190	"4. Predict what information or clarification the Justice would most likely seek next\n"	"about identifying key legal terms, phrases, or concepts from judicial questions.\n\n"	1257
1191	"5. Use your understanding of judicial behavior and legal procedure\n"		1258
1192	"6. Select the most likely next question from the given options\n\n"	"Your task is to:\n"	1259
1193		"1. Carefully analyze the case background and dialogue context\n"	1260
1194	"Focus on:\n"	"2. Focus on the specific judicial question provided\n"	1261
1195	"- Natural flow of judicial inquiry\n"	"3. Identify which legal concept, term, or phrase is most relevant to the justice's inquiry\n"	1262
1196	"- Unresolved issues that need clarification\n"	"4. Consider the legal significance and context of each option\n"	1263
1197	"- The Justice's apparent concerns and priorities\n"	"5. Select the option that best captures the key legal point being made\n\n"	1264
1198	"- Standard legal questioning protocols\n"		1265
1199	"- Follow-up questions that logically flow from the current discussion\n\n"	"Focus on:\n"	1266
1200	"Please provide your reasoning and end with your final answer in the format <code>\boxed{{answer}}</code> ."	"- Core legal concepts and principles\n"	1267
1201		"- Important statutory references and legal authorities\n"	1268
1202	To answer Adaptive Response Generation questions for advocate, we employ the following prompt:	"- Procedural requirements and legal standards\n"	1269
1203		"- Case law principles and precedents\n"	1270
1204		"- Constitutional provisions and interpretations\n"	1271
1205		"- Technical legal terminology essential to the argument\n\n"	1272
1206		"Please provide your reasoning and end with your final answer in the format <code>\boxed{{answer}}</code> ."	1273
1207			1274
1208			1275
1209			1276
1210			1277
1211			1278
1212			1279
1213			1280
1214			1281
1215			1282
1216			1283
1217			1284
1218			1285
1219			1286
1220			1287
1221			1288
1222			1289
1223			1290
1224			
1225		<b>C Analysis of Legal Experts</b>	1291
1226		The correct answer is C of the Question 3, and this question is interesting precisely because it exposes the pressure point of Supreme Court advocacy.	1292
1227		Justice Jackson isn't quibbling over wording for sport. She's performing a jurisdiction-defining move. When she asks what the question "presented in this case actually is," she's signaling something fundamental: the Court can only decide what is properly before it. Everything else-policy arguments, doctrinal cleanup, institutional anxiety-is noise unless tethered to that question.	1293
1228		1. Drawing a hard boundary around the Court's role: Supreme Court justices are acutely sensitive to overreach. By drilling into the Question Presented, Justice Jackson is asking whether the Court is being asked to decide whether a uniquely stringent standard applies at all, or merely whether this	1294
1229			1295
1230			1296
1231			1297
1232			1298
1233			1299
1234			1300
1235			1301
1236			1302
1237			1303
1238			1304
1239			1305
1240			1306
1241			1307
1242			1308
1243			
1244			
1245			
1246			
1247			
1248			
1249			

particular standard (bad faith or gross misjudgment) qualifies as such. That distinction determines the scope of the ruling—narrow case resolution versus sweeping doctrinal shift.

2. Exposing strategic ambiguity in advocacy: Lisa Blatt attempts to present two readings of the petition. Justice Jackson’s questioning highlights that advocates cannot argue narrow when their briefing argues broad. This reveals how credibility and consistency are essential currencies before the Court.

3. Framing everything that follows: Once the Court locks in what the question is, entire lines of argument either survive or collapse—remedies, intent standards, and statutory interpretation all depend on it.

Justice Jackson’s question is interesting because it is a gatekeeping move. It forces clarity, constrains overreach, and reinforces that the Supreme Court is a court of limited questions with lasting consequences.

---

**Case Name:** A.J.T. v. Osseo Area Schools, Independent School District

---

**Question:**  
question: "What was the judge’s likely intent for asking ‘presented in this case actually is.’?"  
A: "To probe the consistency of the advocate’s argument with previous case law.",  
B: "To explore the broader implications of the legal standards being debated.",  
C: "To clarify the specific legal question that the court needs to address in this case.",  
D: "To challenge the advocate’s interpretation of the legal standards involved."

**Answer:** "C"

---

Table 3: Example Analyzed by Legal Experts

**D Representative Dialogue Cases**

We provide a representative oral argument case, as shown in Table 6. The table includes a portion of the preceding context from the argument and the corresponding question generated by our justice agent. This illustrates how our method captures the legal dialogue and generates contextually relevant questions.

**E Details of Difficulty Levels**

---

**Case Name:** A.J.T. v. Osseo Area Schools, Independent School District

---

**Model Name:** o3

---

**Predicted Answer:** A

---

**Explanation:** Justice Sotomayor has just pointed out that intent “is not even an issue here” because the family already obtained IDEA relief, and counsel responded, “They want more.” A natural follow-up for Justice Sotomayor is to refocus the argument on what the Court actually granted certiorari to decide. She is likely to press counsel to confirm that, notwithstanding whatever additional relief the plaintiffs might desire, the only legal question properly before the Court is whether the Eighth Circuit applied the correct intent (versus heightened-intent) standard under Section 504/ADA. Option A mirrors that clarification: it brackets counsel’s “they want more” remark and zeroes in on whether the dispute before the Court is simply the level of intent required. The other options either broaden the discussion too abstractly (B, D) or divert it to remedies under IDEA (C), which is not what she has just highlighted.

---

Table 4: OpenAI o3 Model Output Example

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

---

**Case Name:** A.J.T. v. Osseo Area Schools, Independent School District

---

**Model Name:** Qwen2.5-72B

---

**Predicted Answer:** B

---

**Explanation:** Based on the previous court exchanges and the specific context provided, Justice Sotomayor is most likely to ask:**\*\*Option B: Is there a distinction between seeking damages and obtaining injunctive relief under the IDEA?\*** Reasoning: 1. **\*\*Natural Flow of Judicial Inquiry\*\***: Justice Sotomayor has been focusing on the definitions and implications of "bad faith" and "intent to discriminate." She is concerned about the broader application of these concepts across different statutes. 2. **\*\*Unresolved Issues\*\***: There is a need to clarify how the standards for intent apply differently to various forms of relief, such as damages versus injunctive relief. 3. **\*\*Apparent Concerns and Priorities\*\***: Justice Sotomayor has shown a concern for the procedural and definitional clarity in discrimination cases, particularly in the context of the IDEA. 4. **\*\*Standard Legal Questioning Protocols\*\***: It is common for justices to explore the distinctions between different types of legal remedies to ensure that the court's decision is well-grounded and consistent with existing legal principles. Thus, Option B aligns best with the natural progression of the legal discussion and Justice Sotomayor's apparent concerns.

---

Table 5: Qwen2.5-72B Model Output Example

---

**Court Dialogue**

...

**Justice Sotomayor:** But they didn't pass the policy with antireligion animus.

**Lisa S. Blatt:** If you – let me just give you another example.

**Justice Sotomayor:** You're asking – when you're using the words "bad faith," you're talking about animus.

**Lisa S. Blatt:** No, I'm talking about — and you can — you're in charge, so you can say: Intent to discriminate is the standard. We're not going to use bad faith. We don't like that word...

**Justice Sotomayor:** Counsel, it would have been nice to have known that we were biting off that big a chunk.

**Lisa S. Blatt:** I agree. But in terms of what we had to do when you granted cert was look at the text, and then the blue brief said that there is no intent required.

**Justice Sotomayor:** By the way, intent's not even an issue here because there wasn't an injunction being — or the lack of an injunction — challenged here. They got the injunction under the IDEA, didn't they?

**Lisa S. Blatt:** They want more.

---

**Question for Justice Agent:**

Question: Based on the preceding exchange, what question will Justice Sotomayor most likely ask next?

- A. Well, we can put aside whether they want more. But the only thing before us on the decision below is whether it's an intent standard or a heightened standard, correct?
- B. In determining the appropriate standard for this case, should we consider whether courts must require proof of intentional discrimination instead of negligence or indifference?
- C. Should we focus on the adequacy of the remedies available under the IDEA, rather than changes to the intent standard?
- D. How does the interpretation of intent in discrimination cases differ between educational and employment contexts under Title VII?

---

**Correct Answer: A**

---

Table 6: Court Dialogue and Predicting Actions Example

## Court Dialogue

...

**Justice Sotomayor:** I’m not sure where any of these tests come from, because mens rea is generally willfulness, which requires knowing what the law is, but the statute doesn’t talk about willfulness. Motive, in – intent – we don’t care about motive. We’ve said that repeatedly in a bunch of different contexts. It’s do you know you’re doing the act and are you intending to do the act. If it violates the law, you’re guilty. Pardon the pun. This is a tort, but you’re responsible. Or you do it knowingly, knowing that you’re doing the act. So I don’t know where the bad faith comes from. I don’t know where the gross indifference comes from. I don’t know where the deliberate indifference comes from. Have you figured that out?

### Question for Advocate Agent:

Question: What was the judge’s likely intention for asking ‘I’m not sure where any of these tests come from, because mens rea is generally willfulness, which requires knowing what the law is, but the statute doesn’t talk about willfulness. Motive, in – intent – we don’t care about motive. We’ve said that repeatedly in a bunch of different contexts. It’s do you know you’re doing the act and are you intending to do the act. If it violates the law, you’re guilty. Pardon the pun. This is a tort, but you’re responsible. Or you do it knowingly, knowing that you’re doing the act. So I don’t know where the bad faith comes from. I don’t know where the gross indifference comes from. I don’t know where the deliberate indifference comes from. Have you figured that out?’?

- Clarify the advocate’s stance on the interplay between statutory causation requirements and mens rea.
- Assess the advocate’s understanding and articulation of the statutory basis for different mens rea standards within disability discrimination law.
- Probe the advocate’s ability to distinguish between statutory interpretation and judicially created tests.
- Evaluate the advocate’s interpretation of the statutory language regarding intent requirements.

**Correct Answer: B**

Table 7: Court Dialogue and Understanding Intention Example

Model Names	Understanding Intention			Predicting Actions			Generating Responses for Advocate			Generating Responses for Justice		
	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC
<b>Closed-source Models</b>												
o3 (Reasoning)	<b>94.82</b>	93.53	92.35	<b>77.09</b>	89.26	<b>91.23</b>	<b>88.43</b>	<b>89.31</b>	<b>87.33</b>	<b>83.47</b>	<b>82.61</b>	<b>84.27</b>
GPT-4o	92.06	89.71	90.70	62.72	72.67	88.37	76.66	76.53	77.85	73.04	71.06	75.53
<b>Open-source Models</b>												
Qwen2.5_1.5b_instruct	82.35	83.53	88.53	60.66	66.03	62.72	68.42	70.97	71.54	64.54	66.17	64.89
Qwen2.5_7b_instruct	89.12	89.12	91.76	51.82	78.83	77.41	75.38	75.33	74.38	71.97	72.55	72.77
Qwen2.5_72b_instruct	90.29	92.35	<b>96.18</b>	48.50	78.20	80.57	<b>83.46</b>	<b>85.77</b>	<b>84.82</b>	<b>78.34</b>	<b>78.30</b>	<b>79.15</b>
Qwen3_4b_instruct	86.18	86.76	90.29	45.02	67.14	66.19	72.57	71.24	72.79	67.09	68.30	67.45
Qwen3_30b_instruct	94.41	<b>94.71</b>	95.00	63.35	84.20	85.47	78.32	78.54	78.32	75.80	78.09	76.81
Gemma3_4b_instruct	76.76	76.76	89.12	32.07	52.13	54.82	62.59	61.86	63.19	57.11	58.51	59.36
Phi-4	91.76	91.76	95.88	42.65	74.25	71.41	78.20	77.99	77.61	76.43	77.66	77.66
Llama3.3_70b_instruct	94.71	<b>94.71</b>	<b>96.18</b>	<b>67.14</b>	<b>89.42</b>	<b>89.10</b>	79.32	79.70	80.65	77.49	77.02	78.30
Saul_7b_instruct	79.41	81.18	72.35	35.07	52.13	49.45	53.57	48.77	40.99	43.52	42.55	36.60

Table 8: Performance on **EASY** difficulty questions. **BG** provides the case background, **BG+CC** further includes the current conversation, and **BG+CC+PC** incorporates the preceding conversations. We evaluate Understanding Intention, Action Prediction, and Response Generation using accuracy (%). **Bold-underlined** indicates the strongest among all models; **bold** indicates the strongest among open-source models.

Model Names	Understanding Intention			Predicting Actions			Generating Responses for Advocate			Generating Responses for Justice		
	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC
<b>Closed-source Models</b>												
o3 (Reasoning)	86.06	86.49	84.31	<b>75.88</b>	89.14	<b>91.44</b>	<b>86.09</b>	<b>87.34</b>	<b>86.63</b>	<b>75.48</b>	<b>76.74</b>	<b>78.29</b>
GPT-4o	85.40	83.22	77.94	60.23	70.38	89.22	72.96	75.05	75.50	60.36	58.99	65.25
<b>Open-source Models</b>												
Qwen2.5_1.5b_instruct	67.97	69.93	61.87	58.67	64.46	60.51	52.69	56.13	58.55	49.70	54.75	55.56
Qwen2.5_7b_instruct	81.26	80.83	84.97	49.79	77.86	75.74	63.45	67.29	67.29	54.33	59.60	61.41
Qwen2.5_72b_instruct	86.49	<b>88.45</b>	<b>93.03</b>	44.85	77.86	80.68	<b>76.81</b>	<b>77.88</b>	<b>80.11</b>	<b>70.02</b>	<b>70.51</b>	<b>71.92</b>
Qwen3_4b_instruct	76.91	77.34	79.96	41.47	65.59	63.61	58.97	61.54	65.17	46.88	52.12	52.93
Qwen3_30b_instruct	<b>88.45</b>	<b>88.45</b>	89.76	60.51	83.64	84.91	72.44	72.01	71.79	62.58	64.24	65.45
Gemma3_4b_instruct	61.22	61.22	71.46	28.63	49.22	52.19	45.64	48.14	54.09	37.63	42.63	45.86
Phi-4	84.97	86.27	89.11	37.80	73.77	70.38	72.73	72.86	71.19	65.79	68.08	68.48
Llama3.3_70b_instruct	86.71	86.93	88.02	<b>64.60</b>	<b>89.28</b>	<b>89.28</b>	73.47	74.91	77.88	67.81	68.28	69.09
Saul_7b_instruct	67.32	67.54	59.04	32.86	51.48	47.25	36.92	36.25	36.06	33.60	35.96	34.34

Table 9: Performance on **MEDIUM** difficulty questions. **BG** provides the case background, **BG+CC** further includes the current conversation, and **BG+CC+PC** incorporates the preceding conversations. We evaluate Understanding Intention, Action Prediction, and Response Generation using accuracy (%). **Bold-underlined** indicates the strongest among all models; **bold** indicates the strongest among open-source models.

Model Names	Understanding Intention			Predicting Actions			Generating Responses for Advocate			Generating Responses for Justice		
	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC	BG	BG+CC	BG+CC+PC
<b>Closed-source Models</b>												
o3 (Reasoning)	73.14	71.90	68.18	<b>70.22</b>	84.93	<b>88.89</b>	<b>73.53</b>	<b>74.58</b>	<b>75.00</b>	<b>62.67</b>	<b>62.86</b>	<b>65.18</b>
GPT-4o	61.16	59.09	<b>77.00</b>	55.51	63.60	85.58	59.34	60.67	61.52	43.85	42.66	50.00
<b>Open-source Models</b>												
Qwen2.5_1.5b_instruct	52.89	54.55	43.80	55.70	60.48	56.62	40.91	45.90	49.36	37.10	40.28	41.47
Qwen2.5_7b_instruct	52.48	51.65	50.41	43.20	72.61	71.14	43.64	49.00	51.37	40.67	44.44	46.03
Qwen2.5_72b_instruct	64.46	65.70	64.88	38.79	71.32	74.63	<b>57.45</b>	<b>59.74</b>	<b>62.30</b>	<b>45.44</b>	<b>47.22</b>	<b>50.60</b>
Qwen3_4b_instruct	46.28	46.69	46.28	36.76	60.29	58.64	42.38	44.55	50.10	34.33	35.71	38.49
Qwen3_30b_instruct	60.33	59.92	61.16	53.31	79.04	80.51	50.10	54.26	56.63	40.08	43.25	44.64
Gemma3_4b_instruct	51.24	51.24	49.59	24.26	42.46	44.85	36.91	39.89	44.26	31.75	35.12	34.92
Phi-4	66.12	66.12	66.53	32.35	67.28	64.15	49.82	51.73	52.28	42.26	46.63	50.00
Llama3.3_70b_instruct	<b>71.90</b>	<b>73.55</b>	<b>68.60</b>	<b>59.01</b>	<b>85.29</b>	<b>84.93</b>	50.55	55.01	58.83	44.05	46.23	50.00
Saul_7b_instruct	47.52	47.93	43.80	30.51	47.43	43.93	32.00	31.33	31.15	25.00	26.59	24.60

Table 10: Performance on **HARD** difficulty questions. **BG** provides the case background, **BG+CC** further includes the current conversation, and **BG+CC+PC** incorporates the preceding conversations. We evaluate Understanding Intention, Action Prediction, and Response Generation using accuracy (%). **Bold-underlined** indicates the strongest among all models; **bold** indicates the strongest among open-source models.