

Perceptions to Beliefs: Exploring Precursory Inferences for Theory of Mind in Large Language Models

Chani Jung[♡] Dongkwan Kim[♡] Jiho Jin[♡] Jiseon Kim[♡]
Yeon Seonwoo[♣] Yejin Choi[◇] Alice Oh[♡] Hyunwoo Kim[♠]

[♡]KAIST [♣]Amazon [◇]University of Washington [♠]Allen Institute for AI

{1016chani, dongkwan.kim, jinjh0123, jiseon_kim}@kaist.ac.kr

yseonwoo@amazon.com, yejin@cs.washington.edu, alice.oh@kaist.edu, hyunwook@allenai.org

Abstract

While humans naturally develop theory of mind (ToM), the capability to understand other people’s mental states and beliefs, state-of-the-art large language models (LLMs) underperform on simple ToM benchmarks. We posit that we can extend our understanding of LLMs’ ToM abilities by evaluating key human ToM precursors—*perception inference* and *perception-to-belief inference*—in LLMs. We introduce two datasets, Percept-ToMi and Percept-FANToM, to evaluate these precursory inferences for ToM in LLMs by annotating characters’ perceptions on ToMi and FANToM, respectively. Our evaluation of eight state-of-the-art LLMs reveals that the models generally perform well in perception inference while exhibiting limited capability in perception-to-belief inference (e.g., lack of inhibitory control). Based on these results, we present PercepToM, a novel ToM method leveraging LLMs’ strong perception inference capability while supplementing their limited perception-to-belief inference. Experimental results demonstrate that PercepToM significantly enhances LLM’s performance, especially in false belief scenarios.

1 Introduction

Humans interact with others in various social situations using *theory of mind* (ToM), the cognitive capability to understand other’s mental states (e.g., beliefs, desires, and thoughts; Premack and Woodruff, 1978). While ToM is naturally developed for humans in childhood, large language models (LLMs) are known to exhibit inconsistency in ToM tasks (van Duijn et al., 2023; Trott et al., 2023). Despite some early reports of successful cases (Whang, 2023; Street et al., 2024), studies have shown that even state-of-the-art LLMs significantly lag behind human performance in ToM tasks, particularly in false belief tests (Le et al., 2019; Kim et al., 2023; Gandhi et al., 2023; Wu et al., 2023; Shapira et al., 2024).

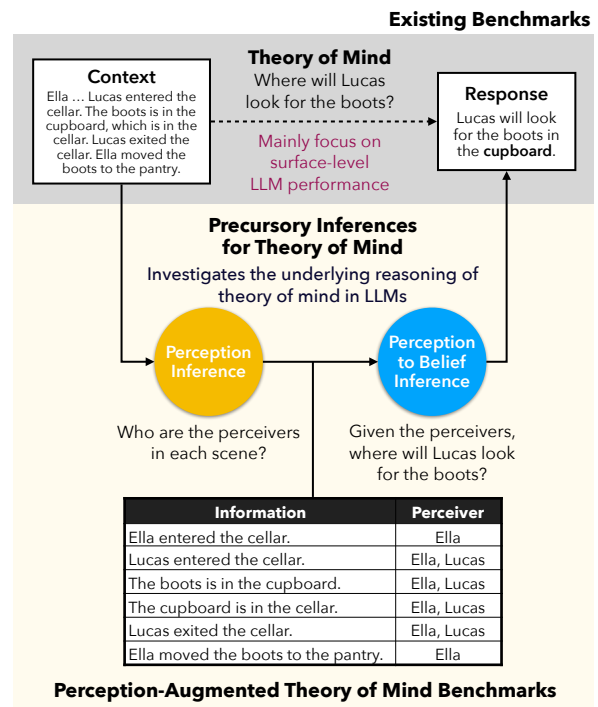


Figure 1: Inspired by children’s developmental trajectory for theory of mind (ToM), our perception-augmented ToM benchmarks test the two precursory inferences of ToM in LLMs in order to examine their underlying social reasoning capabilities: (1) *perception inference* and (2) *perception-to-belief inference* (§2).

However, there are clear limitations to understanding the gaps in LLMs’ underlying ToM abilities based on the evaluation results of existing ToM benchmarks, which only focus on the accuracy of the models’ responses to ToM questions (Ma et al., 2023b). Although some studies have conducted error analysis based on model responses (Ma et al., 2023a; Wu et al., 2023), they rely on qualitative analysis via human inspection.

Psychology literature describes precursory steps to ToM development: *perception inference* (Rakoczy, 2022) and *perception-to-belief inference*—understanding that ‘seeing leads to knowing’ (Pratt and Bryant, 1990; Baron-Cohen and

Goodhart, 1994). These capabilities can be defined in the scenario shown in Figure 1. We refer to the ability to infer others’ perceptions (e.g., *Who are the perceivers of each scene?*) as *perception inference* and the process of deducing others’ beliefs from their perceptions (e.g., *Given the perceivers of each scene, where will Lucas look for the boots?*) as *perception-to-belief inference*.

Inspired by the human developmental stages for ToM, we evaluate the key precursory inference steps of ToM in LLMs. First, we extend the two representative ToM benchmarks, ToMi (Le et al., 2019) and FANToM (Kim et al., 2023), by annotating characters’ perceptions about each piece of information from the input context. Figure 1 illustrates an example of our annotations and tasks on ToMi. Second, using our new benchmarks, we evaluate eight state-of-the-art LLMs and find that models perform generally well in *perception inference* but perform poorly in the *perception-to-belief inference* task (§3.2 and §3.3). We also find that LLMs have weak *inhibitory control* when inferring beliefs – i.e., the capability of suppressing irrelevant information (§4.4).

Based on these findings, we propose PercepToM, a novel framework to enhance the ToM in LLMs by leveraging their perception inference capability. PercepToM first guides LLMs to infer the characters’ perceptions from an input context. Then, it aids LLMs in perception-to-belief inference through the *perspective context extraction* step, which isolates the context perceived by the target character with a simple string-matching algorithm. Finally, LLMs answer to the ToM questions given the isolated context. This approach leads to improved performance on both ToMi and FANToM, particularly on the false belief scenarios (§4.3).

Our contributions are as follows. First, we construct perception-augmented ToM benchmarks which enable the evaluation of the two precursory inferences for ToM in LLMs (§2): *perception inference* and *perception-to-belief inference*. Second, using these benchmarks, we show that current LLMs are good at inferring the perceptions of others but struggle to infer beliefs from the perceptual information (§3.2, §3.3, and §4.4). Lastly, we introduce the PercepToM framework to improve LLMs’ ToM reasoning by leveraging their strong *perception inference* while supplementing their *perception-to-belief inference* (§4). Our method enhances LLMs’ performance on ToMi and FANToM, especially on the false belief scenarios (§4.3).

2 Augmenting Perceptions on Theory of Mind Benchmarks

We construct perception-augmented theory of mind (ToM) benchmarks to evaluate two essential cornerstones for ToM in LLMs: (1) *perception inference* and (2) *perception-to-belief inference* capabilities.

2.1 Perception Inference and Perception-to-Belief Inference

The precursory inferences for ToM (Rakoczy, 2022) can be understood through the Sally-Anne test, a widely recognized psychological assessment for evaluating ToM (Baron-Cohen et al., 1985). In this test, Sally initially observes a marble in a box but does not witness Anne moving the marble to a basket after she leaves the room.

Current ToM benchmarks predominantly assess LLMs based on their surface-level performance on ToM questions (e.g., *Where will Sally look for the marble when she returns?*), leaving their underlying inference capabilities underexplored. To address this gap, drawing from psychology literature, we refer to the ability to infer others’ perceptions (e.g., *Did Sally see Anne moving marble to the basket?*) as *perception inference*. Additionally, we define the process of deducing beliefs from perceptual information (e.g., *Sally did not see the marble being moved. Where will she look for the marble when she returns?*) as *perception-to-belief inference*.

To further investigate these inferences, we construct PercepToMi and PercepToFANToM by annotating each character’s perception of information within the context of the two benchmark datasets ToMi (Le et al., 2019) and FANToM (Kim et al., 2023). Annotation examples are presented in Figure 2.

2.2 The Source Theory of Mind Benchmarks

ToMi (Le et al., 2019) We include ToMi, one of the most widely used ToM benchmarks for reading comprehension tasks. The contexts in ToMi feature narrative scene descriptions, assuming characters acquire information by visual perception. In each story, several characters are present in a room along with an object. The story implicitly presumes that the characters can observe all objects and events taking place within the room. There are four ToM question types in ToMi for a given story: first-order true/false beliefs, and second-order true/false beliefs. In the true belief scenario, all characters ob-

Story in Percept-ToMi 		Conversation in Percept-FANToM 	
Information	Perceivers	Information	Perceivers
Ella entered the cellar.	Ella	Gianna: Guys, I need to change clothes for a meeting later. Talk to you later!	Gianna, Sara, Javier
Lucas entered the cellar.	Ella, Lucas	Sara: Sure thing, Gianna. Take care!	Gianna, Sara, Javier
Benjamin entered the porch.	Benjamin	Javier: Catch you later, Gianna.	Gianna, Sara, Javier
The boots is in the cupboard.	Ella, Lucas	Sara: So Javier, have you ever tried training Bruno?	Sara, Javier
The cupboard is in the cellar.	Ella, Lucas	Javier: Yes, it was a challenge at times, but rewarding nevertheless. How about you?	Sara, Javier
Lucas exited the cellar.	Ella, Lucas	...	
Benjamin exited the porch.	Benjamin	Gianna: Hey guys, I'm back, ... It's amazing how pets further strengthens the bond	Gianna, Sara, Javier
Ella moved the boots to the pantry.	Ella	Sara: Absolutely! The fact that they trust us enough to learn from us is really special.	Gianna, Sara, Javier
The pantry is in the cellar.	Ella	Javier: I can't agree more.	Gianna, Sara, Javier

Figure 2: Example data in Percept-ToMi and Percept-FANToM. For each context, the perceivers of every scene description or utterance are annotated automatically (Percept-ToMi) and manually (Percept-FANToM).

serve everything happening in the room, ensuring that they share identical access to the information. However, in the false belief scenario, a character leaves the room, and then another character moves the object from one container to another, resulting in information asymmetry about the same object.

FANToM (Kim et al., 2023) This recent benchmark reveals a significant performance gap between humans and state-of-the-art LLMs. It consists of multi-party conversations, assuming information transfer through both visual and auditory perceptions. The information asymmetry occurs as some of the characters leave or join the conversation. When a character is absent, the remaining participants share information exclusively among themselves. FANToM also includes true belief scenarios where the absent character gets informed about the conversation upon rejoining the group.

2.3 Perception-Augmented ToM Benchmarks

Percept-ToMi To construct Percept-ToMi, we sample 150 story-question pairs for each of the four ToM question types in ToMi¹: first-order true/false beliefs, and second-order true/false beliefs. We automatically annotate perceivers of the scenes in ToMi using SymbolicToM (Sclar et al., 2023) and manually verify the samples. SymbolicToM tracks the witnesses of each scene by maintaining a graphical representation of the true world state, allowing us to obtain the list of perceivers for each scene from its output. However, upon verifying 50 samples of the SymbolicToM output, we identify two types of errors in the perceiver annotations and correct them across our entire dataset. The details of

¹We use the *Fixed and Disambiguated ToMi* constructed by Sclar et al. (2023), where sentences are inserted to disambiguate the location of containers in the story, and some mislabeled questions are corrected.

this verification and correction of perceiver annotations are explained in Appendix A.1.

Percept-FANToM To build Percept-FANToM, we use all of the short conversations in FANToM, but exclude those that cause errors in our perception annotation format. We assume that a character is the perceiver of all utterances that occur from the time they join the conversation until the time they leave. After two of the authors confirmed the criteria for determining the joining and leaving times (Appendix A.2), each data point was manually annotated by one of the authors. They followed the annotation criteria mechanistically, ensuring no subjectivity was involved. Additionally, the authors randomly selected 20 samples to check for any discrepancies in their annotations. The results confirmed that all annotations were consistent between both authors. Based on the annotations of characters' joining and leaving times, perceivers of each utterance are automatically mapped. The statistics of our perception-augmented ToM benchmarks and the source benchmarks are shown in Appendix B.

2.4 Task and Evaluation

We measure the performance of (1) *perception inference* and (2) *perception-to-belief inference* in both false belief and true belief scenarios.

(1) Perception Inference In order to evaluate the perception inference capability of LLMs, we prompt the models to track characters' perception of each unit of information in the input context. Specifically, we require the models to respond in the format of a JSON array, which consists of JSON objects containing a unit of information from the context as a key and the perceivers of the informa-

tion as a value.² We use individual sentences and utterances as the units of information for ToMi and FANToM, respectively. To ensure the generated answers are in the correct format, we provide an example format of the JSON array using a dummy sentence that does not appear in the datasets. The example input prompt is in Appendix C.1.

(2) Perception-to-Belief Inference To evaluate the perception-to-belief inference capability of the models, we provide them with a ground truth perception inference result and then query ToM questions from the original benchmarks. The ground truth perception inference result is provided in the same JSON array format we use to evaluate the perception inference capability of LLMs. The example and detailed explanation of the input prompt can be found in Appendix C.2.

3 Precursory Inferences of ToM in LLMs

3.1 Experimental Setup

We analyze the *perception inference* and *perception-to-belief inference* (§2.1) performances of LLMs on Percept-ToMi and Percept-FANToM (§2.3) with the following metrics and models.

Perception Inference To evaluate the model-generated perception inference results, we calculate accuracy for a given input context based on the ratio of information units where the model accurately identifies the perceivers. The final *perception inference accuracy* for a dataset is obtained by averaging the accuracies across all contexts in the dataset. In Percept-ToMi, accuracy is calculated across the stories, with each story paired with a single ToM question. In Percept-FANToM, since multiple questions share a single context, we calculate accuracy across these contexts.

Perception-to-Belief Inference and ToM We evaluate the perception-to-belief inference and ToM performance of LLMs using the original questions and answers from ToMi (Le et al., 2019) and FANToM (Kim et al., 2023). For ToMi, we measure accuracy by the ratio of correctly answered questions among all story-question pairs. Note that we do not use the *joint accuracy* proposed in the original ToMi, where a story is counted as correctly answered only if all questions about the story are

²We structure the perception inference results in JSON to leverage its parsability and interpretability. Also, recent works use JSON format to improve language model generation quality (Zhou et al., 2023; OpenAI, 2023).

answered correctly. This is because many of the stories in the Fixed and Disambiguated ToMi (Sclar et al., 2023) do not include all six question types of ToMi. For Percept-FANToM, we report the *set:ALL* score, which requires the model to correctly answer five types of ToM questions³ for the same piece of information within a conversation.

Correlation between LLM’s ToM Performance and Precursory Inference Performance To analyze the relationship between LLMs’ ToM capability and their performance on perception-related ToM precursor tasks (i.e., perception inference and perception-to-belief inference), we measure the Pearson correlation coefficient between models’ performances on ToM and each of these two tasks.

Target Models We examine eight state-of-the-art LLMs: GPT-3.5 Turbo (gpt-3.5-turbo-1106), GPT-4 Turbo (gpt-4-turbo-1106-preview), GPT-4o (gpt-4o-2024-05-13)⁴, Claude 3 (Haiku and Sonnet)⁵, Gemini 1.0 Pro (Gemini-Team, 2024), Llama-3 70B Instruct (AI@Meta, 2024), and Mixtral 8x22B Instruct (Jiang et al., 2024) on Percept-ToMi and Percept-FANToM (§2.3).

3.2 Perception Inference

LLMs generally perform well on perception inference across datasets and scenarios. As shown in Figure 3, most LLMs exhibit high accuracy on perception inference in both Percept-ToMi and Percept-FANToM. The models’ average perception inference accuracy is 0.781 on Percept-ToMi and 0.926 on Percept-FANToM. Also, they exhibit negligible differences in the accuracy between the true belief and false belief scenarios. In ToMi, all models except for GPT-3.5 Turbo and Gemini 1.0 Pro exhibit a gap of less than 0.1 accuracy between the two scenarios. In FANToM, the accuracy gaps between the two scenarios in all models are no greater than 0.014. This result contrasts with the models’ large performance gap in the two scenarios on ToM questions, suggesting that their limited ToM performance in false belief scenarios is not due to the lack of perception inference capability. Detailed results are in Appendix E.

The perception inference and ToM performance do not show a strong correlation. Especially

³BELIEFQ_[CHOICE], ANSWERABILITY Q_[LIST], INFOACCESS Q_[LIST], ANSWERABILITY Q_[Y/N], INFOACCESS Q_[Y/N]

⁴<https://platform.openai.com/docs/models/overview>

⁵<https://www.anthropic.com/product>

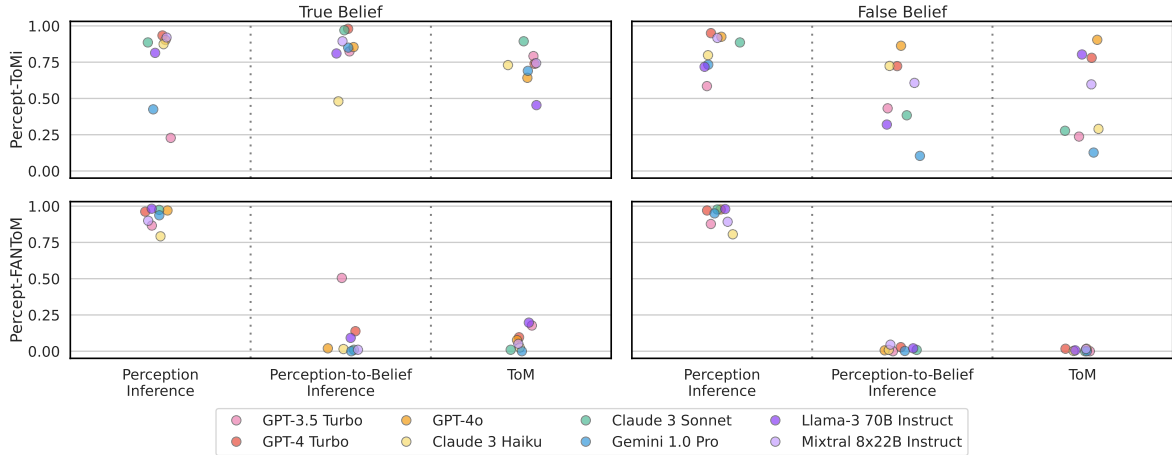


Figure 3: Perception inference, perception-to-belief inference, and ToM performances of LLMs in true and false belief scenarios of Percept-ToMi and Percept-FANToM. Although the models exhibit similar accuracy in perception inference across both scenarios, their performance in perception-to-belief inference and ToM scenarios varies significantly.

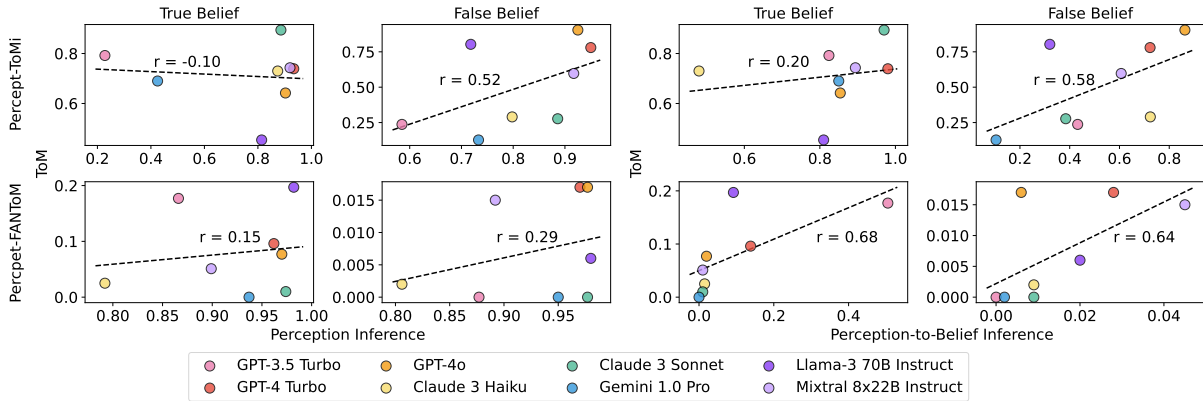


Figure 4: Pearson correlation of LLMs' ToM performance with perception inference (left) and perception-to-belief inference (right) performances. ToM performance shows a positive correlation with perception-to-belief inference performance but exhibits a weak or no correlation with perception inference performance.

in the ToMi's true belief scenarios, the two performances exhibit a near-zero correlation (Figure 4). Although moderate correlations appear in other scenarios, the correlation coefficients are not statistically significant. These results imply that LLMs' perception inference capability is not directly linked to their ToM performance. This contrasts with human adults, where ToM is strictly dependent on perception inference.

3.3 Perception-to-Belief Inference

LLMs struggle with perception-to-belief inference. Surprisingly, although the ground truth perception information for all characters is provided in this task, models still underperform in false belief scenarios compared to true belief scenarios (see Figure 3). This trend is consistent with their ToM

performance. Moreover, their performances on the perception-to-belief inference task are mostly similar to their ToM performances in all scenarios except for the ToMi true belief scenario. The fact that the LLMs hardly benefit from the additional character perception information, which should serve as significant hints for solving ToM questions, suggests that they have limited capability to infer beliefs from perceptions. The exact performances of models are in Appendix E.

The perception-to-belief inference and ToM performance exhibit a positive correlation. This is consistent across all datasets and scenarios (Figure 4). Notably, in FANToM, models exhibit a high correlation between the two performances ($r > 0.6$). This correlation likely arises because the two tasks use the same questions. However,

since LLMs are showing similar performances in both tasks, we can see that they are not fully leveraging the ground truth perception information in the perception-to-belief inference task.

4 PercepToM: Grounding ToM Reasoning on Perception

4.1 Framework

According to our experiment results, LLMs perform adequately well in both true and false belief scenarios on perception inference, while they underperform in perception-to-belief inference (§3). Based on these findings, we propose PercepToM, a framework for improving LLM’s ToM reasoning by grounding it in perception information. PercepToM leverages LLM’s strong perception inference capabilities while enhancing its perception-to-belief inference with a simple string-matching rule. PercepToM consists of the following steps as illustrated in Figure 5:

1. **Perception Inference:** The LLM infers which characters perceived each unit of information in the context (e.g., a scene description or an utterance).
2. **Perspective Context Extraction:** Based on the perception inference result from the LLM, PercepToM extracts the *perspective context* — i.e., the subset of the input context identified by the LLM as perceived by the target character. This process is conducted by a simple string-matching procedure.
3. **Response Generation:** Given the perspective context of the target character, the LLM answers the ToM question.

If the model correctly performs perception inference, the perspective context will only include what the target character has perceived from the original context – that is, what they believe to be true, based on the principle of rational belief (Baker et al., 2011). When given this isolated context along with the ToM question, the scenario becomes a simple true belief scenario, wherein the LLM has access to the same information as the target character (i.e., information symmetry).

SymbolicToM (Sclar et al., 2023) also helps LLM’s ToM reasoning by providing only the context included in the target character’s belief state graph to the model. However, constructing the belief graph in SymbolicToM requires manually

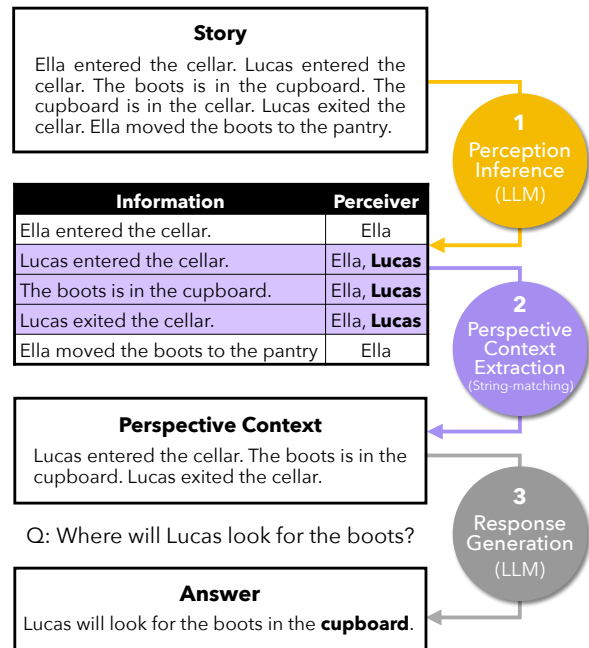


Figure 5: An overview of our PercepToM framework, which enhances LLMs’ ToM reasoning by (1) instructing LLMs to infer the perceivers of each information in the context; (2) aiding their perception-to-belief inference through the *perspective context extraction* step, which isolates the context perceived by the target character; and (3) allowing LLMs to generate responses to ToM questions based on this perspective context.

crafted algorithms tailored to different types of input. In contrast, PercepToM avoids this requirement by leveraging LLM’s perception inference capabilities, which can handle more diverse and complicated contexts, thereby achieving significantly improved generalizability. The example input and output of each step of the algorithm are provided in Appendix D.

4.2 Experimental Setup

Datasets and Metrics We evaluate PercepToM and other baseline models on Fixed and Disambiguated ToMi (Le et al., 2019; Sclar et al., 2023) and FANToM (Kim et al., 2023). As evaluation metrics, we use the ratio of correctly answered story-question pairs for ToMi and the *set:ALL* score for FANToM, as employed in the ToM performance evaluation in the previous section (§3.1).

Baselines We compare Vanilla, Chain-of-Thought (CoT; Wei et al., 2022), and System 2 Attention (S2A; Weston and Sukhbaatar, 2023) with PercepToM. Vanilla involves LLM directly answering questions based on the given context, while CoT adds the prompt “Let’s think step by

step.” to help the model answer ToM questions. S2A improves the reasoning of LLMs by prompting them to extract only the relevant part of the input context before yielding a final response. By using S2A as a baseline, we compare the effectiveness of the perspective context of PercepToM with the relevant context extracted by LLMs using S2A. We also compare the performance of PercepToM to that of SymbolicToM (Sclar et al., 2023) on ToMi. However, we do not extend this comparison to FANToM, as applying SymbolicToM to input formats other than ToMi is not trivial, given that it is specifically tailored to ToMi’s input format.

Target Models Since PercepToM leverages the perception reasoning capability of LLMs, we choose models that show reasonable performance on the perception inference task. Specifically, among the eight models, we exclude the bottom two in terms of perception inference accuracy on Percept-FANToM and Percept-ToMi, which are GPT-3.5 Turbo, Claude 3 Haiku, and Gemini 1.0 Pro. As a result, we apply our PercepToM framework to GPT-4 Turbo, GPT-4o, Claude 3 Sonnet, Llama-3 70B Instruct, and Mixtral 8x22B.

4.3 Results

Table 1 shows that the PercepToM improves overall ToM performance when applied to different LLMs on ToMi and FANToM. Remarkably, GPT-4 Turbo achieves 1.0, a perfect score, on the false belief scenario in ToMi. PercepToM generally outperforms CoT and S2A, suggesting that the perspective context extraction, grounded in LLMs’ perception inference results, is more effective than either CoT reasoning or relevant context extraction in S2A at enhancing LLMs’ ToM reasoning.

PercepToM’s performance improvement is more pronounced in false belief scenarios than in true belief scenarios, likely because only minor parts of the contexts are filtered out during the perspective context extraction in the latter. In the false belief scenario of FANToM, which is recognized as the most complex task, all LLMs equipped with PercepToM achieve the highest performance by a large margin. For instance, Llama-3 70B Instruct achieves 0.147 when its vanilla performance is close to 0.

The performance of PercepToM is also compared with that of SymbolicToM (Sclar et al., 2023) on ToMi (Appendix F). PercepToM performs comparably to SymbolicToM in false belief scenarios

Model	Method	ToMi		FANToM	
		True Belief	False Belief	True Belief	False Belief
GPT-4 Turbo	Vanilla	0.739	0.780	0.096	0.017
	CoT	0.700	0.930	0.066	0.079
	S2A	0.682	0.727	0.015	0.019
	PercepToM	0.824	1.000	0.162	0.306
GPT-4o	Vanilla	0.642	0.904	0.077	0.017
	CoT	0.734	0.987	0.153	0.241
	S2A	0.532	0.933	0.000	0.006
	PercepToM	0.659	0.915	0.117	0.566
Claude 3 Sonnet	Vanilla	0.894	0.277	0.010	0.000
	CoT	0.610	0.880	0.005	0.000
	S2A	0.870	0.354	0.000	0.000
	PercepToM	0.963	0.937	0.035	0.066
Llama-3 70B Inst.	Vanilla	0.454	0.803	0.197	0.006
	CoT	0.644	0.900	0.081	0.046
	S2A	0.410	0.894	0.020	0.037
	PercepToM	0.713	0.744	0.242	0.147
Mixtral 8x22B Inst.	Vanilla	0.743	0.597	0.051	0.015
	CoT	0.567	0.630	0.010	0.007
	S2A	0.750	0.357	0.020	0.007
	PercepToM	0.727	0.964	0.217	0.035

Table 1: PercepToM outperforms the baseline models in most of the scenarios on ToMi and FANToM. Bold indicates the best performance within each language model and scenario (true belief or false belief). Performance comparison between PercepToM and SymbolicToM on ToMi can be found in Appendix F.

across most LLMs. However, in true belief scenarios, SymbolicToM consistently outperforms both PercepToM and PercepToM+Oracle. We speculate that this performance gap arises because SymbolicToM rephrases the ToM questions into simpler reality questions. For example, the ToM question “Where will Bob look for the celery?” gets rephrased into “Where is the celery?” In contrast, PercepToM addresses the ToM questions as is.

4.4 Impact of Irrelevant Information on Perception-to-Belief Inference

We conduct an ablation study on perspective context extraction in PercepToM to demonstrate the impact of irrelevant information on LLMs’ perception-to-belief inference. To remove the impact of LLMs’ perception inference accuracy, we compare their performance on perception-to-belief inference with that of PercepToM+Oracle. Both setups have access to the ground truth perception inference information; however, the PercepToM+Oracle includes

the perspective context extraction step, while the perception-to-belief inference setup does not.

As Table 2 shows, models perform significantly better in the PercepToM+Oracle setup than the perception-to-belief inference setup in most scenarios. This suggests that in the perception-to-belief inference setting, despite the presence of the ground truth perception inference information – which should be a substantial hint – within the context, the inclusion of irrelevant information (e.g., the perception of non-target characters and the context not perceived by the target character) results in suboptimal performance in LLMs. Therefore, we can see LLMs struggle to effectively suppress irrelevant information. This capability, coined ‘*inhibitory control*’ in cognitive science, involves the ability to block out irrelevant stimuli while following a specific cognitive objective (Rothbart and Posner, 1985). Inhibitory control is known to be closely linked to ToM and is considered a crucial component for developing ToM (Carlson and Moses, 2001; Carlson et al., 2002).

5 Related Work

Benchmarks for LLM’s Theory of Mind There has been a growing number of benchmarks aimed to evaluate LLM’s theory of mind (ToM), including ToMi (Le et al., 2019), FANToM (Kim et al., 2023), BigToM (Gandhi et al., 2023), HI-TOM (Wu et al., 2023), ToMChallenges (Ma et al., 2023a), Adv-CSFB (Shapira et al., 2024), and OpenToM (Xu et al., 2024). Most of them adopt the false belief test (Wimmer and Perner, 1983), a famous psychology test developed to assess human ToM capabilities. These benchmarks present scenarios involving a character who holds a false belief about a situation (e.g., not knowing something has changed). Models are then asked to predict the character’s thoughts or actions based on the false belief in the scenario. Many benchmarks also include control scenarios where characters do not hold false belief (i.e., true belief scenarios) – situations where their belief about the world state matches the actual state (Le et al., 2019; Kim et al., 2023; Gandhi et al., 2023; Shapira et al., 2024).

Unlike existing benchmarks that primarily measure performance on (downstream) ToM questions themselves, we aim to inquire into the underlying reasoning abilities of LLM’s theory of mind by examining the precursor of ToM: the concept of *seeing leads to knowing* (Baron-Cohen and Goodhart,

Model	Method	ToMi		FANToM	
		True Belief	False Belief	True Belief	False Belief
GPT-4 Turbo	Perception-to-Belief	0.980	0.723	0.138	0.028
	PercepToM+Oracle	0.885	0.993	0.270	0.336
GPT-4o	Perception-to-Belief	0.854	0.863	0.020	0.006
	PercepToM+Oracle	0.660	0.993	0.102	0.571
Claude 3 Sonnet	Perception-to-Belief	0.970	0.384	0.010	0.009
	PercepToM+Oracle	0.987	0.987	0.031	0.058
Llama-3 70B Inst.	Perception-to-Belief	0.810	0.320	0.092	0.020
	PercepToM+Oracle	0.677	0.980	0.133	0.161
Mixtral 8x22B Inst.	Perception-to-Belief	0.894	0.607	0.010	0.045
	PercepToM+Oracle	0.757	0.970	0.224	0.039

Table 2: LLMs perform significantly better in PercepToM+Oracle than perception-to-belief inference across most scenarios. Since the only difference between the two tasks is the inclusion of a perspective context extraction step in PercepToM+Oracle, the result suggests that LLMs struggle to suppress irrelevant perception information when solving ToM questions.

1994; Pratt and Bryant, 1990). We expand existing datasets to identify the perception inference and perception-to-belief inference capabilities, which are essential for ToM reasoning.

Improving LLM’s Theory of Mind Previous research has explored several methods to enhance LLM’s ToM ability. SymbolicToM (Sclar et al., 2023) tracks multiple characters’ beliefs using graphical representation to provide LLMs the context in the target character’s point of view. However, the necessity to construct the belief state graph restricts its adaptability in complex scenarios involving diverse relationships and interactions between entities. SimToM (Wilf et al., 2023) improves LLM’s ToM ability through prompt tuning and highlights the significance of perspective-taking. ToM-LM (Tang and Belle, 2024) improves performance through LLM fine-tuning, while it requires additional training resources.

While SymbolicToM and SimToM achieve high performance on ToMi, their algorithms are tailored for the ToMi structure. For example, SimToM’s

prompts include ToMi-specific hints, such as specific lists of events a character should be aware of for ToMi and instructions for output to depend on a fixed pattern in BigToM stories. As a result, SimToM is mainly tested on ToMi and a specific subset of BigToM questions that resemble those in ToMi. Similarly, the graph construction pipeline in SymbolicToM is designed specifically for ToMi. In contrast, our method, PercepToM, demonstrates flexibility and effectiveness when applied to data of varying formats, leading to significantly better generalizability compared to the existing methods. Although we do not include these two baselines in our main experiments, we present performance comparisons between PercepToM and these methods in Appendices F and G.

6 Conclusion

Inspired by the psychology literature, we evaluated the precursory inferences for human theory of mind (ToM) in large language models (LLM), aiming to broaden our insight into their ToM capabilities. To this end, we constructed perception-augmented ToM benchmarks, Percept-ToMi and Percept-FANToM, by annotating character perceptions about the contexts. Through evaluations and analyses of eight state-of-the-art LLMs, we found that they perform reasonably well in inferring others' perceptions but struggle with inferring others' beliefs based on that perceptual information. Based on these findings, we proposed a new framework, PercepToM, to improve LLM's ToM reasoning. Our framework leverages LLMs' strength in perception inference and enhances their perception-to-belief inference by extracting the relevant contexts. We expect our work to provide insights and enable further in-depth studies into the extent of LLMs' ToM capabilities and targeted improvements in their weaknesses.

7 Limitations

In this paper, we conduct experiments using only two text-based ToM datasets. While ToM tests in psychology involve visual stimuli (e.g., puppets or image strips), our evaluation of ToM abilities relies on text, requiring the ability to read and understand language. As a result, our models must possess robust language comprehension abilities. Moving forward, we are considering expanding our research to include visual ToM and multimodal ToM evaluations, exploring beyond text-based LLMs.

We compare LLMs' ToM performances between true belief and false belief scenarios, but not those between the different orders of ToM questions (e.g., first-order and second-order). Since higher-order ToM requires more inference steps, it will also be interesting to examine the differences in model behavior and capability in solving different orders of ToM questions in future work.

We analyze the precursory inferences for ToM in state-of-the-art large language models (LLMs) that are trained with the full conventional pipeline – i.e., pretraining, instruction tuning, and preference tuning. To understand whether LLMs follow developmental stages akin to human cognition, it is crucial to conduct experiments across the training phases of LLMs. This would include investigating at which stage LLM's social reasoning abilities emerge. These assessments will help us understand how the models' development of social reasoning aligns with stages observed in human theory of mind (ToM).

8 Societal and Ethical Considerations

Our use of FANToM dataset is consistent with its intended use, which is evaluation. We have adhered to the licenses of the benchmarks, ToMi and FANToM, in processing them to create our benchmarks, Percept-ToMi and Percept-FANToM. We plan to make our benchmarks publicly available with the license of Attribution-Noncommercial 4.0 International (CC BY-NC 4.0), allowing the sharing and adapting of the material.

Although we are analyzing large language models' (LLM) theory of mind (ToM) capabilities and its perception-related precursors, we emphasize that we do not claim these LLMs have a mind or any form of subjective consciousness. Our focus lies on improving the social reasoning capabilities of these models to help them interact better in real-world social situations.

Acknowledgment

We thank Amazon for their gift in support of research on theory of mind. This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government(MSIT) (No. RS-2024-00443251, Accurate and Safe Multimodal, Multilingual Personalized AI Tutors).

References

- AI@Meta. 2024. [Llama 3 model card](#).
- Chris Baker, Rebecca Saxe, and Joshua Tenenbaum. 2011. Bayesian theory of mind: Modeling joint belief-desire attribution. In *Proceedings of the annual meeting of the cognitive science society*, volume 33. escholarship.org.
- Simon Baron-Cohen and Frances Goodhart. 1994. The ‘seeing-leads-to-knowing’ deficit in autism: The pratt and bryant probe. *British Journal of Developmental Psychology*, 12(3):397–401.
- Simon Baron-Cohen, Alan M. Leslie, and Uta Frith. 1985. Does the autistic child have a theory of mind? *Cognition*, 21(1):37–46.
- S M Carlson and L J Moses. 2001. Individual differences in inhibitory control and children’s theory of mind. *Child Dev.*, 72(4):1032–1053.
- Stephanie M Carlson, Louis J Moses, and Casey Breton. 2002. How specific is the relation between executive function and theory of mind? contributions of inhibitory control and working memory. *Infant Child Dev.*, 11(2):73–92.
- Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. 2023. [Understanding social reasoning in language models with language models](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Gemini-Team. 2024. [Gemini: A family of highly capable multimodal models](#). *Preprint*, arXiv:2312.11805.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L elio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th eophile Gervet, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. 2024. [Mixtral of experts](#). *Preprint*, arXiv:2401.04088.
- Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Bras, Gunhee Kim, Yejin Choi, and Maarten Sap. 2023. [FANToM: A benchmark for stress-testing machine theory of mind in interactions](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14397–14413, Singapore. Association for Computational Linguistics.
- Matthew Le, Y-Lan Boureau, and Maximilian Nickel. 2019. [Revisiting the evaluation of theory of mind through question answering](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5872–5877, Hong Kong, China. Association for Computational Linguistics.
- Xiaomeng Ma, Lingyu Gao, and Qihui Xu. 2023a. [ToM-Challenges: A principle-guided dataset and diverse evaluation tasks for exploring theory of mind](#). In *Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL)*, pages 15–26, Singapore. Association for Computational Linguistics.
- Ziqiao Ma, Jacob Sansom, Run Peng, and Joyce Chai. 2023b. [Towards a holistic landscape of situated theory of mind in large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1011–1031, Singapore. Association for Computational Linguistics.
- OpenAI. 2023. [Openai json mode](#).
- Chris Pratt and Peter Bryant. 1990. [Young children understand that looking leads to knowing \(so long as they are looking into a single barrel\)](#). *Child Development*, 61(4):973–982.
- David Premack and Guy Woodruff. 1978. Does the chimpanzee have a theory of mind? *Behavioral and brain sciences*, 1(4):515–526.
- Hannes Rakoczy. 2022. [Foundations of theory of mind and its development in early childhood](#). *Nature Reviews Psychology*, 1(4):223–235.
- Mary K Rothbart and Michael I Posner. 1985. Temperament and the development of self-regulation. In *The neuropsychology of individual differences: A developmental perspective*, pages 93–123. Springer.
- Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr, Yejin Choi, and Yulia Tsvetkov. 2023. [Minding language models’ \(lack of\) theory of mind: A plug-and-play multi-character belief tracker](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13960–13980, Toronto, Canada. Association for Computational Linguistics.
- Natalie Shapira, Mosh Levy, Seyed Hossein Alavi, Xuhui Zhou, Yejin Choi, Yoav Goldberg, Maarten Sap, and Vered Shwartz. 2024. [Clever hans or neural theory of mind? stress testing social reasoning in large language models](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2257–2273, St. Julian’s, Malta. Association for Computational Linguistics.
- Winnie Street, John Oliver Siy, Geoff Keeling, Adrien Baranes, Benjamin Barnett, Michael McKibben, Tatenda Kanyere, Alison Lentz, Blaise Aguera y Arcas, and Robin I. M. Dunbar. 2024. [Llms achieve adult human performance on higher-order theory of mind tasks](#). *Preprint*, arXiv:2405.18870.
- Weizhi Tang and Vaishak Belle. 2024. [Tom-lm: Delegating theory of mind reasoning to external symbolic executors in large language models](#). *Preprint*, arXiv:2404.15515.

- Sean Trott, Cameron Jones, Tyler Chang, James Michaelov, and Benjamin Bergen. 2023. [Do large language models know what humans know?](#) *Preprint*, arXiv:2209.01515.
- Max van Duijn, Bram van Dijk, Tom Kouwenhoven, Werner de Valk, Marco Spruit, and Peter van der Putten. 2023. [Theory of mind in large language models: Examining performance of 11 state-of-the-art models vs. children aged 7-10 on advanced tests](#). In *Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL)*, pages 389–402, Singapore. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. [Chain of thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*.
- Jason Weston and Sainbayar Sukhbaatar. 2023. [System 2 attention \(is something you might need too\)](#). *Preprint*, arXiv:2311.11829.
- Oliver Whang. 2023. [Can a machine know that we know what it knows?](#) *The New York Times*.
- Alex Wilf, Sihyun Shawn Lee, Paul Pu Liang, and Louis-Philippe Morency. 2023. [Think twice: Perspective-taking improves large language models’ theory-of-mind capabilities](#). *Preprint*, arXiv:2311.10227.
- Heinz Wimmer and Josef Perner. 1983. [Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception](#). *Cognition*, 13(1):103–128.
- Yufan Wu, Yinghui He, 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](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10691–10706, Singapore. Association for Computational Linguistics.
- Hainiu Xu, Runcong Zhao, Lixing Zhu, Jinhua Du, and Yulan He. 2024. [Opentom: A comprehensive benchmark for evaluating theory-of-mind reasoning capabilities of large language models](#). *arXiv preprint arXiv:2402.06044*.
- Pei Zhou, Aman Madaan, Srividya Pranavi Potharaju, Aditya Gupta, Kevin R. McKee, Ari Holtzman, Jay Pujara, Xiang Ren, Swaroop Mishra, Aida Nematzadeh, Shyam Upadhyay, and Manaal Faruqui. 2023. [How far are large language models from agents with theory-of-mind?](#) *Preprint*, arXiv:2310.03051.

A Details of Perception-Augmented ToM Benchmarks

A.1 Manual Verification of Perception Annotation in Percept-ToMi

SymbolicToM identifies the perceivers of a scene in ToMi deterministically based on a graphical representation of the world state constructed from a templated description of the scene. However, since it determines the perceiver of a scene as all entities in the same connected component based on the world graph, it produces two types of errors in the perceiver information, as shown in Table 3. By reviewing 50 samples of SymbolicToM output, we discover these error types and correct every occurrence of them throughout our dataset.

The perceiver of distractor sentences in ToMi, which describe a character’s opinion about an object, should be the character themselves, since they do not express their opinion to others in the scenario. However, the SymbolicToM-generated perceiver annotation includes other characters located in the same space. We therefore correct the perceivers for all such distractor sentences. Another error in perceiver annotations occur in the sentences preceding the location-disambiguating sentence, which specifies object locations, where their perceivers are annotated with ‘none.’ We align the perceiver annotations of these sentences with those of the subsequent location-disambiguating sentence, since they are always paired and have the same perceivers.

A.2 Perception Annotation Criteria for Percept-FANToM

The following criteria are used to determine the joining and leaving times of a character in a conversation within the FANToM dataset.

- When a character joins a conversation is determined by the moment the character directly participates in the conversation. If a character enters with an utterance like “you guys are having an interesting conversation,” we consider him/her a perceiver from the moment he/she starts speaking, as the exact point when the character began listening is unclear.
- When a character leaves the conversation is determined by the final farewell utterance. Even if a character disappears mid-utterance (e.g., C: “Bye, A. So, B, what do you think?”), the entire

utterance is still considered as perceived by the departing character.

B Dataset Statistics

Table 4 presents a comparison of data statistics between perception-augmented ToM benchmarks and their corresponding source benchmarks.

C Prompt Examples

This section introduces prompt examples to evaluate perception inference and perception-to-belief inference.

C.1 Perception Inference

The following two boxes are prompt examples using Percept-ToMi and Percept-FANToM, respectively. Some parts are omitted because of the space limit.

```
Story: Ella likes the suit. Ella entered the cellar. Lucas entered the cellar. Benjamin entered the porch. The boots is in the cupboard. The cupboard is in the cellar. Lucas exited the cellar. Benjamin exited the porch. Ella likes the sweatshirt. Lucas entered the porch. Ella moved the boots to the pantry. The pantry is in the cellar.
```

```
Create a JSON array consisting of JSON objects. Each object should contain a sentence from the story and the perceivers of the scene described in that sentence. Assume that characters in the story can perceive every scene occurring in their location but not scenes occurring elsewhere. Also, include the actant of any action as a perceiver of that action. Provide only a JSON array in the following format. Do not include any explanation. [{"Noah exited the living room.": ["Noah", "Emma"]}].
```

```
Gianna: Guys, I've really enjoyed sharing our pet stories, but I need to excuse myself. I need to change clothes for a meeting later. Talk to you later!  
Sara: Sure thing, Gianna. Take care!  
Javier: Catch you later, Gianna.  
Sara: So Javier, have you ever tried training Bruno?  
Javier: Yes, I did actually. It was a challenge at times, but rewarding nevertheless. How about you? Did you try training Snowflake?
```

```
...  
Gianna: Hey guys, I'm back, couldn't miss out on more pet stories. Speaking of teaching and training pets, it is amazing how that further strengthens the bond between us and our pets, right?  
...  
Create a JSON array consisting of JSON objects. Each object should include an utterance from the dialogue and the audience for that utterance. Assume that characters in the story can hear every utterance that occurs while they are involved in the dialogue, but not those that occur when they are absent. Also, ensure that the speaker of each utterance is included in the audience. Provide only the JSON array in the following format. Do not include any explanations. [{"Noah: Hi, Emma.": ["Noah", "Emma"]}].
```

C.2 Perception-to-Belief Inference

The following two boxes are prompt examples using Percept-ToMi and Percept-FANToM, respectively. Some parts are omitted because of the space limit.

Sentence Type	Information	SymbolicToM Output	Final Annotation
Object Location	The slacks is in the pantry.	None	Ella, Benjamin
	The pantry is in the master bedroom.	Ella, Benjamin	Ella, Benjamin
Distractor	Olivia loves the skirt.	Olivia, James, Lily	Olivia

Table 3: The example perceiver annotations in ToMi corrected by manual verification.

Datasets	True Belief					False Belief				
	# Ctx.	# Q.	Avg. # Q. per Ctx.			# Ctx.	# Q.	Avg. # Q. per Ctx.		
			P.I.	PtoB	ToM			P.I.	PtoB	ToM
Disambiguated ToMi - Test Set	2793	2793	-	-	1.0	1210	1210	-	-	1.0
Percept-ToMi	300	1802	4.0	1.0	1.0	300	1804	4.0	1.0	1.0
FANToM	402	3432	-	-	8.5	642	8530	-	-	13.3
Percept-FANToM	340	13161	24.2	7.3	7.3	539	25279	23.3	11.8	11.8

Table 4: Comparison of dataset statistics between the source datasets and our proposed datasets. Number of contexts, the total number of questions, and the average number of questions per context for each task are compared. Our benchmarks expand upon the sampled contexts from the source datasets by incorporating two precursory inference tasks for Theory of Mind (ToM)—perception inference and perception-to-belief inference.

Each JSON object in the following list contains the description of a consecutive scene in a story and its perceivers.

```
[{"Ella likes the suit": ["Ella"]},
{"Ella entered the cellar": ["Ella"]},
{"Lucas entered the cellar": ["Lucas", "Ella"]},
{"Benjamin entered the porch": ["Benjamin"]},
{"The boots is in the cupboard": ["Ella", "Lucas"]},
{"The cupboard is in the cellar": ["Ella", "Lucas"]},
{"Lucas exited the cellar": ["Lucas", "Ella"]},
{"Benjamin exited the porch": ["Benjamin"]},
{"Ella likes the sweatshirt": ["Ella"]},
{"Lucas entered the porch": ["Lucas"]},
{"Ella moved the boots to the pantry": ["Ella"]},
{"The pantry is in the cellar": ["Ella"]}]
```

Question: Where will Lucas look for the boots? State the most detailed position possible. (e.g., in A in B) Answer in one sentence without explanation.
Answer:

Each JSON object in the following list contains consecutive utterances in a dialogue and its audiences.

```
[{"Gianna: Guys, I've really enjoyed sharing our pet stories, but I need to excuse myself. I need to change clothes for a meeting later. Talk to you later!": ["Gianna", "Sara", "Javier"]},
{"Sara: Sure thing, Gianna. Take care!": ["Sara", "Gianna"]},
{"Javier: Catch you later, Gianna.": ["Javier", "Gianna"]},
{"Sara: So Javier, have you ever tried training Bruno?": ["Sara", "Javier"]},
{"Javier: Yes, I did actually. It was a challenge at times, but rewarding nevertheless. How about you? Did you try training Snowflake?": ["Javier", "Sara"]},
...
{"Gianna: Hey guys, I'm back, couldn't miss out on more pet stories. Speaking of teaching and training pets, it is amazing how that further strengthens the bond between us and our pets, right?": ["Gianna", "Sara", "Javier"]},
...]
```

Target: Who discussed their experiences training their pets, Bruno and Snowflake?
Question: Does Javier know the precise correct answer to this question? Answer yes or no.
Answer:

D Input and Output Examples of PercepToM Pipeline

This section presents examples of input prompts and intermediate outputs of PercepToM steps. Note that PercepToM consists of three steps: perception inference, perspective context extraction, and reading comprehension.

First, the following two boxes are prompts for character perception inference on ToMi and FANToM, respectively.

[Input Prompt]: Story: Ella likes the suit. Ella entered the cellar. Lucas entered the cellar. Benjamin entered the porch. The boots is in the cupboard. The cupboard is in the cellar. Lucas exited the cellar. Benjamin exited the porch. Ella likes the sweatshirt. Lucas entered the porch. Ella moved the boots to the pantry. The pantry is in the cellar.

Create a JSON array consisting of JSON objects. Each object should contain a sentence from the story and the perceivers of the scene described in that sentence. Assume that characters in the story can perceive every scene occurring in their location but not scenes occurring elsewhere. Also, include the actant of any action as a perceiver of that action.
Provide only a JSON array in the following format. Do not include any explanation.
[{"Noah exited the living room.": ["Noah", "Emma"],}]

[Input Prompt]: Gianna: Guys, I've really enjoyed sharing our pet stories, but I need to excuse myself. I need to change clothes for a meeting later. Talk to you later!
Sara: Sure thing, Gianna. Take care!
Javier: Catch you later, Gianna.
Sara: So Javier, have you ever tried training Bruno?
Javier: Yes, I did actually. It was a challenge at times, but rewarding nevertheless. How about you? Did you try training Snowflake?
Sara: Oh gosh, trying to train a cat is a whole different ball game. But I did manage to teach her a few commands and tricks. She was quite an intelligent little furball.
Gianna: Hey guys, I'm back, couldn't miss out on more pet

stories. Speaking of teaching and training pets, it is amazing how that further strengthens the bond between us and our pets, right?
...

Create a JSON array consisting of JSON objects. Each object should contain a sentence from the story and the perceivers of the scene described in that sentence. Assume that characters in the story can perceive every scene occurring in their location but not scenes occurring elsewhere. Also, include the actant of any action as a perceiver of that action. Provide only a JSON array in the following format. Do not include any explanation.
[{"Noah exited the living room.": ["Noah", "Emma"]}]]

Here are the outputs of perception inference, which will be the input for perspective context extraction.

```
[Out]: [{"Ella likes the suit.": ["Ella"],
"Ella entered the cellar.": ["Ella"],
"Lucas entered the cellar.": ["Ella", "Lucas"],
"Benjamin entered the porch.": ["Benjamin"],
"The boots is in the cupboard.": ["Ella", "Lucas"],
"The cupboard is in the cellar.": ["Ella", "Lucas"],
"Lucas exited the cellar.": ["Lucas"],
"Benjamin exited the porch.": ["Benjamin"],
"Ella likes the sweatshirt.": ["Ella"],
"Lucas entered the porch.": ["Lucas"],
"Ella moved the boots to the pantry.": ["Ella"],
"The pantry is in the cellar.": ["Ella"]}]]
```

```
[Out]: [{"Gianna: Guys, I've really enjoyed sharing our pet stories, but I need to excuse myself. I need to change clothes for a meeting later. Talk to you later!": ["Gianna", "Sara", "Javier"]},
{"Sara: Sure thing, Gianna. Take care!": ["Sara", "Gianna"]},
{"Javier: Catch you later, Gianna.": ["Javier", "Gianna"]},
{"Sara: So Javier, have you ever tried training Bruno?": ["Sara", "Javier"]},
{"Javier: Yes, I did actually. It was a challenge at times, but rewarding nevertheless. How about you? Did you try training Snowflake?": ["Javier", "Sara"]},
{"Sara: Oh gosh, trying to train a cat is a whole different ball game. But I did manage to teach her a few commands and tricks. She was quite an intelligent little furball.": ["Sara", "Javier"]},
{"Gianna: Hey guys, I'm back, couldn't miss out on more pet stories. Speaking of teaching and training pets, it is amazing how that further strengthens the bond between us and our pets, right?": ["Gianna", "Sara", "Javier"]}]]
```

The perspective context extraction selects the subset of context perceived by the target character. The outputs will be as follows:

```
[Out]: Lucas entered the cellar. The boots is in the cupboard. The cupboard is in the cellar. Lucas exited the cellar. Lucas entered the porch.
```

```
[Out]: Gianna: Guys, I've really enjoyed sharing our pet stories, but I need to excuse myself. I need to change clothes for a meeting later. Talk to you later!
Sara: Sure thing, Gianna. Take care!
Javier: Catch you later, Gianna.
Gianna: Hey guys, I'm back, couldn't miss out on more pet stories. Speaking of teaching and training pets, it is amazing how that further strengthens the bond between us and our pets, right?
...
```

Lastly, based on the extracted perspective contexts, we build prompts to answer the ToM question.

```
[Input Prompt]: Here are the past scenes in sequence that Lucas knows about.
```

```
Lucas entered the cellar. The boots is in the cupboard. The cupboard is in the cellar. Lucas exited the cellar. Lucas entered the porch.
```

```
Question: Where will Lucas look for the boots? State the most detailed position possible (e.g., in A in B). Answer in one sentence without explanation.
Answer:
```

```
[Input Prompt]: Here are the past utterances in sequence that Gianna is aware of.
```

```
Gianna: Guys, I've really enjoyed sharing our pet stories, but I need to excuse myself. I need to change clothes for a meeting later. Talk to you later!
Sara: Sure thing, Gianna. Take care!
Javier: Catch you later, Gianna.
Gianna: Hey guys, I'm back, couldn't miss out on more pet stories. Speaking of teaching and training pets, it is amazing how that further strengthens the bond between us and our pets, right?
...
```

```
Question: What does Gianna believe about who discussed their experiences training their pets, Bruno and Snowflake? Choose between (a) and (b). Do not include any explanation.
(a) Gianna believes that Sara and Javier discussed their experiences training their pets, Bruno and Snowflake.
(b) Gianna knows that Javier discussed training his pet, Bruno. However, Gianna will not know training a pet named Snowflake.
```

E LLM Performances on Percept-ToMi and Percept-FANToM

Table 5 presents the exact performance of Percept-ToMi and Percept-FANToM in perception inference, perception-to-belief inference, and ToM, which is also depicted in Figure 3.

F Performance Comparison Between PercepToM and SymbolicToM

Table 6 shows the performances of PercepToM, PercepToM+Oracle, and SymbolicToM on ToMi.

G Performance Comparison Between PercepToM and SimToM

Table 7 shows the performances of PercepToM and SimToM on ToMi.

Dataset	Model	True Belief			False Belief		
		Perception	Perception-to-Belief	ToM	Perception	Perception-to-Belief	ToM
Percept-ToMi	GPT-3.5 Turbo	0.228	0.824	0.792	0.585	0.432	0.237
	GPT-4 Turbo	0.934	0.980	0.739	0.950	0.723	0.780
	GPT-4o	0.903	0.854	0.642	0.925	0.863	0.904
	Claude 3 Haiku	0.874	0.480	0.730	0.798	0.724	0.290
	Claude 3 Sonnet	0.886	0.970	0.894	0.886	0.384	0.277
	Gemini 1.0 Pro	0.425	0.850	0.690	0.733	0.104	0.127
	Llama-3 70B Instruct	0.814	0.810	0.454	0.718	0.320	0.803
	Mixtral 8x22B Instruct	0.920	0.894	0.743	0.917	0.607	0.597
Percept-FANToM	GPT-3.5 Turbo	0.866	0.505	0.177	0.877	0.000	0.000
	GPT-4 Turbo	0.962	0.138	0.096	0.970	0.028	0.017
	GPT-4o	0.970	0.020	0.077	0.977	0.006	0.017
	Claude 3 Haiku	0.792	0.015	0.025	0.806	0.009	0.002
	Claude 3 Sonnet	0.974	0.010	0.010	0.977	0.009	0.000
	Gemini 1.0 Pro	0.937	0.000	0.000	0.950	0.002	0.000
	Llama-3 70B Instruct	0.982	0.092	0.197	0.980	0.020	0.006
	Mixtral 8x22B Instruct	0.899	0.010	0.051	0.892	0.045	0.015

Table 5: LLM performances for perception inference, perception-to-belief inference, and Theory of Mind (ToM), as illustrated in Figure 3 for Percept-ToMi and Percept-FANToM.

Model	Method	True Belief	False Belief
GPT-4 Turbo	PercepToM	0.824	1.000
	PercepToM+Oracle	0.885	0.993
	SymbolicToM	0.997	0.977
GPT-4o	PercepToM	0.659	0.915
	PercepToM+Oracle	0.660	0.993
	SymbolicToM	1.000	0.977
Claude 3 Sonnet	PercepToM	0.963	0.937
	PercepToM+Oracle	0.987	0.987
	SymbolicToM	1.000	0.977
Llama-3 70B Inst.	PercepToM	0.713	0.744
	PercepToM+Oracle	0.677	0.980
	SymbolicToM	1.000	0.977
Mixtral 8x22B Inst.	PercepToM	0.727	0.964
	PercepToM+Oracle	0.757	0.970
	SymbolicToM	1.000	0.977

Table 6: Performance comparison of PercepToM, PercepToM+Oracle, and SymbolicToM on the ToMi dataset. PercepToM+Oracle and PercepToM show comparable performance to SymbolicToM in false belief scenarios across most models. In true belief scenarios, SymbolicToM consistently outperforms PercepToM+Oracle, likely due to its question rephrasing process.

Model	Method	True Belief	False Belief
GPT-4 Turbo	SimToM	0.657	0.873
	PercepToM	0.824	1.000
GPT-4o	SimToM	0.797	0.450
	PercepToM	0.659	0.915
Llama-3 70B Inst.	SimToM	0.644	0.770
	PercepToM	0.713	0.744
Mixtral 8x22B Inst.	SimToM	0.677	0.660
	PercepToM	0.727	0.964

Table 7: Performance comparison between SimToM and PercepToM on Fixed and Disambiguated ToMi (Sclar et al., 2023). Overall, PercepToM shows more robust performance across different models in both of true and false belief scenarios.