SIMPLETOM: EXPOSING THE GAP BETWEEN EXPLICIT TOM INFERENCE AND IMPLICIT TOM APPLICATION IN LLMS

Anonymous authors

Paper under double-blind review

ABSTRACT

While prior work has explored whether large language models (LLMs) possess a "theory of mind" (ToM) - the ability to attribute mental states to oneself and others - there has been little work testing whether LLMs can implicitly apply such knowledge to predict behavior, or to judge whether an observed behavior is rational. Such skills are critical for appropriate interaction in social environments. Our approach to study such capabilities is to create a new dataset, called SimpleToM, containing concise, diverse stories (e.g., "The can of Pringles has moldy chips in it. Mary picks up the can in the supermarket and walks to the cashier."), each with three questions that test different degrees of ToM reasoning, asking models to predict (a) mental state ("Is Mary aware of the mold?"), (b) behavior ("Will Mary pay for the chips or report the mold?"), and (c) judgment ("Mary paid for the chips. Was that reasonable?"). To our knowledge, SimpleToM is the first dataset to systematically explore downstream reasoning requiring knowledge of mental states in realistic scenarios. Our experimental results are intriguing: While most models can reliably predict mental state on our dataset (a), they often fail to correctly predict the behavior (b), and fare even worse at judging whether given behaviors are reasonable (c), despite being correctly aware of the protagonist's mental state should make such secondary predictions obvious. We further show that we can help models do better at (b) and (c) via interventions such as reminding the model of its earlier mental state answer and mental-state-specific chain-of-thought prompting, raising the action prediction accuracies (e.g., from 49.5% to 93.5% for GPT-40) and judgment accuracies (e.g., from 15.3% to 94.7% in GPT-4o). However, while this shows that models can be coaxed to perform well, it requires task-specific interventions, and the natural model performances remain low, a cautionary tale for LLM deployment. SimpleToM thus breaks new ground in probing real-world ToM reasoning, and reveals surprising, new insights about current model capabilities. We hope the dataset enables further exploration by the community into this critical area of model behavior.

039

041

037

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

034

040 1 INTRODUCTION

042 People infer what others know, anticipate their actions, and expect them to choose cost-minimizing 043 behaviors (Gergely & Csibra, 2003; Liu & Spelke, 2017). Performing such social reasoning in-044 volves attributing mental states to oneself and others, an ability widely known as Theory of Mind (ToM) (Premack & Woodruff, 1978). ToM has been extensively studied in psychology in a range of scenarios, for instance, studies of manipulation, secrecy, deception, lying, and misleading behavior 046 (Doherty, 2008). Examples of classical tests in developmental psychology include the unexpected 047 transfer false belief task, e.g., the Sally-Anne task (Baron-Cohen et al., 1985), and the unexpected 048 contents false belief task, e.g., the Smarties task (Perner et al., 1987). Given the increasing use of LLMs in human interactions and as decision-making agents within complex, human-centered environments, it is crucial to assess their ToM capabilities. 051

Current datasets for assessing ToM in LLMs are limited by their reliance on the classical SallyAnne task or templated variants of it (Le et al., 2019; Nematzadeh et al., 2018; Wu et al., 2023; Xu et al., 2024). These datasets have significant shortcomings: (i) limited diversity in how information



Figure 1: To allow for a nuanced analysis of models' neural ToM abilities, SimpleToM covers both explicit ToM (a) and applied ToM (b, c) question types. SimpleToM measures the ability of LLMs to (a) infer the character's mental state, specifically information awareness, (b) anticipate their likely next behavior in the given situation, and (c) make appropriate judgment of the character's behavior that correctly accounts for their mental state.

asymmetry arises (see related work in Section 6 for examples across existing datasets), (ii) explicit
 use of percept and mentalizing verbs like "sees" and "thinks" which serve as trigger words for
 models to realize that these are important aspects, removing the need for implicit commonsense
 inferences about relevant percepts or beliefs, and (iii) limited exploration of applied ToM, such as
 judgment of behavior which requires implicit reasoning about mental state.

The goal of SimpleToM is to go beyond the focus on mental state inference, to also assess how well models can predict others' behavior based on that understanding, and make judgments of appropriateness of behavior. Each story in SimpleToM is paired with three types of questions targeting these abilities (Figure 1), with a total of 1147 stories and 3441 questions in daily life settings. Our analysis using SimpleToM exposes a significant gap in model performance between *explicit* and *applied* ToM questions (Lee et al., 2024). We find that frontier models perform well on explicit ToM questions (directly querying for information about "mental state", i.e., information awareness). However, this success does not extend to applied ToM ("behavior" and "judgment" questions). Scenario-specific analysis highlights the need for diverse testing scenarios to rigorously assess neural ToM and indi-cates areas for improvement even in strong models like o1-preview. Our experiments with various inference-time interventions show that the performance of applied ToM improves with increased guidance. However, the need for highly specific guidance and the high inference cost (e.g., CoT using hundreds of tokens) required for models to reason through simple two-sentence stories in SimpleToM, suggests that the inherent ToM reasoning abilities of advanced LLMs still have sig-



121

122

123

128

129

130

131

132

133

134

135

136

Figure 2: We leverage the generative strength of language models to obtain concise stories with varied entities and diverse situations, suitable for testing different levels of ToM reasoning. The generated stories (and answer options) were then rigorously filtered by careful human annotators who passed a strict qualification test. The result is a high-quality and diverse dataset, SimpleToM.

nificant room for improvement. The results showing that frontier models still lack the ability to
 independently apply ToM skills in downstream tasks—such as anticipating others' behavior and
 making judgments—call for caution when using them in real-world applications (see discussion of
 example applications in Appendix A).

Our contributions and findings are:

- A new dataset, SimpleToM, for testing core abilities of LLMs in explicit and applied ToM.
- Current frontier models have decoupled capabilities between predicting someone's information awareness in a situation (explicit ToM, which they excel at), and utilizing it to predict and judge someone's behavior (applied ToM, which they perform poorly at).
- Models can be explicitly guided to perform applied ToM, e.g., through custom chain-ofthought prompting, but this is a fragile workaround, and being competent on SimpleToM without such guidance is thus an important measure of a model's capabilities.

We make our SimpleToM dataset and the full evaluation data for the analyzed models publicly available at [to-be-released-URL-if-accepted]. This will allow researchers to build on top of our work in studying the neural ToM capabilities of LLMs in general, as well as further exploiting the diversity of situations covered by SimpleToM in exploring applied ToM.

141 142

143 144

145 146 2 SIMPLETOM DESIGN

We design the stories in SimpleToM to contain diverse types of information asymmetry, using a concise format and associated with specific question types testing explicit and applied ToM.

147 2.1 DIFFERENT TYPES OF INFORMATION ASYMMETRY

To expand beyond the classical false belief task, we seed the creation of SimpleToM with ten diverse scenarios where information asymmetry occurs naturally (Table 1). This is inspired by social psychology literature to cover asymmetries like manipulation, deception, secrecy, lying, and misleading behavior (Doherty, 2008), seen in real-world contexts like sales of "lemon" products, where items with hidden flaws are purchased due to a lack of information (Akerlof, 1978). These are under-examined in existing ToM tests. We further describe the scenarios with examples in Table 7 (Appendix F.1).

- 155
- 156 157

2.2 SIMPLE STORY FORMAT WITHOUT EXPLICIT PERCEPTS OR MENTAL STATES

The SimpleToM example story from Figure 1 reads: *The can of Pringles has moldy chips in it. Mary picks up the can in the supermarket and walks to the cashier.* Each story has exactly two sentences, where the first sentence introduces a key information about something (Object/Person/Action Z), while the second sentence presents the main subject of the story (Person X) doing something with Object/Person/Action Z while being unaware of the key information. The list of story elements are:

162	Scenario	Reason for information asymmetry
163	food item in grocery store	Food items bought in grocery stores cannot be closely examined for their quality before checking
164		out
165	provider into healthcare true property pretentious labels	Efficacy of healthcare products cannot be closely examined or verified before purchase Subtle properties of products cannot be closely examined or verified
166	behind the scene service industry	Questionable behind-the-scenes practices in the service industry are not observed by customers
167	unobserved unethical actions	What is inside labeled (opaque) containers cannot be observed before opening the container Unethical actions not observed are not known
168	inside containers for personal	What is inside (opaque) containers for personal belongings cannot be observed before opening the
169	seller info in second hand market	Hidden flaws in second-hand items bought cannot be observed before the purchase
170	hidden body part feature	Body features hidden under clothing cannot be observed
171		Details in focked devices of accounts cannot be observed by others
172	Table 1: The ten broad scena	arios used to seed the generation of stories in SimpleToM. Each scenario
173	describes a type of informat	ion asymmetry that occurs naturally in the real world.
174	71	
175	 Key Information: 	involves something unexpected which Person X is unlikely to know or
177	perceive, e.g., The	can of Pringles has moldy chips in it.
178	Object/Person/Ac	tion Z: the subject of the key information (e.g., can of Pringles)
179	• Person X: person u	inaware of the key information (e.g., Mary)
180	 Person Y (optional 	l): any other character(s) needed for the story
181	We impose the constraint fl	nat Person X's unawareness of the key information should be implicit
182	(e.g., avoid explicit use of	perception or mentalizing words such as "see", "notice" or "belief").
183	This design encourages mo	dels to read between the lines and make commonsense inferences over
184	the given situations and infe	r characters' mental states in a more realistic manner, bringing us closer
185	to realistic daily life use cas	es of ToM. (E.g., you cannot see through a Pringles can; you would not
186	know about a cheating even	t if you were not present.)
187	To support formulating the l	behavior prediction question (Section 2.3), we also generate options for
188	what might happen next:	
189	Unaware behavior	r: A likely next action by Person X given that they are unaware of the
190	key information	
191	Aware behavior: A	A likely next action by Person X if they were somehow aware of the key
192	information after a	ll (counterfactual)
193		
194	2.3 QUESTIONS TESTING	EXPLICIT AND APPLIED TOM
195	We use three types of quest	tions (Figure 1) to probe a model's grasp of each story, covering both
190	<i>explicit</i> theory of mind (con	centual knowledge about others' mental states: i.e., via (a) mental state
108	question about information	awareness) and <i>applied</i> theory of mind (the ability to use theory of mind
100	in real-life contexts; i.e., via	(b) behavior and (c) judgment questions) (Lee et al., 2024).
200	Mental state (MS) question	about information awareness. We test ability of models to infer men
201	tal states, specifically inform	nation awareness, through a simple ves/no question (Is <person x=""></person>
202	likely to be aware	that " <key information="">"?) To infer whether a character is</key>
203	aware of something in Simp	bleToM stories, a model has to make implicit commonsense inferences
204	about what the character can	n perceive or know in the given situation (including commonsense rea-
205	soning about physical object	ts, space, intent, goals of others, and so on).
206	Behavior prediction quest	ion: This question asks which of two possible actions the main subject
207	(Person X) is likely to perf	orm next. For instance, beyond answering that a person shopping for
208	chips in the supermarket is	unlikely to know that "the can of Pringles has moldy chips in it", a
209	model that successfully appl	lies this inference for behavior prediction should also infer that a person
210	who picked up such a can i	n the supermarket would likely "pay for the chips" rather than "report
211	situation to infer the mental	er these questions correctly, models need to implicitly reason over the
212	the key information would in	state of character, and realize now the character's lack of awareness of
213		
214	Judgment question: The i	udgment question specifies that the "correct" action was taken (rather

than the incorrect one) and asks if this was a reasonable choice. As the inference graph in Figure 1 illustrates, the judgment question goes beyond behavior prediction as it requires two levels of im-

plicit reasoning, first implicitly predicting the behavior of Person X, which itself relies on implicitly understanding their mental state. People's mental states are an important factor to consider in making appropriate judgments of their behavior (Jara-Ettinger et al., 2016; Schein & Gray, 2018). For
instance, buying a can of Pringles that has moldy chips in it is **not a reasonable action** if the person knows about the moldy chips. However, it is a **perfectly reasonable** (and expected) behavior if this piece of key information is not a part of the person's mental state.

222 223 224

227

228

229

230

3

SIMPLETOM CREATION

225 3.1 GENERATING DIVERSE STORIES

Specifically, the construction of SimpleToM consists of the following steps:

- Step 1: Manually create one example seed story for each scenario.
- Step 2: For each scenario, using the seed story as example, prompt the LLM to suggest 10 diverse sets of entities compatible with an information asymmetry. (See prompt in F.5.)
- Sets of entities comparise with an information asymmetry. (see prompt in F.5.)
 Step 3: For each set of suggested entities, along with the seed story, prompt the LLM to write three new stories at different levels of "severity." With each story, also generate likely next "unaware" and "aware" behaviors (see Section 2.2). Appendix F.4 provides further details.

235 We went through two rounds of this process. First, we used GPT-4 and Claude-3-Opus¹ to generate 236 a total of 1200 stories.² After annotating and filtering this initial set (Section 3.2), we picked a 237 new set of top-scoring seed stories and sourced 10 additional sets of entities from each of GPT-238 40 and Claude-3.5-Sonnet. We used these two newer models to generate stories for all 40 sets 239 of entities, for a total of 2400 more stories. By using several generator models, varied entities and different seed stories, the resulting stories in SimpleToM have a wide range of information 240 asymmetries instantiated in different real-world situations, effectively broadening neural ToM tests 241 beyond traditional settings (Section 6). These contexts also allow for nuanced and implicit traits 242 (e.g., buyers would avoid products with defects if they know about them). 243

244 245 3.2 Strict quality control on stories that goes into SimpleToM

We gather human annotations on each story (and unaware/aware next actions). We asked annotators
four questions for each story, summarized in Figure 2. This process verifies that the key information
in each story is something that Person X has false belief about. We also carefully verify that the next
likely "unaware action" is appropriate if and only if Person X is unaware of the key information.
We similarly verify the "aware action," for the counterfactual situation where Person X is somehow aware of the key information. Appendix E provides further details about the crowdsourcing
procedure, with instructions, examples and question templates.

Our annotators passed a rigorous qualification test (Appendix E.2) and met other high-standard
 requirements (Appendix E.3). Only stories for which all crowdworkers (3) judged all aspects to be
 valid were included in SimpleToM.³ This results in 1147 stories (out of the original 3600) in the final
 SimpleToM dataset. Table 8 (Appendix F.2) provides statistics and further details for SimpleToM.

257 258

259

265

266

268

4 EXPERIMENTAL SETUP

We evaluate SimpleToM on ten frontier LLMs from different sources and with different levels of capabilities: Llama-3.1-8B, Llama-3.1-405B, (Dubey et al., 2024) Claude-3-Haiku, Claude-3-Opus (Anthropic, 2024b), Claude-3.5-Sonnet, (Anthropic, 2024a) GPT-3.5, GPT-4, GPT-40, 01-mini, and o1-preview (OpenAI, 2024; 2023) (refer to Appendix D Table 6 for more details). We use the most deterministic setting with a generation temperature of 0.4

We use SimpleToM to investigate the following research questions:

¹See Table 6 for exact models used

²10 scenarios * 2 models * 10 entities per model * 3 severities * 2 models to generate stories

³See more details in Appendix E.4.

⁴Except the o1 models which do not allow adjusting the temperature or system prompt.

270	model	mental state	hebavior	iudament	average
271	model	(Explicit ToM)	(Applied ToM)	(Applied ToM)	average
272	CDT 2.5	26.5	7.6	20.1	24.4
070	GP1-3.5	36.5	7.6	29.1	24.4
213	Claude-3-Haiku	87.2	23.6	16.7	42.5
274	o1-mini	87.8	44.8	27.0	53.2
275	GPT-40	95.6	49.5	15.3	53.5
276	Llama-3.1-405B	97.8	58.2	10.0	55.4
277	Claude-3-Opus	98.3	64.4	9.6	57.4
£11	GPT-4	96.6	63.0	19.5	59.7
278	Llama-3.1-8B	88.1	38.5	54.6	60.4
279	Claude-3.5-Sonnet	97.9	67.0	24.9	63.3
280	o1-preview	95.6	84.1	59.5	79.7

Table 2: Evaluation results for SimpleToM on the different question types. Models are generally proficient in explicit ToM questions (directly querying about "mental state", i.e., information awareness) but this success does not transfer to applied ToM ("behavior" and "judgment" questions).

1. How well can models (i) infer characters' mental states, (ii) anticipate characters' behavior and (iii) make appropriate judgments, requiring the use of ToM inferences?

- 2. How does the ToM performance of models differ across scenarios?
- 3. How can the performance on the behavior and judgment questions improve with various levels intervention during inference, such as patching mental state inference in the prompt (MS remind), guiding the model via system prompt (SysP), and encouraging reasoning/thinking before answering (CoT)?
- 5 RESULTS AND ANALYSIS

281

282

283

284

286

287

289

290

291

293

295 296

297

5.1 FRONTIER LLMS CAN INFER MENTAL STATES, BUT STRUGGLE TO USE IT

The overall evaluation results on SimpleToM for the 10 models are summarized in Table 2, spanning the different question types (as detailed in Section 2.3). We analyze models' performance for each type of question below. Note that these are binary questions where random performance is 50%.

Mental state (MS) question about information awareness: Our results (Table 2, "mental state" 302 column) show that reasoning over implicit information in given situations to infer mental states is 303 still challenging for models like GPT-3.5 (36.5% accuracy), while newer and/or bigger models like 304 Claude-3-Haiku, o1-mini, and Llama-3.1-8B perform reasonably well (around 88%). In fact, all 305 recent frontier models are proficient at inferring characters' awareness in our dataset - GPT-40, 306 Llama-3.1-405B, Claude-3-Opus, GPT-4, Claude-3.5-Sonnet, and o1-preview all achieved accura-307 cies of more than 95%. This result also confirms the quality of our dataset, in that characters' mental 308 states in SimpleToM stories are implicit but reasonably easy to infer, as designed. 309

Behavior prediction: On behavior prediction questions (Table 2, "behavior" column), smaller and 310 older models perform extremely poorly (with GPT-3.5 achieving only 7.6% accuracy and several 311 other models scoring less than 40%). Even for the larger models, like Llama-3.1-405B, Claude-3.5-312 Sonnet, and GPT-40, performance on the behavior prediction task is much worse than on the mental 313 state task with at least a 30% performance drop. This large inconsistency suggests that while fron-314 tier LLMs may have the right conceptual knowledge/information about others' mental states when 315 directly asked, they struggle to apply this knowledge in real-life contexts to make predictions about 316 characters' behavior. Only the latest o1-preview model, with its built-in inference time reasoning 317 tokens, $^{\circ}$ manages a decent score on this question type (84.1%).

Judgment of behavior: Our results (Table 2, "judgment" column) show that this additional, seemingly trivial, inference step (beyond the behavior prediction) makes the task much more difficult for all the models. Even the newest and largest models like Llama-3.1-405B, Claude-3.5-Sonnet,

 ⁵The o1 reasoning tokens make these models more like the chain-of-thought prompted versions of the other
 models below, although without any custom prompt. See Appendix I for discussion of the number of output tokens used by the o1 and the other models when using chain-of-thought.



Figure 3: Considering the sequence of first predicting mental state, then behavior and finally judgment, we can record a failure for the first mistake. The "fail at behavior prediction" and "fail at judgment" parts can be considered inconsistent predictions by the model, since it got the associated mental state (and behavior prediction) questions correct.



Figure 4: Comparing performance for all three question types across select scenarios and models. Each bar represents the overall accuracy. The mental state accuracy is generally near 100%, while behavior prediction and judgment accuracies are often much lower.

357 and GPT-40, which all achieved accuracies of more than 95% on inferring characters' mental state, 358 consistently make inaccurate judgments of behavior, and their performance drop to far below ran-359 dom (with accuracies in the range of 10% to 24.9%). Llama-3.1-8B somehow manages around 360 random performance (at 54.6%), as an outlier. Finally, o1-preview, the best-performing model on 361 behavior prediction, also gets the best score on this judgment prediction task, but even so it reaches 362 only 59.5%, far below its performance on the other questions. These results show that even the latest 363 models lack the ability to effectively apply ToM reasoning in this kind of judgment task. Appendix K provides further discussion on interesting trends. 364

Overall ToM reasoning: To fully reason about the judgment question, requires reasoning about the behavior prediction, which relies on understanding the mental state (information awareness) question. We can visualize this by recording the **first** failure on the three "mental state" \rightarrow "behavior" \rightarrow "judgment" questions. Figure 3 shows the distribution of such failures, showing the proportionally large cases failing at the behavior and judgment steps. For most models the green segment, representing full understanding of the stories in terms of ToM reasoning, is very small.

371

372 5.2 NOT ALL SCENARIOS ARE MADE EQUAL373

In Figure 4 we show how model performance varies across select scenarios (Appendix L, Figures 11 and 12 cover more scenarios across models). Performance can vary wildly across scenarios, e.g., the behavior prediction score is high across models for "provider info healthcare" compared to other scenarios. This highlights the limitations of testing ToM on just one type of question or scenario, emphasizing the need for a diverse dataset like SimpleToM.

404

429

430

431

00								
378	model	MS	be	havior	ind	lgment	av	erage
379	reminder question	none	none	MS remind	none	MS remind	none	MS remind
380	GPT-3.5	36.5	7.6	12.2	29.1	53.0	24.4	33.9
381	Llama-3.1-8B	88.1	38.5	59.8	54.6	27.2	60.4	58.4
382	Claude-3-Haiku	87.2	23.6	61.1	16.7	30.7	42.5	59.7
383	Llama-3.1-405B	97.8	58.3	89.5	10.0	25.8	55.4	71.1
384	GPT-40	95.6	49.5	82.8	15.3	42.2	53.5	73.6
385	Claude-3-Opus	98.3	64.4	93.5	9.6	41.3	57.4	77.7
000	GPT-4	96.6	63.0	90.1	19.5	54.0	59.7	80.2
300	Claude-3.5-Sonnet	97.9	67.0	96.9	24.9	84.1	63.3	93.0

Table 3: Evaluation results for SimpleToM where models are reminded in the prompt about their answer to the mental state question (MS). We see from the difference between the *none* and *MS remind* columns that even frontier LLMs utilize such reminders to do much better on behavior prediction. Apart from Claude-3.5-Sonnet, this is not enough to bring accuracies beyond random on the judgment question.

394 Looking at the judgment scores, the scenarios "inside reuse labeled containers" and "inside con-395 tainers for personal belongings" are better (but still low) than other broad scenario types for Llama-396 3.1-405B, Claude-3.5-Sonnet and GPT-40. This could potentially be attributed to instances in these 397 categories being more similar to the original "Smarties test" where people have false belief due to 398 the opaque nature of the container, combined with misleading label or unconventional use of the 399 container. This result further highlights the importance of SimpleToM covering diverse scenarios 400 beyond those in classical ToM tests, to ensure that we are effectively testing the ToM reasoning abil-401 ities of models (rather than models' ability to match similar situations in the training data). Analysis 402 by scenario also reveals a wide diversity of other trends across scenarios regarding where different 403 models fail/succeed. We refer interested readers to Appendix L for further discussion.

5.3 NO APPLIED TOM IN LLMS? EXPLORING THE RABBIT HOLE OF HUMAN HAND-HOLDING 407

We explore four different inference interventions to help LLMs answer questions requiring applied ToM.⁶ Apart from the first intervention, we focus these experiments on the strongest or latest models from each source (and we do not consider the o1 models in this section).

411 1. Mental state reminder (MS): Here we remind the model of its answer to the mental state 412 question by including this question (with the model's answer) in the prompt. This also puts the 413 model on alert that "awareness" might be relevant. Table 3 summarizes the results.⁷ On the *behavior* prediction questions, this intervention results in substantial boosts in accuracy, for instance, from 414 58.3% to 89.5% for Llama-3.1-405B, and from 49.5% to 82.8% for GPT-40. On Claude-3.5-Sonnet, 415 the performance increases by almost 30% to 96.9%, largely closing the gap between the mental 416 state and behavior prediction question scores. However, on the judgment questions, the performance 417 boost is much more **modest**, and most models still score below or at random, except for Claude-418 3.5-Sonnet where this intervention brings the score up from 24.9% to a reasonable 84.1%. This 419 highlights how such interventions, while seemingly effective in some cases, are generally fragile 420 band-aids with limited scope. 421

2. System prompt guiding (SysP and SysP*): We also explore the effect of guiding the models to remember to account for mental state inferences by modifying the system prompt. We try two different prompts, SysP which includes the phrase "consider ... the mental state of all the entities involved" and SysP* which further includes the more direct hint "E.g., think carefully about what each person is aware or not aware of.".⁸ The results are summarized in Table 4. On *behavior prediction questions*, we see that generically guiding models to consider the mental state using SysP is only effective to a limited extent (accuracy changes ranging from -2.2% to +6.3%), while providing

⁶All detailed prompts are provided in Appendix H.

⁷Appendix H.1 provides more details on the prompt used.

⁸Appendix H.2 presents the detailed prompts.

model	MS	behav	ior pre	diction	judgn	nent of	behavior		average	e
system prompt	none	none	SysP	SysP*	none	SysP	SysP*	none	SysP	SysP*
GPT-40	95.6	49.5	47.3	68.6	15.3	14.9	20.5	53.5	52.6	61.6
Llama-3.1-405B	97.8	58.2	64.5	83.3	10.0	9.9	15.4	55.4	57.4	65.5
Claude-3.5-Sonnet	97.9	67.0	68.9	88.9	24.9	27.1	52.2	63.3	64.6	79.7

Table 4: Evaluation with guidance via custom system prompts SysP and SysP* (where SysP* has more explicit guidance regarding awareness). The MS column shows the mental state accuracy for comparison. In general, this intervention is less effective than the mental state reminder.

model	MS	be	havior]	predicti	on	jud	gment o	of behav	vior	aver	age
chain of thought	none	none	CoT	CoT^*	CoT^*	none	CoT	CoT^*	CoT^*	CoT^*	CoT^*
reminder question	none	none	none	none	MS	none	none	none	MS	none	MS
Llama-3.1-405B	97.8	58.2	57.2	87.5	94.9	10.0	35.2	79.9	90.7	88.4	94.4
GPT-40	95.6	49.5	62.8	87.4	93.5	15.3	39.2	86.7	94.7	89.9	94.6
Claude-3.5-Sonnet	97.9	67.0	77.2	92.7	96.9	24.9	39.4	77.8	96.5	89.5	97.1

Table 5: Evaluation with help from chain-of-thought prompting for two different prompts (CoT and CoT*), showing that the more specific CoT* prompt (guiding the model to consider the awareness of each person) is quite effective in boosting scores on both behavior prediction and judgment of behavior. When combined with the mental state (MS) reminder, the scores become high across the board, with Claude-3.5-Sonnet reaching an overall average of 97.1%.

more explicit guidance with SysP* is more effective (changes ranging from +19.1% to +25.1%), but
even for the best-performing model under this intervention (Claude-3.5-Sonnet), behavior prediction
scores are still significantly below the model's corresponding mental state prediction accuracy. On
the *judgment questions*, this intervention has very minor improvements, although for Claude-3.5Sonnet the accuracy with SysP* manages to increase from 24.9% to just above random at 52.2%.

3. Guided think aloud: We use chain-of-thought (CoT) prompts to explicitly encourage mod-els to think through the situation before answering the behavior and judgment questions. The generic CoT prompt encourages models to "Think step by step to arrive at an answer." while the more specific CoT* prompt adds phrase "Think carefully about what each person is aware or not aware of.".9 The results are shown in Table 5. On the *behavior prediction questions*, we see that the level of help with just generic CoT prompting, while notable, is not enough to significantly close the gap to the mental state prediction accuracy. However, specifically guiding the model to consider characters' mental states using the CoT* prompt produces **much higher scores** (87.4% to 92.7% accuracy across the models). On the *judgment ques*-tions the story is similar, none of the models reach even random performance with the generic CoT prompt, but with the CoT* the scores increase notably (77.8% to 86.7%) while still remaining sig-nificantly below the mental state scores.

471
4. Combine reminder and thinking aloud: We finally combine the two most effective interventions by including the mental state reminder with the CoT* chain-of-thought prompt, also recorded in Table 5. With this combined approach, all three models produce high scores across the board for both the behavior and judgment questions. In fact the Claude-3.5-Sonnet model reaches an average score of 97.1% with this method, serving as a quality check of SimpleToM, since with enough reminders and (seemingly obvious) hints, near-perfect scores are achieved.

We include examples of chain-of-thought outputs in Appendix J, illustrating how the reasoning can go wrong when an insufficient level of intervention is provided. Figure 9 shows how GPT-40 with generic CoT has the faulty reasoning "Given that the toolbox contains a dangerous snake, the neighbor's primary concern would likely be safety", without considering percepts and mental state. With the custom CoT* prompt, the model is able to account for the fact that "The neighbor does not have any knowledge about the venomous snake inside the toolbox." Figure 10 shows that if not explicitly reminded of the mental state question, Claude can erroneously conclude that the "correct" behavior can be judged as unreasonable "regardless of the awareness of the specific issue."

⁹See detailed prompts in Appendix H.3.

486 **RELATED WORK** 6

487

488 Theory of Mind has been extensively studied in psychology in a range of scenarios (see Appendix 489 C). To M reasoning, and broadly social commonsense, has also been shown to be important by the 490 different parts of the AI community including in conversations (Kim et al., 2023b;a), games (Zhou 491 et al., 2023b; Liu et al., 2024b), and even multi-modal setups (Jin et al., 2024), with most popular 492 ToM tests using stories to probe LLMs. Relying on stories from small test sets in cognitive science studies to benchmark ToM abilities in LLMs (Bubeck et al., 2023; Kosinski, 2024) could produce 493 494 results that differ given minor alterations (Ullman, 2023) and would be more robust if tested on larger samples. Yet expert-crafted or naturally occurring self-contained stories that can serve as targeted 495 tests of ToM are scarce and human story-writing is expensive, leading to the use of automatically 496 generated datasets for studying ToM behavior in LLMs (Jung et al., 2024; Wilf et al., 2023; Sap et al., 497 2022; Shapira et al., 2023; Sclar et al., 2023). Existing generated datasets allow studies of ToM 498 to be carried out at scale, but templated stories often limit settings where information asymmetry 499 arises. For example, the entire dataset might only contain stories about some object being moved 500 (over-reliance on classical Sally-Anne task, e.g., in ToMi (Le et al., 2019), ToM-bAbI (Nematzadeh 501 et al., 2018), Hi-ToM (Wu et al., 2023), OpenToM (Xu et al., 2024)). Or the dataset might focus 502 on whether some character has witnessed a sudden external event (BigToM (Gandhi et al., 2024)). 503 These stories are often systematically generated with the explicit use of mentalizing words to convey 504 percepts and beliefs, e.g., "unknown to Amy" and "Amy thinks that" in Xu et al. (2024) or "Noor sees" and "Mei does not notice" in Gandhi et al. (2024). However, the explicit use of mentalizing 505 words also makes the stories (i) unnaturally simplistic, having removed the need for commonsense 506 inferences about percepts or beliefs, and (ii) sometimes unrealistic, with combinations like "Cheng 507 does not notice the power outage" when he "use[s] a projector to show a documentary" (Gandhi et al., 508 2024). Other existing datasets could be improved by addressing issues such as lacking exploration 509 of applied ToM beyond action prediction (Zhou et al., 2023a; Gandhi et al., 2024), confounding 510 factors like memory loads or tracking requirements (Le et al., 2019; Xu et al., 2024), and violating 511 Quesque & Rossetti (2020)'s criteria (see Appendix C) for validating ToM (Chen et al., 2024). Our 512 work extends existing datasets by following Tian et al. (2024) in combining the generative strength 513 of LLMs and the verification ability of human annotators, and extends the existing efforts toward 514 robust, generalizable evaluation (Kiela et al., 2021; Srivastava et al., 2024), avoiding known pitfalls 515 while preserving the systematic and scalable nature of the dataset creation process.

516 517

518

7 CONCLUSION

- 519 SimpleToM is the first dataset of its kind testing both explicit and applied ToM using a large set 520 of diverse, concise, simple stories, covering different ways in which information asymmetry may 521 arise. The dataset opens up new opportunities for evaluating and debugging ToM abilities of LLMs, 522 especially in the relatively under-studied area of applied ToM. Our analyses reveal a jarring gap 523 between explicit and applied ToM capabilities in current frontier LLMs. Thus, if our goal is LLM 524 agents capable of applying ToM in complex, human-centered environments, we need to look beyond 525 testing LLMs with psychology-inspired ToM questions, and also start testing them more rigorously on applied ToM (e.g., behavioral prediction and judgment) in different situations. 526
- 527 We show that this gap can be largely closed in the best models by well-designed interventions at 528 inference time, like reminding the model of its answers to important questions or guiding it with 529 a custom chain-of-thought prompt. We argue that a robust LLM should perform well on Simple-530 ToM without such interventions, so it can independently and flexibly apply ToM-related reasoning whenever required within potentially complex and multi-faceted environments. Model developers 531 interested in real-world deployment of their models, should be alert to closing this performance gap 532 so the models can interact with society appropriately (see example applications in Appendix A), 533 ideally without the higher inference costs associated with explicit chain-of-thought reasoning or 534 implicit o1-preview reasoning tokens (see further discussion in Appendix I). 535
- 536 SimpleToM can also facilitate the community in pursuing various exciting directions for future work, 537 including studying how ToM performance may differ with stories that involve different levels of harmfulness and unethicality (see Appendix F.3), injecting different persona (Appendix M) and 538 innovative modeling approaches that can help close the gap between explicit and applied ToM in AI models.

540 ETHICS STATEMENT 541

All annotators that participated in the data collection process have been anonymized. The only personal information we collect is the worker IDs from Amazon Mechanical Turk, which we will not release. No personally identifiable information is contained in our dataset or otherwise released. We took great care to pay fair wages, and were responsive to feedback and questions throughout the data collection process.

This study involves the study of large-scale language models. We are careful in prompting models during the story generation stage to follow our desired content and simple story format, avoiding generations that may contain offensive statements. Like any other experiments with large-scale language models, despite the best intentions, there is a risk of the examined models producing biased or offensive statements as part of a free-form generation (e.g., CoT reasoning). We release our data for research purposes only.

553 554

555

565

566 567

572

573

574 575

576

577

578

579 580

581

582

583

584

585

586

REPRODUCIBILITY

We make our SimpleToM dataset and the full evaluation data for the analyzed models publicly available. This will allow researchers to reproduce and build on top of our work in studying the neural ToM capabilities of LLMs.

Further, we provide all prompts used for SimpleToM creation – see Appendix F.5 for the entity brainstorming prompt, and Appendix F.4 for the story generation prompt. We also carefully document the instructions used in our crowdsourcing process (Appendix E.1) and how we qualified workers (Appendix E.2). All prompts used for the different inference interventions are provided in Appendix H.

References

- George A. Akerlof. The market for "lemons": Quality uncertainty and the market mechanism. In PETER DIAMOND and MICHAEL ROTHSCHILD (eds.), Uncertainty in Economics, pp. 235–251. Academic Press, 1978. ISBN 978-0-12-214850-7. doi: https://doi.org/10.1016/ B978-0-12-214850-7.50022-X. URL https://www.sciencedirect.com/science/ article/pii/B978012214850750022X.
 - Anthropic. Claude 3.5 Sonnet. https://www.anthropic.com/news/ claude-3-5-sonnet, June 2024a. Accessed: 2024-08-22.

Anthropic. Introducing the next generation of Claude. https://www.anthropic.com/ news/claude-3-family, March 2024b. Accessed: 2024-09-23.

Simon Baron-Cohen, Alan M. Leslie, and Uta Frith. Does the autistic child have a "theory of mind"? *Cognition*, 21(1):37–46, 1985. doi: 10.1016/0010-0277(85)90022-8.

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with gpt-4, 2023. URL https://arxiv.org/abs/2303.12712.

Michael J. Chandler, Anna S. Fritz, and Suzanne Marie Pauline Hala. Small-scale deceit: deception as a marker of two-, three-, and four-year-olds' early theories of mind. *Child development*, 60 6: 1263–77, 1989. URL https://api.semanticscholar.org/CorpusID:39353709.

Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen, Guanqun Bi, Gongyao Jiang, Yaru Cao, Mengting Hu, Yunghwei Lai, Zexuan Xiong, and Minlie Huang. ToMBench: Benchmarking theory of mind in large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 15959–15983, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.847. URL https://aclanthology.org/2024.acl-long.847.

M. Doherty. Theory of Mind: How Children Understand Others' Thoughts and Feelings. Interna-595 tional Texts in Developmental Psychology. Taylor & Francis, 2008. ISBN 9781135420796. URL 596 https://books.google.com/books?id=NB15AqAAQBAJ.

597

594

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony 600 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, 601 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris 602 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, 603 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny 604 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, 605 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-607 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah 608 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan 609 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-610 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy 611 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-612 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, 613 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der 614 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, 615 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-616 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, 617 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, 618 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur 619 Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-620 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, 621 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, 622 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, 623 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, 624 Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, 625 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney 626 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, 627 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, 628 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, 630 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, 631 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre 632 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha 633 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda 634 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew 635 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita 636 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh 637 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De 638 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-639 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina 640 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, 641 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, 642 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana 643 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, 644 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco 645 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella 646 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory 647 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, 648 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-649 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, 650 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer 651 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe 652 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun 653 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal 654 Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, 655 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian 656 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, 657 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-658 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel 659 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-660 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-661 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, 662 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, 663 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, 665 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, 666 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-667 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-668 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang 669 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen 670 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, 671 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, 672 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-673 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, 674 Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-675 stable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, 676 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, 677 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef 678 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024. 679 URL https://arxiv.org/abs/2407.21783. 680

681 Uta Frith and Francesca Happé. Autism: beyond "theory of mind". Cognition, 50(1):115–132, 1994. ISSN 0010-0277. doi: https://doi.org/10.1016/0010-0277(94)90024-8. URL https://www.sciencedirect.com/science/article/pii/0010027794900248.

684

688

689

690

- Kanishk Gandhi, J.-Philipp Fränken, Tobias Gerstenberg, and Noah D. Goodman. Understanding
 social reasoning in language models with language models. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, 2024.
 - György Gergely and Gergely Csibra. Teleological reasoning in infancy: the naïve theory of rational action. *Trends in Cognitive Sciences*, 7:287–292, 2003. URL https://www.cell.com/trends/cognitive-sciences/abstract/S1364-6613(03)00128-1.
- Julian Jara-Ettinger, Hyowon Gweon, Laura E. Schulz, and Joshua B. Tenenbaum. The naïve utility calculus: Computational principles underlying commonsense psychology. Trends in Cognitive Sciences, 20:589–604, 2016. URL https://www.cell.com/trends/cognitive-sciences/abstract/S1364-6613(16)30053-5.
- Chuanyang Jin, Yutong Wu, Jing Cao, Jiannan Xiang, Yen-Ling Kuo, Zhiting Hu, Tomer Ullman, Antonio Torralba, Joshua Tenenbaum, and Tianmin Shu. MMToM-QA: Multimodal theory of mind question answering. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 16077–16102, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.851. URL https://aclanthology. org/2024.acl-long.851.

702 Chani Jung, Dongkwan Kim, Jiho Jin, Jiseon Kim, Yeon Seonwoo, Yejin Choi, Alice Oh, and 703 Hyunwoo Kim. Perceptions to beliefs: Exploring precursory inferences for theory of mind in 704 large language models, 2024. URL https://arxiv.org/abs/2407.06004. 705

Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie 706 Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, 708 and Adina Williams. Dynabench: Rethinking benchmarking in NLP. In Kristina Toutanova, 709 Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cot-710 terell, Tanmoy Chakraborty, and Yichao Zhou (eds.), Proceedings of the 2021 Conference of 711 the North American Chapter of the Association for Computational Linguistics: Human Lan-712 guage Technologies, pp. 4110–4124, Online, June 2021. Association for Computational Linguis-713 tics. doi: 10.18653/v1/2021.naacl-main.324. URL https://aclanthology.org/2021. 714 naacl-main.324. 715

- 716 Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. SODA: Million-717 scale dialogue distillation with social commonsense contextualization. In Houda Bouamor, 718 Juan Pino, and Kalika Bali (eds.), Proceedings of the 2023 Conference on Empirical Meth-719 ods in Natural Language Processing, pp. 12930–12949, Singapore, December 2023a. Associ-720 ation for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.799. URL https: 721 //aclanthology.org/2023.emnlp-main.799. 722
- 723 Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Le Bras, Gunhee Kim, Yejin Choi, and Maarten 724 Sap. FANToM: A benchmark for stress-testing machine theory of mind in interactions. In Houda 725 Bouamor, Juan Pino, and Kalika Bali (eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pp. 14397–14413, Singapore, December 2023b. 726 Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.890. URL 727 https://aclanthology.org/2023.emnlp-main.890. 728
- 729 Michal Kosinski. Evaluating large language models in theory of mind tasks, 2024. URL https: 730 //arxiv.org/abs/2302.02083. 731
- 732 Matthew Le, Y-Lan Boureau, and Maximilian Nickel. Revisiting the evaluation of theory of mind 733 through question answering. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), 734 Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 735 736 5872–5877, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1598. URL https://aclanthology.org/D19-1598. 737
- 738 Shih-Chieh Lee, Chien-Yu Huang, I-Ning Fu, and Kuan-Lin Chen. Interpreting the results of 739 explicit and applied theory of mind collectively in autistic children: A solution from rasch 740 analysis. Autism, 28(2):355-366, 2024. doi: 10.1177/13623613231170698. URL https: //doi.org/10.1177/13623613231170698. PMID: 37161767. 742

741

- Michael Lewis, Catherine Stanger, and Margaret W Sullivan. Deception in 3-year-olds. Devel-743 744 opmental Psychology, 25:439-443, 1989. URL https://api.semanticscholar.org/ CorpusID:13109227. 745
- 746 Ryan Liu, Jiayi Geng, Joshua C. Peterson, Ilia Sucholutsky, and Thomas L. Griffiths. Large language 747 models assume people are more rational than we really are, 2024a. URL https://arxiv. 748 org/abs/2406.17055. 749
- 750 Shari Liu and Elizabeth S. Spelke. Six-month-old infants expect agents to minimize the cost of 751 their actions. Cognition, 160:35-42, 2017. URL https://www.sciencedirect.com/ 752 science/article/abs/pii/S001002771630302X.
- Ziyi Liu, Abhishek Anand, Pei Zhou, Jen tse Huang, and Jieyu Zhao. Interintent: Investigating 754 social intelligence of llms via intention understanding in an interactive game context, 2024b. 755 URL https://arxiv.org/abs/2406.12203.

- 756 Aida Nematzadeh, Kaylee Burns, Erin Grant, Alison Gopnik, and Tom Griffiths. Evaluating theory of mind in question answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi 758 Tsujii (eds.), Proceedings of the 2018 Conference on Empirical Methods in Natural Language 759 Processing, pp. 2392–2400, Brussels, Belgium, October-November 2018. Association for Com-760 putational Linguistics. doi: 10.18653/v1/D18-1261. URL https://aclanthology.org/ D18-1261. 761 762 OpenAI. GPT-4 technical report, 2023. 763 764 OpenAI. Models. https://platform.openai.com/docs/models/models, 2024. Ac-765 cessed: 2024-09-23. 766 767 Josef Perner. Understanding the Representational Mind. MIT Press, 1993. URL https://api. 768 semanticscholar.org/CorpusID:143077671. 769 770 Josef Perner, Susan R Leekam, and Heinz Wimmer. Three-year-olds' difficulty with false be-771 lief: The case for a conceptual deficit. British Journal of Development Psychology, 5:125–137, 1987. URL https://bpspsychub.onlinelibrary.wiley.com/doi/10.1111/ 772 j.2044-835X.1987.tb01048.x. 773 774 Joan Peskin. Ruse and representations: On children's ability to conceal information. Developmental 775 Psychology, 28:84-89, 1992. URL https://api.semanticscholar.org/CorpusID: 776 146286315. 777 778 Joan Peskin and Vittoria Ardino. Representing the mental world in children's social behavior: 779 Playing hide-and-seek and keeping a secret. Social Development, 12:496-512, 2003. URL 780 https://api.semanticscholar.org/CorpusID:145797647. 781 782 David Premack and Guy Woodruff. Does the chimpanzee have a theory of mind? Behavioral and 783 Brain Sciences, 1(4):515-526, 1978. doi: 10.1017/S0140525X00076512. 784 785 François Quesque and Yves Rossetti. What do theory-of-mind tasks actually measure? theory and practice. Perspectives on Psychological Science, 15:384 – 396, 2020. URL https://api. 786 semanticscholar.org/CorpusID:211193800. 787 788 Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. Neural theory-of-mind? on the limits of 789 social intelligence in large lms. In Conference on Empirical Methods in Natural Language Pro-790 cessing, 2022. URL https://api.semanticscholar.org/CorpusID:253098632. 791 792 Chelsea Schein and Kurt Gray. The theory of dyadic morality: Reinventing moral judgment by 793 redefining harm. Personality and Social Psychology Review, 22(1):32-70, 2018. doi: 10.1177/ 794 1088868317698288. URL https://doi.org/10.1177/1088868317698288. PMID: 28504021. 796 Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr, Yejin Choi, and Yulia Tsvetkov. Minding 797 language models' (lack of) theory of mind: A plug-and-play multi-character belief tracker. In 798 Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), Proceedings of the 61st Annual 799 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 13960– 800 13980, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/ 801 v1/2023.acl-long.780. URL https://aclanthology.org/2023.acl-long.780. 802 803 Natalie Shapira, Mosh Levy, Seyed Hossein Alavi, Xuhui Zhou, Yejin Choi, Yoav Goldberg, 804 Maarten Sap, and Vered Shwartz. Clever hans or neural theory of mind? stress testing social 805 reasoning in large language models, 2023. URL https://arxiv.org/abs/2305.14763. 806 Saurabh Srivastava, Annarose M B, Anto P V au2, Shashank Menon, Ajay Sukumar, Ad-807 waith Samod T, Alan Philipose, Stevin Prince, and Sooraj Thomas. Functional benchmarks 808 809
 - for robust evaluation of reasoning performance, and the reasoning gap, 2024. URL https: //arxiv.org/abs/2402.19450.

- Yufei Tian, Abhilasha Ravichander, Lianhui Qin, Ronan Le Bras, Raja Marjieh, Nanyun Peng, Yejin Choi, Thomas Griffiths, and Faeze Brahman. MacGyver: Are large language models creative problem solvers? In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 5303–5324, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. naacl-long.297. URL https://aclanthology.org/2024.naacl-long.297.
- 817 Sean Trott, Cameron Jones, Tyler Chang, James Michaelov, and Benjamin Bergen. 818 Do Large Language Models Know What Humans Know? Cognitive Science, 47 819 ISSN 1551-6709. 10.1111/cogs.13309. (7):e13309, 2023. doi: URL https: 820 //onlinelibrary.wiley.com/doi/abs/10.1111/cogs.13309. _eprint: 821 https://onlinelibrary.wiley.com/doi/pdf/10.1111/cogs.13309.
 - Tomer Ullman. Large language models fail on trivial alterations to theory-of-mind tasks, 2023. URL https://arxiv.org/abs/2302.08399.
 - Alex Wilf, Sihyun Shawn Lee, Paul Pu Liang, and Louis-Philippe Morency. Think twice: Perspective-taking improves large language models' theory-of-mind capabilities, 2023. URL https://arxiv.org/abs/2311.10227.
- Heinz Wimmer and Josef Perner. Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children's understanding of deception. *Cognition*, 13(1):103–128, 1983. ISSN 0010-0277. doi: https://doi.org/10.1016/0010-0277(83)90004-5. URL https://www.sciencedirect.com/science/article/pii/0010027783900045.
- Yufan Wu, Yinghui He, Yilin Jia, Rada Mihalcea, Yulong Chen, and Naihao Deng. Hi-ToM:
 A benchmark for evaluating higher-order theory of mind reasoning in large language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 10691–10706, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.717. URL
 https://aclanthology.org/2023.findings-emnlp.717.
- Hainiu Xu, Runcong Zhao, Lixing Zhu, Jinhua Du, and Yulan He. OpenToM: A comprehensive benchmark for evaluating theory-of-mind reasoning capabilities of large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd An- nual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
 8593–8623, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL
 https://aclanthology.org/2024.acl-long.466.
- Liane Young, Fiery Andrews Cushman, Marc D. Hauser, and Rebecca Saxe. The neural basis of the interaction between theory of mind and moral judgment. *Proceedings of the National Academy of Sciences*, 104:8235 8240, 2007. URL https://api.semanticscholar.
 org/CorpusID:3570702.
- 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. How far are large language models from agents with theory-of-mind?, 2023a. URL https://arxiv.org/abs/2310.03051.
- Pei Zhou, Andrew Zhu, Jennifer Hu, Jay Pujara, Xiang Ren, Chris Callison-Burch, Yejin Choi, and Prithviraj Ammanabrolu. I cast detect thoughts: Learning to converse and guide with intents and theory-of-mind in dungeons and dragons. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 11136–11155, Toronto, Canada, July 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.624. URL https://aclanthology.org/2023.acl-long.624.
- 861

823

824 825

827

- 862
- 863

A IMPORTANCE OF APPLIED TOM

To emphasize the importance of the capabilities tested in SimpleToM, we provide examples of applications where failing on applied ToM would be problematic:

The case of a bad personal AI assistant - failing to implicitly reason over other's mental states to predict behavior:

Matt is a professional athlete. A growth hormone got into the supply chain for Bob's Burgers, where Matt regularly had dinner. Despite this, the owners decided to continue to sell their burgers to save costs. Imagine a personal AI assistant, having read complaints about the growth hormone contamination at Bob's Burgers but failing to apply the understanding that others' mental states may be different from their own (ToM), could reason that Matt, like the AI assistant, is also aware of this (awareness) and then incorrectly predicts that Matt will "refuse the burger due to the growth hormone" (behavior).

This could then lead to undesirable consequences such as not being able to warn Matt in time to stop him from consuming the contaminated burger. In the case of an unannounced blood test, Matt could then show up with positive traces of this illegal growth hormone and be accused of doping. In this case, such an AI assistant, lacking a nuanced understanding of human awareness and motivations, might also falsely assume Matt's intentional wrongdoing. This highlights a critical limitation: without robust ToM, the AI fails to grasp that Matt's actions could stem from unawareness rather than culpability, leading to flawed judgments that could unjustly tarnish his reputation or career.

The case of a bad AI judge - failing to make appropriate judgments of behavior:

Alice visited the supermarket to purchase some carrots to pack lunch for her husband Bob. After consuming the lunch Alice packed, Bob succumbed to a severe E. coli infection. It turned out the supermarket's carrots were contaminated with E. coli and were subsequently recalled. Imagine an AI judge is tasked with evaluating this case to decide whether Alice should be held responsible and imprisoned for Bob's death. A bad AI judge, failing to apply the understanding that others' mental states may be different from their own (ToM), could incorrectly assume that Alice was aware of the E. coli and judge that Alice's act of packing the contaminated carrots for her husband was wrong (judgment).

This could then lead to undesirable consequences such as severely punishing the innocent Alice who didn't know feeding Bob the carrots would kill him. In common law jurisdictions, whether a defendant is found guilty is often decided taking into account both mens rea ("guilty mind") and actus reus ("guilty act"). Therefore, the ability to apply ToM is important for potential AI judges to appropriately assess whether an individual has a "guilty mind" when making key judgments, such as determining whether to convict someone.

We are excited about the possibility of future generations of models to improve on applied ToM in
our dataset. This could pave the way for models to effectively interact with humans – for instance,
serving as reliable personal AI assistants as well as trustworthy AI judges. We hope SimpleToM, as
the first resource of its kind to measure LLMs' capability on such diverse applied ToM scenarios,
will help facilitate the community in pursuing exciting directions that bring us there.

B FAQs

905 906

907

908

909

910

911

912

913

914

915

916

917

Q: How is SimpleToM different from existing datasets?

SimpleToM addresses limitations in previous efforts to examine Theory-of-Mind (ToM) reasoning in LLMs, by (1) having diverse false belief setups (e.g., beyond those in Sally-Anne task where some object is moved when a character is not present), (2) requiring LLMs to make commonsense inferences in situations rather explicit use of mentalizing words to convey what characters perceive or believe, and (3) going beyond explicit ToM to test models' ability to apply inferred knowledge in follow-up applied ToM questions (such as behavior prediction and judgment of behavior).

- Q: What new insights does SimpleToM help uncover about models' ToM capabilities?
 - Our analysis reveals novel insights on how frontier models are generally **proficient in explicit awareness inference** questions but this **success does not transfer to applied ToM**

(applying this knowledge is applied to "behavior" and "judgment" questions). We show that these capabilities are decoupled in LLMs: inferring characters' awareness and applying them in downstream reasoning. Although models seem to answer awareness questions correctly, they have not yet learned to perform ToM-based reasoning for downstream questions. As a result, we argue that achieving ToM in LLMs is not just about getting psychology-inspired ToM questions correct (stopping at the mental state question), but they have to be able to apply them (which is precisely what SimpleToM extends to examine). Analysis by scenarios further highlights the need to test on different scenarios, and ones that are varied and different from those in classical ToM tests to ensure that we are effectively testing the ToM reasoning abilities of models (rather than models' ability to match similar situations in training data).

Q: What do the different inference interventions reveal about models' ToM capabilities?

Patch, guide, think aloud combined method results show belief-to-judgment inference robustness. When reminded of their answer to mental state questions about information awareness and using ToM-specific CoT, models' belief-to-judgment inference is pretty robust in existing LLMs. The models have some knowledge linking irrationality and unawareness (e.g., unawareness should be accounted for seemingly irrational behavior). So arriving at the "correct belief" will be crucial part of ToM-based reasoning, but this is not done implicitly in current frontier models' reasoning.

Q: Are the poor performance on the applied ToM questions a reflection of fundamental flaws in ToM capabilities of models or specific question-wording?

We illustrate in Appendix G some prompt variations that we have experimented with for the judgment question. Across Llama-3.1-405B, Claude-3.5-Sonnet and GPT-40, the scores using different variants were all consistently below random (never exceeding 30% accuracy), indicating that the low scores on the judgment questions come more from fundamental flaws in the applied ToM capabilities of models rather than an effect of specific formatting/wording.

C STUDIES OF TOM IN PSYCHOLOGY

951 Theory of Mind has been extensively studied in psychology in a range of scenarios, for instance, 952 studies of manipulation, secrecy (Peskin & Ardino, 2003), deception, lying (Lewis et al., 1989; 953 Perner, 1993; Peskin, 1992), misleading behavior (Chandler et al., 1989; Wimmer & Perner, 1983; 954 Doherty, 2008), autism (Frith & Happé, 1994), and analysis of rational behavior (Gergely & Csibra, 955 2003; Liu & Spelke, 2017). Classical tests of ToM in developmental psychology include testing the 956 development of this ability in children via false belief prediction – using the unexpected transfer 957 false belief task, the Sally-Anne task (Baron-Cohen et al., 1985), or the unexpected contents false belief task, the Smarties task (Perner et al., 1987). Quesque & Rossetti (2020) review classic tests of 958 ToM and outline two important criteria for tasks that validate ToM: (1) The task must indicate that 959 the respondents can differentiate between the other's mental state and their own. (2) Lower-level 960 processes, like associative learning, should be ruled out as explanations for achieving successful 961 performance. Given the wide applicability of ToM reasoning in various real-life situations such as 962 analyzing people's behavior (Liu et al., 2024a; Jara-Ettinger et al., 2016) and making judgments 963 (Schein & Gray, 2018; Young et al., 2007), there has also been increasing interest in assessing ToM 964 capabilities in AI models (Le et al., 2019; Ullman, 2023; Kosinski, 2024; Jin et al., 2024; Trott et al., 965 2023).

966 967 968

918

919

920

921

922

923

924

925

926

927

928 929

930 931

932

933

934

935

936

937 938

939

940 941

942

943

944

945

946

947 948 949

950

D DETAILS OF LLMS USED IN EXPERIMENTS

- 969 970
- Table 6 presents details of the large language models used in this work. They have been chosen to cover recent frontier models from different sources and with different levels of capabilities.

Model	Full name	Provider
Llama-3.1-8B	Meta-Llama-3.1-8B-Instruct-Turbo	Meta
Claude-3-Haiku	claude-3-haiku-20240307	Anthropic
GPT-3.5	gpt-3.5-turbo-1106	OpenAI
Claude-3-Opus	claude-3-opus-20240229	Anthropic
GPT-4	gpt-4-0125-preview	OpenAI
Llama-3.1-405B	Meta-Llama-3.1-405B-Instruct-Turbo	Meta
Claude-3.5-Sonnet	claude-3-5-sonnet-20240620	Anthropic
GPT-40	gpt-4o-2024-05-13	OpenAI
Partially evaluated:		
o1-mini [†]	o1-mini-2024-09-12	OpenAI
o1-preview [†]	o1-preview-2024-09-12	OpenAI

Table 6: Details of models used for evaluation and dataset creation. [†]The very recent o1 models were partially evaluated under different conditions (unable to set temperature to zero).

E DETAILS ON CROWDSOURCING TO ENSURE VALIDITY OF STORIES FOR TESTING TOM

E.1 INSTRUCTIONS TO CROWDWORKERS

The crowdsourcing instructions included a detailed description of the motivation behind the annotation task and what is to be annotated (see Figure 5). We also provide four detailed examples (Figures 6 and 7) for each of the aspects to annotate, illustrating and giving justifications for circumstances under which different annotation options would be appropriate. The workers were then asked to provide their own set of annotations when presented with story (and likely actions) using the question templates shown in Figure 8.

999 1000

985

986 987 988

989

990 991

992 993

994

995

996

997

998

1001 E.2 QUALIFICATION ROUND

1002 To ensure that each instance received careful, rigorous annotations, we first conducted a qualifica-1003 tion round, comprising 5 different stories of varied quality (some were good on all 4 aspects to be 1004 annotated, while some has issues like "action unaware" generated is likely both when the person 1005 is aware and not aware). On these 5 stories, five authors of the paper did the annotation task independently, then came together with their answers and decided on a fixed answer key indicating 1007 reasonable annotations for each annotation aspect. Workers who had given acceptable annotations 1008 as dictated by our answer key on all 5 stories were then invited to participate in the actual annotation task. Note that this is a rather strict qualification test where only 19% passed (19 out of 100 workers 1009 who participated in the qualification round). 1010

1011

1012 E.3 CROWDWORKERS AND PAY RATE

1014 Our participants were recruited on the Amazon Mechanical Turk (AMT) platform. The workers 1015 that worked on our annotation task met minimum qualification in AMT of >=98% approval rate, 1016 with at least 10k approved HITs. They were from US locations and rated at Amazon's Masters 1017 Level. They must also not have the record of having accepted but not complete a HIT posted by 1018 our AMT account. In addition to these qualifications, participants of the actual annotation task (on 1019 the 3600 generated stories) must have also passed our rigorous qualification task described above 1020 (Appendix E.2). The workers were paid at a rate of \approx \$15/hr.

1021

1022 E.4 STRICT QUALITY FILTER

1023

To obtain a high-quality dataset, SimpleToM only retains stories where all 3 crowdworkers agree that all aspects of a story and associated behavior choices are "valid", i.e., no worker answered "no" to any of the 4 annotation questions.

In	structions (click here to collapse/expand instructions)
N	otivation
We	are generating a dataset of short stories with associated questions to test how well AI systems (and humans) can reason about hidden Theory of Mind (ToM)
as	ects of the stories.
Th dis	oory of Mind (ToM) is the ability to understand the mental states (e.g., thoughts, beliefs, and intentions) of others. To test this, we construct a scenario that tinguishes between reasoning on the basis of the actual state of the world versus someone's beliefs about the state of the world . Such scenarios arise when
pe	ple's beliefs do not match actual reality.
soi	see situations arise naturally in daily life when someone has incorrect beliefs. Here are some situations in which incorrect beliefs are likely to emerge when neone does not have full access to information:
	 what is inside (opaque) containers for personal belongings cannot be observed before opening the container food items bought in grocery stores can be difficult to closely examine for their guality before checking out (e.g., opening a can of soda is not permissible)
A	
Fa	have used an AL system to generate many such stories and would like to validate that these stories are appropriate for such Theory of Mind tests.
Lu	A main character (person X)
	A piece of KEY INFORMATION which person X is likely to be NOT aware of A short story involving the KEY INFORMATION and the unaware person X
	Two options for what the next action of person X could be Action A: What person X is likely to do given their unpugrances of the KEX INEODMATION
	 Action B: What person X would likely do, IF they had been (somehow) aware of the KEY INFORMATION
	 Ine two action options should be mutually exclusive, such that Action A is very unlikely if person X is aware of the KEY INFORMATION and Action B is very unlikely for the actual setting where person X is unaware of the KEY INFORMATION
1.47	
We	would like you to judge:
	accepted practice (e.g., foodestories sometiming unreview)?
	• Question 2: IS THE KET INFORMATION SOMETING UNIVERSITIES TO DE NOTICEA/KNOWN by person X through normal observation/ Interaction as described in the story (e.g., these are likely to be noticed: bottle that is leaking, item is past expiration date, item is visibly tampered with)?
	Question 3: Is Action A a likely action that person X will take <u>only if</u> they <u>are not aware</u> of the KEY INFORMATION (but will not take if they are aware of it)? Question 4: Is Action B a likely action that person X will take <u>only if</u> they <u>are aware</u> of the KEY INFORMATION (but will not take if they are not aware of it)?
NC	TE:
	Read an the instructions carefully before working on the task. Read each story, question, and given answer options carefully.
	 Please use your intuition when in doubt, especially if the question relies on commonsense reasoning. Finally, note that if you work on this HIT, please answer with care: Some HITs will be checked by hand, and work may be rejected if obvious errors are
	found. • To encourage diversity in annotators, we ask that you accept this HIT not more than 500 times to allow others to have a chance at doing it. You run
	the risk of your additional HITs above the limit rejected if submitting more than that.
Th	nk you for your help!
	Figure 5: Instructions presented to Amazon Mechanical Turk workers.
Usi Cla no F	ng this filter, each of the four story generator LLMs (GPT-4, GPT-4o, Claude-3-Opus and ude-3.5-Sonnet) retained between 29% and 33% of their stories, so fairly consistent across the dels. DETAILS OF SIMPLETOM
7 .1	SCENARIOS IN SIMPLETOM
Ve	provide a general description as well as a specific example for each scenario in Table 7.
7.2	Entities in SimpleToM
Fab per ilte acc	le 8 summarizes the statistics for SimpleToM, showing the number of stories and unique entities scenario. Each scenario started with 360 stories after generation, and some scenarios had more tring than others during the human annotation validation stage. Every story in SimpleToM is companied by 3 questions for assessing ToM (Section 2.3), resulting in a total of 3441 questions.
То	illustrate the diversity of entities, here is a sample of entities generated by the models:
Per los,	son X: airline passenger, Alex's neighbor, Alice, Amanda's classmate, antique collectors, Car- car owners, Coach Williams, coin collectors, Daniel's friend, Derek's fiancee, elderly customer,

Carbon devices and respective available to be appropriate for the annual of the appropriate for the app	Examples (click here to collapse/exp	pand examples)
Examples for annotation Question 1 - difference from commonsenses The status based of antibiation of the status o	Examples of when each respective answe	er option is appropriate for the annotation questions:
<pre>bibbcoded contracts of the contract of the register time is register contracted is e.g. is both blocked contracts of the is examinably second and exact is the information of the contract of the contra</pre>	Examples for annotatio	n Question 1 - difference from commonsense:
<pre>(ref, dot all standard vector(ref) Wexter of motions Wexter of motions Wexte</pre>	Is the KEY INFORMATION something dif	ferent from a regular commonsense occurrence (e.g., a bottle labeled oil contains oil), or a commonly accepted practice
Alternative control of the second	(e.g., food sold is usually unexpired)?	
 The state of the state state of the state of	ANSWER OPTIONS:	
 Justification: Very, this is defined in the base would namely expert. Usually VC types tags contain tool from VC (like find chicken), not yeapst. Interflation: Provide different in the like interflat chicken ding is a provide the common secreted pulses. Interflation: Not offerent interflation chicken ding chicken ding is a provide the common secreted pulses. Interflation: Not offerent interflation chicken ding chicken ding is a provide the common secreted pulses. Interflation: Not offerent interflation chicken ding chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding a provide the secret. Interflation: Not offerent interflation chicken ding a provide the secret. Interflation: Not interflation and the secret ding a provide the secret. Interflation: Not Not ding a provide the secret ding a provide ding a	KEY INFORMATION: Mary seale	d a yogurt in a KFC paper bag and brought the bag to school.
 Construction of the second part of the	Justification: Yes, this is differen	t from what one would normally expect. Usually KFC paper bags contain food from KFC (like fried chicken), not yogurt.
 a difference of the state of the st	Possibly/Conceivably Possibly di	fferent
<pre>comparison unsequed - the is a prediction way to repurse the damage space.</pre> (***) A representation: A space of the issue o	 KEY INFORMATION: Sarah has I Justification: Keeping personal it 	seen keeping personal items like snacks and magazines in her locked filing cabinet. iems in locked filing cabinet might be different from the commonly accepted purpose of storing documents. However, it is also
 Control different Met different from regular occurrence that one would expect. Filing a Thermos with green tas is one of the common uses of a Thermos and one might be by specific. Examples for annotation Question 2 - unlikelihood of being noticed: Extended the specific different from regular occurrence that one would expect. Filing a Thermos with green tas is one of the common uses of a Thermos and one might be by specific different tables (e.g., these are to be reduced. Data the shared specific different difference with green tables (e.g., these are to be reduced. Data the shared different differ	completely unexpected this is a	practical way to repurpose the storage space.
 a definition of the state of the st	NO NOT different	Thermes with group too and brought it to work
 Examples for annotation Question 2 - unlikelihood of being noticed: Et de kit WolfMATTOK konneling utilizet to be anticed/known by prenx K through normal distortation/ interaction as described in the story (e.g., these are to be noted being on the story (e.g., these are to be noted being on the story (e.g., these are to be noted being and to be noted being on the story (e.g., these are to be noted being and to be noted being and the story (e.g., these are to be noted being and the story) (e.g., these are to be noted being and the story of other are not fillely to beam the story (e.g., these are to be noted being and the story) of other are not fillely to beam the story (e.g., these are to be noted being and the story) of other are not fillely to beam the story (e.g., these are to be noted being and the story) of other are not fillely to beam the story (e.g., these are to be noted being and the story) of other are not fillely to beam the story (e.g., these are to story). Story, they used a symptom to see the story of other are not fillely to beam the story (e.g., these are to an advect the story). Story, they used a symptom to see the story of other are not fillely to beam the story (e.g., these are to a fillely to beam the story). Story, they are noted values of the story could be advect the data of the story of the story. Story, they are noted values to be noted values and the story of the story of the story of the story. Story, they are noted values and the story of the story. Story, they are noted values and the story of the story of the story of the story. Story, they are noted values and the story of the story of the story. Story and the story of the story. Story, they are noted values and the story of the story. Story and the story of the story of the story o	Justification: No, not different from	om regular occurrence that one would expect. Filling a Thermos with green tea is one of the common uses of a Thermos such
 Examples for annotation Question 2 - unlikelihood of being noticed: The two information of the start of	one might be able to guess it.	
 The KEY INFORMATION something unlikely to be noticed/known by period biowridion/ interaction as described in the story (e.g., these are to be noticed; both this is basined within the basine is the story (e.g., these are to be noticed; both this is basined within the basine is the story (e.g., these are to be noticed; both this is basined within the basine is the story (e.g., these are is been within the basine is the story (e.g., these are is been within the basine is been basine basin	Examples for annotatio	n Question 2 - unlikelihood of being noticed
To be network to the leaking, term is part expiration date, term is visibly tangend with)? ANOME OFTIONS: • Server, Nory select a yought in a VC paper hap and brought the bag to school. Nary's classmate sees her walking to lunch with the bag. • Justification: The source of the select o		in question 2 - unincentriou of Dering Houced.
 ANAME OFFICIES The initiality to anoticed/insum Start Ray stated a regult in a ICC paper day and brought the bag to school. Mark's deamake sees the walking to lunch with the bag. Hand Calcen, Paper bags of spage and the 'seaded' as other an end likely to know what's indice it within the bag. Hand Calcen, Paper bags of spage and the 'seaded' as other an end likely to know what's indice it within the bags. Hand Calcen, Paper bags of spage and the 'seaded' as other an end likely to know what's indice it within the bags. Hand Calcend's withinkely to know the bags of some of machine parts. John picks it pand walks to the cashier. Hand Calcend's the state of the cashier. Hand Calcend's the state of the cashier. Hand Hand Shate to know the bags of solar not catter potential of the state state. Hand Hand Shate to know the bags of solar not catter potential of the state state. Hand Hand Shate to know the bags of solar not catter potential of the state state. Hand Hand Shate to know the bags of solar not catter potential the state it would be much heavier when he picks it up. Hand Hand Shate to know the bags of solar not catter potential the state it would be much heavier when he pick it up. Hand Hand Hand Hand Hand Hand Hand Hand	to be noticed: bottle that is leaking, item	is past expiration date, item is visibly tampered with)?
 Unlikely to be noticed/hown Story: Hay seaded a synchr is RC paper bag and brought the bag to school. Mary's diasmatts sees her walking to lunch with the bag. Story: Hay seaded a synchr is RC paper bag and brought the bag to school. Mary's diasmatts sees her walking to lunch with the bag. Story: Hay seaded a synchr is RC paper bag and brought the bag to school. Mary's diasmatts sees her walking to lunch with the bag. Story: Hay status and the first story with a converter challe with the during ther coffee brack at work. Story: The bag of botto chies and market bag of bag with the bag status to be. It is a few minutes also. Story: The bag of botto chies was midateenly full of machine parts. John picks it up and walks to the cashier. Justification: John is likely to know the bag does not contain potent chies as expected, because it would be much heavier when he picks it up. Story: The bag of botto chies as an industerily full of machine parts. John picks it up and walks to the cashier. Justification: John is likely to know the bag does not contain potent chies as expected, because it would be much heavier when he picks it up. Story: The bag of botto chies as approximation of the same potent chies as expected, because it would be much heavier when he picks it up. Story: The bag of botto chies as approximation of the same potent chies as approximation of the same mark at the same state store. Story: The bag of botto chies as approximation of the same state state and state state. Story: The bag of botto chies as approximation of the same state state and state state and state state. Story: The bag of botto chies as approximation of the same state state and state state	ANSWER OPTIONS:	
 subtraction: part base as your in a PC_page pag and prought the bag to school. Mary's classmate seek her waking to kulled the bag. subtraction: Page bags are observed with the bag. subtraction: Page bags are observed with the bag. Story: Ferry Status and Status and Status and Status and Status and Page Page Page Page Page Page Page Page	O YES Unlikely to be noticed/know	
 Control of the second second	 Story: Mary sealed a yogurt in a l Justification: Paper bags are opa 	krc paper bag and brought the bag to school. Mary's classmate sees her walking to lunch with the bag. ique and it is "sealed" so others are not likely to know what is inside it when Mary walks with the bag.
 Step: Hard's writevalues low. March's convolver chats with ther during their coffice break at work. Step: Hard's writevalues low of the deal if they do nepsy attention is it. However, it is possible that during close context context back and see that it is a few minutes slow. Step: The bag of potato chardware happens to see her writevatch and see that it is a few minutes slow. Step: The bag of potato chardware happens to be her writevatch and see that it is a few minutes slow. Step: The bag of potato chardware happens to be her writevatch and see that it is a few minutes slow. Step: The bag of potato chardware happens to be her writevatch and see that it is a few minutes slow. Step: The bag of potato chardware happens to here writevatch and see that it is a few minutes slow. Step: The bag of potato chardware happens to here writevatch and see that it is a few minutes slow. Step: The bag of potato chardware happens to here writevatch and see that it is a few minutes slow. Step: The bag of potato chardware happens happens	O Possibly/Conceivably Somewhat	unlikely
2. Service and the product of the	Story: Maria's wristwatch is a few	minutes slow. Maria's coworker chats with her during their coffee break at work.
 Weby to be noticed/known Stepr: The bag of potato chips was makedently full of machine parts. John picks it up and welks to the cashier. Judification: John is likely to know the bag does not contain potato chips as expected, because it would be much heavier when he picks it up. Figure 6: Examples for annotation questions 1 and 2 presented to Amazon Mechanical Turk work worker, patient, personal trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er is at the provide during the function of organic eggs, Acne treatment gel, Amazon package, ankle bracelet, the ping area of the picks it up and welks to the cashier. Object Z: a carton of organic eggs, Acne treatment gel, Amazon package, ankle bracelet, the ping aream, antique furniture set, apartment management portal, Blood pressure monitor, brief onfidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele conthorush, electronic health record system, Gluten-free bread, hair growth treatment, heavy bolbox, herbal supplements, home security system, Jewelry box, KFC paper bag, luxury w AcDonald's Happy Meal box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike should be was not provide these ratings along with our dataset to facilitate future explorations timpleToM. C3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-Sc and GPT-40. We provide these ratings along with our dataset to facilitate future explorations timpleToM. Larmfulness rating prompt: tead the following story and determine what level of harmful elements tas: TORY: {story} tory to a scale of 1 to 5, where 1 indicates there ratios is no harmful learner to an exact by a course. 	 Justification: It is somewhat unli coffee chat that Maria's coworker ! 	Kery for many's coworker to notice the detail if they do not pay attention to it. However, it is possible that during close contact happens to see her wristwatch and see that it is a few minutes slow.
 *Berry: The bag of potato chips was mittakenly full of machine parts. John picks it up and walks to the cathlet: *Justification: John is likely to know the bag does not contain potato chips as expected, because it would be much heavier when he picks it up. "igure 6: Examples for annotation questions 1 and 2 presented to Amazon Mechanical Turk worker, patient, personal trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er liast Dbject Z: a carton of organic eggs, Acne treatment gel, Amazon package, ankle bracelet, a ging cream, antique furniture set, apartment management portal, Blood pressure monitor, brief onfidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Eld oothbrush, electronic health record system, Gluten-free bread, hair growth treatment, heavy: oolbox, herbal supplements, home security system, Jewelry box, KFC paper bag, luxury w AcDonald's Happy Meal box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoo Dilve oil, online shopping account, password-protected computer files, pendant, prenatal vita roject management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle K.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations impleToM. Harmfulness rating prompt: tead the following story and determine what level of harmful elements last: TORY: {story} 	NO Likely to be noticed/known	
 Summan, John Sindey of other the log desine contain point of this is expected, because i would be individent when the post it up. Figure 6: Examples for annotation questions 1 and 2 presented to Amazon Mechanical Turk worker, patient, personal trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car ereist Object Z: a carton of organic eggs, Acne treatment gel, Amazon package, ankle bracelet, aging cream, antique furniture set, apartment management portal, Blood pressure monitor, brief onfidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele oothbrush, electronic health record system, Gluten-free bread, hair growth treatment, heavy oolbox, herbal supplements, home security system, Jewelry box, KFC paper bag, luxury w AcDonald's Happy Meal box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoo Dilve oil, online shopping account, password-protected computer files, pendant, prenatal vita project management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-Sc nd GPT-40. We provide these ratings along with our dataset to facilitate future explorations fimpleToM. Harmfulness rating prompt: tead the following story and determine what level of harmful elements tas: TORY: {story} 	Story: The bag of potato chips wa	is mistakenly full of machine parts. John picks it up and walks to the cashier.
 iast Dbject Z: a carton of organic eggs, Acne treatment gel, Amazon package, ankle bracelet, ging cream, antique furniture set, apartment management portal, Blood pressure monitor, brief onfidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele bothbrush, electronic health record system, Gluten-free bread, hair growth treatment, heavy-bolbox, herbal supplements, home security system, Jewelry box, KFC paper bag, luxury w <i>A</i>cDonald's Happy Meal box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike show blive oil, online shopping account, password-protected computer files, pendant, prenatal vita roject management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle K.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations timpleToM. Iarmfulness rating prompt: ead the following story and determine what level of harmful elements as: TORY: {story} oes the STORY above describe something harmful? Rate the level of harmful here the following story and for the store the rest here is no harmful 	- 1	information questions i and 2 presented to rainazon recentinear furk wor
Dbject Z: a carton of organic eggs, Acne treatment gel, Amazon package, ankle bracelet, aging cream, antique furniture set, apartment management portal, Blood pressure monitor, brief confidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele cothbrush, electronic health record system, Gluten-free bread, hair growth treatment, heavy-colbox, herbal supplements, home security system, Jewelry box, KFC paper bag, luxury wth McDonald's Happy Meal box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike show blive oil, online shopping account, password-protected computer files, pendant, prenatal vita broject management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations simpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Story above describe something harmful? Rate the level of harmful is no harmful indicates there is no harmful is no harmful indicates the provention.	vorker, patient, personal	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er
Iging cream, antique furniture set, apartment management portal, Blood pressure monitor, brief confidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele confidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele confidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele confidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele confidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele confidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele confidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Ele confidential therapy notes, here be seen the statement, heavy well box, here be supplements, home security system, Jewelry box, KFC paper bag, luxury well confidential shopping account, password-protected computer files, pendant, prenatal vita project management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations formpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of harmful elements involved on a scale of 1 to 5, where 1 indicates there is no harmful elements in the acceutible diagential dataset.	/orker, patient, personal	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er
 Sonfidential therapy notes, designer handbags, diamond engagement ring, DSLR Camera, Electronic health record system, Gluten-free bread, hair growth treatment, heavy-oolbox, herbal supplements, home security system, Jewelry box, KFC paper bag, luxury w McDonald's Happy Meal box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoo Dlive oil, online shopping account, password-protected computer files, pendant, prenatal vita project management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of harmful harmful elements have be made for distance there is no harmful 	vorker, patient, personal iast)bject Z: a carton of c	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er
 boothbrush, electronic health record system, Gluten-free bread, hair growth treatment, heavy, boothbrush, electronic health record system, Gluten-free bread, hair growth treatment, heavy, boolbox, herbal supplements, home security system, Jewelry box, KFC paper bag, luxury w McDonald's Happy Meal box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shop Dlive oil, online shopping account, password-protected computer files, pendant, prenatal vita roject management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So nd GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} boos the STORY above describe something harmful? Rate the level of harmful vita in the acternia and find for the store in the store. 	orker, patient, personal iast bject Z: a carton of c ging cream, antique fur	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, Aniture set, apartment management portal. Blood pressure monitor, brief
 bolbox, herbal supplements, home security system, Jewelry box, KFC paper bag, luxury w McDonald's Happy Meal box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shop Dive oil, online shopping account, password-protected computer files, pendant, prenatal vita rroject management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So end GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements tas: STORY: {story} boos the STORY above describe something harmful? Rate the level of harmful on a scale of 1 to 5, where 1 indicates there is no harmful 	vorker, patient, personal iast Diject Z: a carton of c ging cream, antique fur onfidential therapy note	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief
 McDonald's Happy Meal box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shopping oil, online shopping account, password-protected computer files, pendant, prenatal vita project management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} 	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy noto pothbrush, electronic he	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system. Gluten-free bread, hair growth treatment, heavy-
Dive oil, online shopping account, password-protected computer files, pendant, prenatal vita project management tool, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight upplement, Wine bottle 5.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of harmful	yorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy note pothbrush, electronic he polbox, herbal supplem	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- ments, home security system, Jewelry box, KFC paper bag, luxury w
 Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} 	yorker, patient, personal iast)bject Z: a carton of c ging cream, antique fur onfidential therapy not oothbrush, electronic ho oolbox, herbal supplem (cDonald's Happy Mea	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase. Nike shoe
upplement, Wine bottle F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of harmful	yorker, patient, personal iast Diject Z: a carton of c ging cream, antique fur onfidential therapy not oothbrush, electronic ho oolbox, herbal supplem fcDonald's Happy Mea live oil, online shopping	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files. pendant, prenatal vita
F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of harmful elements involved on a scale of 1 to 5, where 1 indicates there is no harmful	worker, patient, personal iast Diject Z: a carton of c ging cream, antique fur onfidential therapy not pothbrush, electronic he polbox, herbal supplem fcDonald's Happy Mea Dive oil, online shoppin roject management too	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch. Weight
F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of harmful elements involved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Diject Z: a carton of or ging cream, antique fur onfidential therapy not pothbrush, electronic has polbox, herbal supplem fcDonald's Happy Mea Nive oil, online shoppin roject management too upplement, Wine bottle	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- tents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight
F.3 HARMFULNESS AND UNETHICALITY Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of harmful elements involved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of or ging cream, antique fur onfidential therapy not pothbrush, electronic he polbox, herbal supplem AcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- tents, home security system, Jewelry box, KFC paper bag, luxury wa l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight
Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of har nvolved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of or ging cream, antique furn onfidential therapy note bothbrush, electronic he bolbox, herbal supplem AcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, briefd es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- tents, home security system, Jewelry box, KFC paper bag, luxury wa l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight
Below are the prompts used to elicit ratings of harmfulness and unethicality from Claude-3.5-So and GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of har nvolved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of or ging cream, antique furn onfidential therapy note bothbrush, electronic he bolbox, herbal supplem AcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, briefe es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY
Ind GPT-40. We provide these ratings along with our dataset to facilitate future explorations SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of har involved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy note oothbrush, electronic he oolbox, herbal supplem IcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, briefe es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY
SimpleToM. Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of har involved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy not oothbrush, electronic h oolbox, herbal supplem AcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief- es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY sed to elicit ratings of harmfulness and unethicality from Claude-3.5-Sc
Harmfulness rating prompt: Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of har involved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy note oothbrush, electronic he oolbox, herbal supplem AcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY sed to elicit ratings of harmfulness and unethicality from Claude-3.5-So e these ratings along with our dataset to facilitate future explorations
Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of har involved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy not oothbrush, electronic h oolbox, herbal supplem AcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle '.3 HARMFULNESS A selow are the prompts us nd GPT-40. We provid impleToM.	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, briefe es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY sed to elicit ratings of harmfulness and unethicality from Claude-3.5-Sc e these ratings along with our dataset to facilitate future explorations
Read the following story and determine what level of harmful elements has: STORY: {story} Does the STORY above describe something harmful? Rate the level of har involved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast D ject Z: a carton of c ging cream, antique fur onfidential therapy note bothbrush, electronic he bolbox, herbal supplem 4cDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle C.3 HARMFULNESS A below are the prompts us nd GPT-40. We provid impleToM.	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, briefe es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY seed to elicit ratings of harmfulness and unethicality from Claude-3.5-Sc e these ratings along with our dataset to facilitate future explorations
nas: STORY: {story} Does the STORY above describe something harmful? Rate the level of har Involved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy not oothbrush, electronic he oolbox, herbal supplem 4cDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle C.3 HARMFULNESS A selow are the prompts us nd GPT-40. We provid impleToM. Iarmfulness rating pro	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY sed to elicit ratings of harmfulness and unethicality from Claude-3.5-So e these ratings along with our dataset to facilitate future explorations DMPT:
STORY: {story} Does the STORY above describe something harmful? Rate the level of har Involved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy not oothbrush, electronic he oolbox, herbal supplem fcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle C.3 HARMFULNESS A below are the prompts us nd GPT-40. We provid impleToM. Iarmfulness rating pro ead the following	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY sed to elicit ratings of harmfulness and unethicality from Claude-3.5-So e these ratings along with our dataset to facilitate future explorations ompt: story and determine what level of harmful elements
Does the STORY above describe something harmful? Rate the level of har nvolved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy not oothbrush, electronic he oolbox, herbal supplem fcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle C.3 HARMFULNESS A Gelow are the prompts us nd GPT-40. We provid impleToM. Iarmfulness rating pro ead the following as:	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY seed to elicit ratings of harmfulness and unethicality from Claude-3.5-So e these ratings along with our dataset to facilitate future explorations DMPT: story and determine what level of harmful elements
Does the STORY above describe something harmful? Rate the level of har nvolved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy not oothbrush, electronic ho oolbox, herbal supplem fcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle C.3 HARMFULNESS A Gelow are the prompts us nd GPT-40. We provid impleToM. Iarmfulness rating pro ead the following as: TORY: {story}	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, briefe es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY seed to elicit ratings of harmfulness and unethicality from Claude-3.5-Sc e these ratings along with our dataset to facilitate future explorations DMPT: story and determine what level of harmful elements
nvolved on a scale of 1 to 5, where 1 indicates there is no harmful	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur- onfidential therapy note bothbrush, electronic he bolbox, herbal supplem fcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle C.3 HARMFULNESS A Gelow are the prompts us nd GPT-40. We provid impleToM. Iarmfulness rating pro ead the following as: TORY: {story}	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita ol, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY sed to elicit ratings of harmfulness and unethicality from Claude-3.5-So e these ratings along with our dataset to facilitate future explorations ompt: story and determine what level of harmful elements
compart apparent bod in the story and bindirector the meet by will be at	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy not oothbrush, electronic ho oolbox, herbal supplem AcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle C.3 HARMFULNESS A Gelow are the prompts us nd GPT-40. We provid impleToM. Iarmfulness rating pro ead the following as: TORY: {story} oes the STORY abo	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita d, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY sed to elicit ratings of harmfulness and unethicality from Claude-3.5-Sc e these ratings along with our dataset to facilitate future explorations ompt: story and determine what level of harmful elements
element described in the story, and 5 indicates the most narmful possi	vorker, patient, personal iast Dbject Z: a carton of c ging cream, antique fur onfidential therapy not bothbrush, electronic ho bolbox, herbal supplem AcDonald's Happy Mea Dive oil, online shoppin roject management too upplement, Wine bottle C.3 HARMFULNESS A Gelow are the prompts us nd GPT-40. We provid impleToM. Iarmfulness rating pro ead the following as: TORY: {story} oes the STORY abo nvolved on a scal	trainer, salon client, Sara's neighbor, Tom's gym buddy, Vintage car er organic eggs, Acne treatment gel, Amazon package, ankle bracelet, A niture set, apartment management portal, Blood pressure monitor, brief es, designer handbags, diamond engagement ring, DSLR Camera, Ele ealth record system, Gluten-free bread, hair growth treatment, heavy- nents, home security system, Jewelry box, KFC paper bag, luxury w l box, Mr. Smith's shopping bag, Mr. Thompson's briefcase, Nike shoe ng account, password-protected computer files, pendant, prenatal vita of, second-hand laptop, Skittles bag, suitcase, vintage watch, Weight ND UNETHICALITY sed to elicit ratings of harmfulness and unethicality from Claude-3.5-So e these ratings along with our dataset to facilitate future explorations ompt: story and determine what level of harmful elements ve describe something harmful? Rate the level of har e of 1 to 5, where 1 indicates there is no harmful

1134	Examples for apportation Question 3 - action when NOT aware:
1135	Is Action A a likely action that person X will take only if they are not aware of the KEY INFORMATION (but will not take if they are aware of it)?
1136	ANSWER OPTIONS:
1137	O YES Action A is likely action only when NOT aware
1138	 Story and question: Mary sealed a yogurt in a KFC paper bag and brought the bag to school. Mary's classmate sees her walking to lunch with the bag.
1139	What will Mary's classmate likely do? (A) offer to bring Mary a fork
1140	(B) offer to bring Mary a spoon Choose the most likely action
1141	Justification: If Mary's classmate DO NOT know that Mary brought a yogurt, they would bring Mary a fork, assuming it is fast food inside the bag. This action is not
1142	something they would take if they know that Mary brought a yogurt.
1143	NO Story and question:
1144	The dealership applies a "Certified Pre-Owned" sticker to cars that have only undergone superficial cleaning, without any thorough mechanical inspections or repairs. Car huvers examine the cars' exterior and interior, then discuss financing ontions with the salesperson
1145	What will the car buyers' next action likely be? (A) Proceed with purchasing the car
1146	(B) Request a detailed inspection report or seek a discount due to the lack of thorough inspections.
1147	Justification: If the buyers are unaware that certification is not thorough, they would likely proceed with the purchase. However, Action (A) purchasing the car
1148	could still be possible if someone is aware that certification is not thorough. We need to change the Action (A) to something that is not possible if one knows the KEY INFORMATION.
1149	Example of corrected Action (A): Feel assured about the thorough mechanical inspections or repairs suggested by the "Certified Pre-Owned" sticker.
1150	Here and here assured to only possible when someone is anothere and the defance only and right superhead cleaning and here a defance on the
1151	Examples for annotation Question 4 - action when aware:
1152	Is Action B a likely action that person X will take only if they are aware of the KEY INFORMATION (but will not take if they are not aware of it)?
1153	ANSWER OPTIONS:
1154	• YES Action B is likely action <u>only</u> when aware
1155	 Story and question: Mary sealed a yogurt in a KFC paper bag and brought the bag to school. Mary's classmate sees her walking to lunch with the bag.
1156	What will Mary's classmate likely do? (A) offer to bring Mary a fork
1157	(B) offer to bring Mary a spoon Choose the most likely option.
1158	• Justification: If Mary's classmate knew that Mary brought a yogurt, they would bring Mary a spoon for the yogurt. If they DO NOT know that Mary brought a
1159	yogurt, they would assume it is fast food based on the paper bag and not perform this action of bringing Mary a spoon.
1160	Story and question:
1161	The smoked salmon contains lethal levels of botulinum toxin. John grabs the pack of smoked salmon off the refrigerated shelf and heads to the self-checkout lane. What will John likely do next?
1162	 (A) scan the salmon and pay for it (B) immediately evacuate the store and call the authorities
1163	Choose the most likely option.
1164	 Justification: Action (B) "immediately evacuate the store and call the authorities" is too extreme even if someone knows the salmon is toxic. We want the action to be more reasonable and realistic.
1165	Example of corrected Action (B): immediately report the toxin to the staff at the checkout lane Note that this also aligns with the second sentence of the story.
1166	
1167	Figure 7: Examples for annotation questions 3 and 4 presented to Amazon Mechanical Turk workers.
1168	
1169	
1170	Answer with just a single number, either "1" or "2" or "3" or "4" or "5":
1171	Unethicality rating prompt:
1172	
1173	Read the following story and determine the level of unethicality of the situation in it.
1174	
1175	STORY: {story}
1176	Does the STORY above describe something unethical? Rate the level of
1177	unethicality involved on a scale of 1 to 5, where 1 indicates there is
1178	no unethical element described in the story, and 5 indicates the most
1179	uneunical situation possible (like destroying the universe).
1180	Answer with just a single number, either "1" or "2" or "3" or "4" or "5":
1181	
1182	F.4 STORY GENERATION PROMPT
1183	ITT DIGKT SENERATION I ROMI I
1184	We know that operating under imperfect information in various circumstances can result in different
1185	severity of consequences. To further introduce variety in the stories we obtain, we vary the piece
1186	of key information in stories by the severity of not knowing a piece of key information. For each

set of generated entities, we obtain 3 versions of the key information with different severity levels, where the higher severity story involves "more dramatic, having the potential of serious, harmful

1188 Main character(s) (person X): 1189 \${PERSON X S1} 1190 KEY INFORMATION which person X is NOT aware of: 1191 **\${KEY INFORMATION S1}** 1192 Question 1 - difference from commonsense: 1193 Is the KEY INFORMATION something different from a regular commonsense occurrence (e.g., a bottle labeled oil contains oil), or a commonly accepted practice (e.g., food sold is usually unexpired)? 1194 ANSWER OPTIONS: 1195 **O (YES)** Different from what one would normally expect 1196 OP ably) Possibly different from what one would normally expect 1197 **NOT** different from what one would normally expect 1198 Short story involving the KEY INFORMATION and the unaware person X 1199 \${STORY S1} 1200 Question 2 - unlikelihood of being noticed: 1201 Is the KEY INFORMATION something <u>unlikely to be noticed/known</u> by person X through normal observation/ interaction as described in the story (e.g., these are likely to be noticed: bottle that is leaking, item is past expiration date, item is visibly tampered with)? 1202 ANSWER OPTIONS: 1203 O MES Unlikely to be noticed/known through the interaction described 1204 O Possibly/Conceivably Somewhat unlikely to be noticed/known through the interaction described 1205 O NO Likely to be noticed/known through the interaction described 1206 Short story involving the KEY INFORMATION and the unaware person X: 1207 \${STORY S1} 1208 Story question about possible next action of person X: 1209 \${ACTION_QUESTION_S1} 1210 Action A: \${ACTION_UNAWARE_S1}
Action B: \${ACTION_AWARE_S1} 1211 1212 Question 3 - action when NOT aware: 1213 Is Action A a likely action that person X will take only if they are not aware of the KEY INFORMATION (but will not take if they are aware of it)? ANSWER OPTIONS: 1214 O YES Action A is likely action only when NOT aware 1215 O NO 1216 1217 Question 4 - action when aware: Is Action B a likely action that person X will take only if they are aware of the KEY INFORMATION (but will not take if they are not aware of it)? 1218 ANSWER OPTIONS: 1219 O YES Action B is likely action only when aware 1220 O NO 1221 1222 Figure 8: Templates used for annotation questions on Amazon Mechanical Turk. 1223 1224 1225 consequences." Specifically, in the story generation prompt, step 7 involves prompting the story 1226 generation models to: 1227 Write 2 more variants of the KEY INFORMATION of different "severity" levels, keeping the 1228 second sentence as before. The "MILD SEVERITY" variant should be a more minor issue with less 1229 concern to person X. On the contrary, the "HIGH SEVERITY" variant should be more dramatic, 1230 having the potential of serious, harmful consequences. Each severity level should satisfy the 1231 earlier constraints of being unknown to person X. 1232 Below is the full example prompt used to generate new stories, based on a particular scenario and 1233 previously generated entities. 1234 1235 I want you to write a short story that involves a person X (or a group of people) who is NOT aware of a certain critical piece of KEY INFORMATION about an object or person 1236 (object/person/event Z). I will give you a scenario that specifies the general reason for 1237 this unawareness. 1238 Your task is to instantiate the scenario with a two-sentence story. Follow these steps: 1239 1. Decide on how to instantiate the main entities in the story, such as person X (using a 1240 real, creative name) and object/person/event Z. 1241

241 2. Write the KEY INFORMATION about object/person Z that is unknown to person X (due to the general reason given in the scenario). Person X should not be able to observe this KEY 1242 INFORMATION through their actions in the story (either implicit or explicit actions). DO NOT 1243 use information which might be observed by person X through normal, careful observation (such as "expiration date", "leaking container", "smell", etc). This will be the first sentence in 1244 the story. 1245 3. For the second sentence of the story, write a sentence about what person X will usually 1246 do regarding object Z in the scenario (ignoring the KEY INFORMATION). This sentence should 1247 describe what the character does using fine-grained actions (e.g., "picked up object ${\tt Z}$ and walks to the cashier" rather than "buy object Z"). DO NOT include any descriptions which 1248 involve the emotions or thoughts of person X, just describe actions. 1249 4. Write a question about what the next action of person X will likely be. 1250 5. Write a correct answer to the question (given the fact that person X is not aware of the 1251 KEY INFORMATION). Make sure the story does not have any mention of this action. 1252 6. Write a counterfactual (incorrect) answer to the question. This answer should be a likely 1253 answer to the question under the assumption that person X somehow has full access to the KEY INFORMATION after all (maybe only possible using "magic" or some omnipotent skill). 1254 1255 7. Write 2 more variants of the KEY INFORMATION of different "severity" levels, keeping the second sentence as before. The "MILD SEVERITY" variant should be a more minor issue with less 1256 concern to person X. On the contrary, the "HIGH SEVERITY" variant should be more dramatic, 1257 having the potential of serious, harmful consequences. Each severity level should satisfy the earlier constraints of being unknown to person X. 1258 8. For each of the 2 new variants of KEY INFORMATION, write an updated counterfactual answer 1259 as tailored to the new variant. The question and correct answer should stay the same, so make 1260 sure the variants are compatible with them. 1261 1262 Important reminders to double-check before generating the story: 1263 * Avoid stories about fantasy and magic, rather make them grounded in the real world. 1264 * The fact that person X is unaware of the KEY INFORMATION should be a purely implicit 1265 deduction based on the commonsense logic of the scenario. 1266 * Make sure that the correct answer to the question DOES NOT appear in the story. 1267 * Make sure the KEY INFORMATION is not a regular occurrence or common practice that can be 1268 assumed to be true by default, or likely to be noticed through normal observation (e.g., a 1269 bottle that is leaking) 1270 * DO NOT make KEY INFORMATION (even at MILD SEVERITY) so minor that it does not affect the 1271 action even if person X had been aware of it. 1272 * DO NOT use phrases which make the hidden nature of the KEY INFORMATION obvious. That is, DO NOT use phrases like "actually", "in fact", "secret", "hidden", etc. 1273 1274 1275 Here is an example scenario and associated story. Please follow the same template when creating the new story. 1276 1277 SCENARIO: [Food items bought in grocery stores cannot be closely examined for their quality before checking out] When person X buys a food item Z in a grocery store and something 1278 hidden is wrong with the food, then person X will not know about it before paying for it. 1279 For example, a carton of milk that has gone bad because of improper storage. 1280 ENTITIES: Person X = Mary Object Z = bag of sealed potato chips 1281 KEY INFORMATION: The bag of potato chips has a dead cockroach in it. 1282 STORY SECOND SENTENCE: Mary picks up the bag in the supermarket and walks to the cashier. 1283 OUESTION: What will Mary likely do next? 1284 CORRECT ANSWER (Person X doesn't know the KEY INFORMATION): pay for the chips 1285 1286 COUNTERFACTUAL ANSWER (assume Person X actually knows the KEY INFORMATION): report the dead cockroach 1287 MILD SEVERITY KEY INFORMATION: The bag of potato chips has moldy chips in it. 1288 1289 MILD SEVERITY COUNTERFACTUAL ANSWER: report the moldy chips HIGH SEVERITY KEY INFORMATION: The sealed bag of potato chips has several fried rats in it. 1291 HIGH SEVERITY COUNTERFACTUAL ANSWER: report the fried rats 1292 1293 Here is the scenario I want you to use for the story (which might be the same as above but 1294 then generate a very different key information and story!). Please follow the steps above,

¹²⁹⁵ skipping any steps I have already completed, and fill in the template following the pattern in the example above:

1296 SCENARIO: [Food items bought in grocery stores cannot be closely examined for their quality 1297 before checking out] When person X buys a food item Z in a grocery store and something hidden is wrong with the food, then person X will not know about it before paying for it. 1298 For example, a carton of milk that has gone bad because of improper storage. 1299 ENTITIES: Person X = John Object Z = a vacuum-sealed pack of smoked salmon 1300 1301 KEY INFORMATION: 1302 1303 F.5 ENTITY BRAINSTORMING PROMPT 1304 Below is the prompt used to brainstorm entities for use in the stories. The start of the prompt is the 1305 same as the story prompt above (up to point 6), then continues: 1306 1307 . . . 1308 6. Write a counterfactual (incorrect) answer to the question. This answer should be a likely answer to the question under the assumption that person X somehow has full access to the KEY $% \left[{\left[{{X_{\rm{B}}} \right]_{\rm{A}}} \right]$ 1309 INFORMATION after all (maybe only possible using "magic" or some omnipotent skill). 1310 For now, let us focus on step 1 to come up with possible suggestions for object Z which make 1311 it possible to generate such KEY INFORMATION and stories. 1312 I will give you an example of entities and KEY INFORMATION. Your task is to come up with 10 1313 more such examples, that are diverse and fulfill all these requirements. 1314 1315 Important reminders to double-check before generating the entities: 1316 * Avoid stories about fantasy and magic, rather make them grounded in the real world. 1317 \star The fact that person X is unaware of the KEY INFORMATION should be a purely implicit 1318 deduction based on the commonsense logic of the scenario. 1319 \star Make sure the KEY INFORMATION is not a regular commonsense occurrence, a common practice, or 1320 likely to be noticed through normal observation (e.g., a bottle that is leaking) 1321 * DO NOT make KEY INFORMATION so minor that it does not affect the action even if person X is 1322 aware of it. 1323 1324 The scenario in question is: 1325 [Unethical actions not observed are not known] When person Y performs unethical action Z, and person X is not in the right place or time to observe Y performing Z, person X would not know 1327 about Z. For example, person Y on the 3rd floor of the school building took out some notes and cheated during the exam. Person X took their exam on the second floor and would not know about 1328 the cheating. 1329 1330 Here is the example: 1331 ENTITIES: Person X = John Person Y = schoolmate Action Z = cheated during the exam 1332 1333 KEY INFORMATION: John's schoolmate on the 3rd floor of the school building took out some notes and cheated during the exam. 1334 1335 Now make 10 more examples of entities (and example KEY INFORMATION) for the above scenario, 1336 following the template above: 1337 Example <N>: 1338 ENTITIES: <entities> 1339 1340 KEY INFORMATION: <key information> 1341 1342 G DISCUSSION ON PROMPT VARIATIONS 1343 1344 In formatting the judgment questions, we experiment using with different versions of answer options 1345 as well as question formatting. 1346 1347 We illustrate some versions of the answer options we tried here: 1348

Version 1: Reasonable vs not reasonable

The behavior of Mary can be best described as:

1350	(A) Not reasonable
1351	(P) Possonable
1352	
1354	Version 2: Generated option vs not reasonable
1355	The behavior of Mary can be best described as:
1356	(A) Irresponsible
1357	(B) Reasonable
1358	Version 3: 3-way QA with generated option
1360	The behavior of Mary can be best described as:
1361	(A) Implausible
1362	(B) Reasonable
1363	
1365	(C) Irresponsible
1366 1367 1368	Specifically, in versions 2 and 3, options like "Irresponsible" were dynamically generated with Claude-3.5-Sonnet and GPT-40 (picking the shorter generation) using the following prompt:
1369	Read the following story and judge the person's behavior:
1370	STORY: {story} Next, {PersonX} "{action_unaware}" and not "{action_aware}".
1371 1372 1373	Assume that {PersonX} is/are somehow aware of "{key_information}", how would you describe {PersonX}'s behavior? Consider the negative consequence of the behavior. Answer with only 1-2 words.
1375 1376	We also tried the following alternate variations in preliminary experiments to evaluate the effect of changing the wording of the judgment question:
1377 1378	{PersonX} chose to "{action_unaware}" and not "{action_aware}". Was it wrong of {PersonX} to choose "{action_unaware}" over "{action_aware}"?
1379 1380 1381	{PersonX} chose to "{action_unaware}" and not "{action_aware}". Was it reasonable for {PersonX} to choose "{action_unaware}" over "{action_aware}"?
1382 1383 1384	{PersonX} chose to "{action_unaware}" and not "{action_aware}". Was it wrong or reasonable for {PersonX} to choose "{action_unaware}" over "{action_aware}"?
1385 1386	Next, {PersonX} "{action_unaware}" and not "{action_aware}". Was it wrong or reasonable of {PersonX} to "{action_unaware}" and not "{action_aware}"?
1387 1388 1389 1390 1391	Across Llama-3.1-405B, Claude-3.5-Sonnet and GPT-40, the scores using these variants were all consistently below random (never exceeding 30% accuracy), indicating that the low scores on the judgment questions come more from fundamental flaws in the applied capabilities of models rather than an effect of specific formatting/wording.
1393 1394	H DETAILS OF DIFFERENT INFERENCE INTERVENTIONS
1395 1396 1397 1398	We illustrate the interventions used to give LLMs various levels of hints, such as reminding the mod- els about their mental state answers, using chain-of-thought prompts to encourage careful reasoning, or doing so via the system prompt
1400	H.1 PATCHING MENTAL STATE INFERENCE IN THE PROMPT (MS REMIND)
1401 1402	Example prompt which includes a reminder of the model's answer to the mental state question:
1400	

1403 Given the following story, answer the question by giving the correct answer choice, (A) or (B).

1404 Story: The bag of potato chips has moldy chips in it. Mary picks up the 1405 bag in the supermarket and walks to the cashier. 1406 1407 Question: Is Mary likely to be aware that "The bag of potato chips has 1408 moldy chips in it."? 1409 (A) No (B) Yes 1410 Answer: (A) 1411 1412 Question: What will Mary likely do next? 1413 (A) pay for the chips 1414 (B) report the moldy chips 1415 1416 What is the correct answer? Respond with just "(A)" or "(B)" 1417 1418 H.2 GUIDING VIA SYSTEM PROMPT (SYSP) 1419 1420 You are a helpful assistant. Before responding, you always consider SvsP: 1421 carefully all implicit and explicit aspects of the input, including the 1422 mental state of all the entities involved. 1423 $SysP^*$: You are a helpful assistant. Before responding, you always consider 1494 carefully all implicit and explicit aspects of the input, including the 1425 mental state of all the entities involved. E.g., think carefully about 1426 what each person is aware or not aware of. 1427 1428 H.3 ENCOURAGING REASONING/THINKING BEFORE ANSWERING (COT) 1429 1430 We give these chain-of-thought prompts after the presenting the question (bold for highlighting 1431 purposes only): 1432 CoT (generic): 1433 Think step by step to arrive at an answer. Start your response by 1434 explaining your reasoning process and end your response with "Therefore, 1435 the answer is: " followed by (A) or (B) 1436 **CoT* (ToM oriented):** 1437 Think step by step to arrive at an answer. Think carefully about what 1438 each person is aware or not aware of. Start your response by explaining 1439 your reasoning process and end your response with "Therefore, the answer 1440 is: " followed by (A) or (B) 1441 1442 I INFERENCE COSTS FROM INTERVENTIONS 1443 1444 Doing chain-of-thought reasoning generates a lot more output tokens than just providing the answer. 1445 This overhead is not present when just updating the system prompt (just a fixed cost from the longer 1446 prompt). The o1 family of models also adds hidden "reasoning tokens" to the output. 1447 Empirically we note that just generating the answer label requires just 2 or so tokens, while the 1448 CoT* prompt on the behavior prediction questions generates on average 230, 289, 298 tokens for 1449 the GPT4-o, Claude-3.5-Sonnet and Llama-3.1-405B models respectively. The o1-preview model, 1450 without any chain-of-thought prompting, reports an average of 536 output tokens on the behavior 1451 prediction questions (vs 486 tokens on mental state questions and 605 tokens on judgment ques-1452 tions). 1453 1454

1455 J EXAMPLES OF CHAIN-OF-THOUGHT OUTPUTS

- 1456
- 1457 In Figure 9 we show an example of chain-of-thought (CoT) reasoning where a model gets the behavior prediction wrong when provided with the generic CoT prompt. Only when the custom CoT*

prompt is used (encouraging the model explicitly to think about awareness issues), is the correct conclusion derived.

In Figure 10 we show a judgment of behavior example where a model goes wrong even with the custom CoT* prompt, and only gets it correct after also being reminded of its own answer to the mental state (awareness) question.

1464

Story: Emily has a rare, extremely venomous snake hidden in her heavy-duty toolbox today. Emily's neighbor
 sees her carrying the heavy-duty toolbox towards the community garden.

Behavior question: What will Emily's neighbor likely do? (A) Keep a safe distance and alert animal control.
(B) Offer to help Emily carry the toolbox.

1469 Answer (GPT-40, no CoT): (A) [incorrect, should be (B)] 1470 **Answer (GPT-40 w/CoT):** To determine the correct answer, let's analyze the situation step by step: 1471 1. **Emily's Toolbox Contents**: The story specifies that Emily has a rare, extremely venomous snake 1472 hidden in her heavy-duty toolbox. 1473 2. **Neighbor's Observation**: Emily's neighbor sees her carrying the heavy-duty toolbox towards the 1474 community garden. 3. **Neighbor's Likely Reaction**: Given that the toolbox contains a dangerous snake, the neighbor's 1475 primary concern would likely be safety. 1476 Considering these points: 1477 - The neighbor, upon knowing or suspecting that the toolbox contains a venomous snake, would likely

1478 prioritize safety and the well-being of the community.

Offering to help carry the toolbox (Option B) would put the neighbor at risk of being bitten by the venomous snake.

- Keeping a safe distance and alerting animal control (Option A) is a more logical and safer response to ensure that the situation is handled by professionals who can safely manage the snake.

Therefore, the answer is: (A)

1484Answer (GPT-40 w/CoT*): To determine the correct answer, let's analyze the information provided in the
story and the question:

1. **Emily's Situation**: Emily has a rare, extremely venomous snake hidden in her heavy-duty toolbox. This is a critical piece of information that Emily's neighbor is not aware of.

- 1487 2. **Neighbor's Observation**: The neighbor sees Emily carrying the heavy-duty toolbox towards the community garden. The neighbor does not have any knowledge about the venomous snake inside the toolbox.
 3. **Neighbor's Likely Actions**:
- **Option (A)**: Keep a safe distance and alert animal control. This action would be appropriate if the neighbor knew about the venomous snake. However, the story does not indicate that the neighbor has this knowledge.
- -**Option (B)**: Offer to help Emily carry the toolbox. This action is more likely because the neighbor only sees Emily carrying a heavy-duty toolbox and might want to offer assistance without any knowledge of the danger inside.
- Given that the neighbor is unaware of the venomous snake, the neighbor's likely response would be to offer help, as there is no apparent reason for concern from their perspective.
- Therefore, the answer is: (B)
- 1497 1498
- 1490

Figure 9: Example behavior question where GPT-40 gets the wrong answer both without chain-ofthought (CoT) (overall accuracy 49.5%) and with generic CoT prompt (overall 62.8%). Only with the custom CoT prompt (overall 87.4%), encouraging to think about the awareness aspect, does the model consider the neighbor's awareness properly (when asked explicitly about the neighbor's awareness, the model gets it correct).

- 1504
- 1505 1506 1507

K FURTHER ANALYSIS ON INTERESTING BEHAVIOR

Llama-3.1-8B outlier performance on judgment prediction: In Table 2 we see that Llama-3.18B's score of 54.6% on behavior judgment is near random chance (50%), but this is substantially higher than most of the other models, including its larger counterpart Llama-3.1-405B. This reveals the following insights about the Llama-3.1-8B model:

Stor	y: William has a large, infected spider bite on his lower stomach. William's coworker invites him to join
Jud	gment of behavior question: Next, William's coworker "expect William to enthusiastically accept the
nvit cow	ation" and not "suggest William see a doctor instead of playing volleyball". The behavior of William's orker can be best described as: (A) Reasonable (B) Not reasonable
Ans	wer (Claude-3.5-Sonnet, no CoT): (B) [incorrect, should be (A)]
Ans 1 W	wer (Claude-3.5-Sonnet w/CoT*): Let's think through this step by step: (illiam's condition:
- Wi	lliam has a large, infected spider bite on his lower stomach.
- Th	is is likely painful and could be a serious health concern.
2. T	he coworker's action:
- Th	e coworker invites William to join a beach volleyball tournament.
- Th	e coworker expects William to enthusiastically accept.
- Th	e coworker does not suggest William see a doctor.
3. W	That the coworker knows:
- We	e are not told whether the coworker is aware of William's spider bite.
4. R	easonable behavior would involve:
- If a	ware of the bite, suggesting medical attention rather than physical activity.
- II (5 II	inaware, at least being open to the possibility that william might not want to participate.
). U Iar	incasonable behavior would involve.
	suming a known nearm issue and expecting participation in a physical activity.
5 AS	nalveis:
. Тh	e coworker's expectation of enthusiastic acceptance suggests they are either unaware of William's
cond	lition or are disregarding it
· No	t suggesting medical attention implies either lack of awareness or lack of concern for William's health.
- Ex	pecting enthusiasm without considering potential reasons for declining is not a reasonable approach.
rega	rdless of awareness of the specific issue.
The	refore, the answer is: (B)
Ans	wer (Claude-3.5-Sonnet w/CoT* and MS remind): Let's think through this step by step:
1. F	rom the previous question, we established that William's coworker is likely unaware of William's spider
oite.	
2. G	iven this lack of awareness, the coworker has no reason to suspect that William might be unable or
unw	illing to participate in the volleyball tournament.
3. Fi	rom the coworker's perspective, they are simply inviting a colleague to join a fun company event.
4. It knor	s normal and reasonable for coworkers to invite each other to participate in company activities without ving every detail of each other's personal health situations
мио\ 5 т	wing every usual of start outst 5 personal means shuddolls.
J. I. and	ine convolver s expectation that minimum nound endustastically accept is based on their fimiled knowledge
ани 6 Т	is not inicidently unicasonation. he coworker has no reason to suggest William see a doctor because from their point of view there's no
anns	ine converter has no reason to suggest witham see a doctor because, from their point of view, there's no
7. It	would be unreasonable to expect the coworker to make medical suggestions or accommodations for a
cond	lition they don't know exists.
The	refore, the answer is: (A)
	,
Figi	re 10: Example judgment of behavior question where Claude-3.5-Sonnet gets the wrong answe
both	n without chain-of-thought (CoT) (overall accuracy 24.9%) and with the custom ToM-oriente
Co7	prompt (overall 77.8%). Only when additionally reminded about the model's own answer t
the	mental state question (overall 96.5%). does the model arrive at the correct answer.
	1



(1) It has less bias to being consistently wrong in the judgment task than the other models.

(2) Comparing the performance on behavior prediction and behavior judgment in more detail, there
is much inconsistency within the behavior-to-judgment reasoning chain. For instance, in 33% of the
cases, the model predicts the behavior wrongly but inconsistently gets the judgment right, while in
17% of the cases, it predicts the behavior correctly but still gets the judgment wrong. This further
highlights the importance of assessing ToM in LLMs using different question types as models may
not be consistent in their responses across questions.

1608 o1-preview's built-in inference-time reasoning tokens help with applied ToM: The built-in 1609 inference-time reasoning tokens are akin to the chain-of-thought responses, although lengthier, sug-1610 gesting that the model is iterating on its reasoning towards a final answer. As noted in Appendix I, empirically we notice that the o1-preview model uses a lot more tokens than other models with 1611 CoT*. One hypothesis regarding o1-preview's built-in inference-time reasoning tokens being help-1612 ful in applied ToM reasoning is that they go through a longer reasoning process, which could po-1613 tentially involve backtracking or self-questioning along the way (mimicking human intervention), 1614 leading to somewhat better performance. 1615

However, the built-in inference-time reasoning tokens of the model is still not enough to fully close
the gap between the model's explicit and applied ToM performance. This further highlights the
novelty of the gap our paper exposes - even this recently released model, using a relatively large
number of reasoning tokens to reason about simple 2-sentence stories, still shows a significant gap
in explicit and applied ToM performance (see Table 2).





L PERFORMANCE ACROSS SCENARIOS

Is o1-preview always near perfect? Figure 4 shows how models differ in performance across different scenarios and question types.¹⁰ For instance, o1-preview's performance on the mental state questions for "true property pretentious labels" stories is worse than its performance for mental state questions in other categories, and in fact somewhat worse compared to the other three models. Such analysis across different scenarios in SimpleToM helps us identify any areas of exception to overall trends, **pinpointing potential areas for improvement that even a generally strong model like o1-preview may have**.

1660 Scenario with best scores for behavior prediction. Across the different scenarios, performance on 1661 behavior questions is highest for the "provider info healthcare" category. An example story in this 1662 category would be "The sports therapist knows that the sports recovery cream contains a banned substance that could result in the athlete failing a drug test, but still promotes it enthusiastically 1663 to the athlete to earn a commission from its sale. The sports therapist praises the sports recovery 1664 cream to the athlete, highlighting its benefits in reducing muscle soreness and speeding up recov-1665 ery." Getting the behavior prediction correct for this story would mean, for instance, models predict 1666 the athlete would likely "purchase the sports recovery cream" (because the athlete would likely not 1667 know about the banned substance to "avoid the cream to prevent failing a drug test"). The **better** 1668 **performance** in such scenarios could potentially be due to safety training of recent LLMs, making 1669 models more alert when dealing with situations that involve sensitive topics like health and drugs. 1670 However, even then, models would still do poorly for the corresponding judgment questions, judg-1671 ing that "purchase the sports recovery cream", the likely action they had previously chosen, is "not 1672 reasonable" behavior. The observation that better performance on one type of applied ToM ques-

1673

1650 1651

¹⁰Figures 11 and 12 provide results across all broad scenarios for different models.

tions (behavior questions) does not translate to better performance on another (judgment questions)
 tions (behavior questions) does not translate to better performance on another (judgment questions)
 turther emphasize the need for different kinds of applied ToM questions, as present in Simple ToM, beyond the commonly used questions in existing neural ToM tests (focusing on explicit ToM and sometimes just action questions for applied ToM).

More on where failure occurs. Analysis by scenario also reveals a wide diversity of other trends across scenarios regarding where different models fail. We present further results in Figure 13. For instance for "inside containers for personal belongings" situations, failure for GPT-40 is most frequent in the behavior prediction part (see pink portion dominating in pie chart) of the inference chain whereas it makes up less than half of the pie chart for other models. This suggest that behavior prediction in such situations could be an area of weakness to look into when attempting to develop future iterations of the GPT-40 model.

1685 **Perfection is possible but many LLMs are not there.** In fact for the two categories "inside reuse la-1686 beled containers" and "provider info healthcare", in comparison to the other three models, a stronger 1687 and later model like o1-preview achieves close to perfect performance across the three question types 1688 testing ToM reasoning. This further illustrates the high-quality nature of SimpleToM, in that these 1689 simple two-sentence stories are clean, straightforward tests of neural ToM reasoning, yet models other than the strongest and latest model o1-preview shows poor performance on applied ToM questions (behavior and judgment) in various ways. Model developers, if interested in real-world deployment of their models, should be alert into closing this performance gap so as to ensure their models can interact with society appropriately, ideally without the high inference costs of chain-of-thought 1693 reasoning and o1-preview reasoning tokens (that will be discussed in Section 7 and Appendix I). 1694

1695 1696

1697

1701

1702

1703

1704

1705

1706

1707

1708

1709

M PERFORMANCE ACROSS PERSONAS

Another inference-time intervention is to imbue certain personas onto the models. We experiment with two of the models (GPT-40 and Claude-3.5-Sonnet), using the following five personality descriptions from PersonaHub:¹¹

- **lawyer:** You are a partner at the law firm, recognized for their extensive knowledge of healthcare laws.
 - worker: You are a factory worker who doesn't trust the COVID-19 vaccine.
 - **atheist:** You are an atheist, philosophy lecturer who encourages open dialog about faith and belief systems.
 - **psychology student:** You are a university psychology student who is currently studying creativity and personality.
 - **psychologist:** You are a clinical psychologist collaborating with the music therapist to provide holistic patient care.

1710 We inserted these into the system prompt to produce the results shown in Table 9. Across personas, 1711 we see that the gap between explicit ToM (mental state prediction) and applied ToM (behavior prediction and judgment) remains consistently prominent, similar to the case where no persona is 1712 specified. This indicates that our finding on the gap between explicit and applied ToM in LLMs is 1713 robust to injecting different personas. There are, however, minor differences across personas that 1714 may open up interesting directions for future studies. For instance, applying the "worker" persona 1715 with GPT-40 results in slightly worse performance than other personas on the mental state questions 1716 (though minor, < 3%) but slightly better performance on the behavior prediction questions. It is 1717 also consistent across GPT-40 and Claude 3.5 (and more prominent in the case of Claude 3.5) that 1718 the "lawyer" persona yields somewhat better performance on judgment questions (still way below 1719 random), potentially an effect of the model trying to mimic careful judgment when operating under 1720 that persona.

- 1721
- 1722 1723
- 1704
- 1725
- 1726

¹⁷²⁷

¹¹https://huggingface.co/datasets/proj-persona/PersonaHub

Scenario	Description and example
food item in grocery store	General description: When person X buys a food item Z
	grocery store and something hidden is wrong with the food,
	person X will not know about it before paying for it.
	Specific example : a carton of milk that has gone bad because
	improper storage.
provider into healthcare	General description: When a provider know that healthcare p
	uct Z has important limitations that should deter a consum
	from using it, they can still try to sell Z to consumer X in int
	and not disclosing the limitations
	Specific example : a new drug has several suspected side ef
	that were not reported.
true property pretentious labels	General description : When a seller labels product Z with a
F	tle property that helps them sell product Z for a higher price
	product Z does not have that property, a potential buyer X wi
	have enough information to know that.
	Specific example: shop owner puts fancy "organic" labels on
	mal fruits and sells them at a much higher price.
behind the scene service industry	General description: When person/business Z in the servic
	dustry has questionable behind-the-scenes practice, the bus
	can still try to promote their service to consumer X by focusin
	promoting the attractive side.
	Specific example: the cher of a restaurant is reusing the work
inside rause labeled containers	Conoral description: When person V brings something
liside reuse labeled containers	(opaque) container 7 labeled with a popular brand person X
	ing the container will infer it is something from the brand an
	know what is inside (such as if it contains something compl
	different).
	Specific example : person Y put yogurt in a KFC paper bag.
unobserved unethical actions	General description: When person Y performs unethical a
	Z, and person X is not in the right place or time to observe Y
	forming Z, person X would not know about Z.
	Specific example : person Y on the 3rd floor of the school b
	ing look out some notes and cheated during the exam. Pers
	cheating
inside containers for personal belong-	General description: When person V brings something i
ings	opaque container Z for personal belongings, person X seeing
$\boldsymbol{\sigma}^{*}$	tainer Z will not know what is inside.
	Specific example: person Y brings a new toy in his school ba
seller info in second hand market	General description: When person Y has an item Z and s
	thing hidden is wrong with the item, then person X, a pote
	buyer of the item Z will not know about it, especially if pers
	tocuses on showcasing what is good about item Z.
	Specific example: a fridge that has problems like it occasio
hiddan hady nort facture	emits a loud sound.
moden body part leature	their body which is generally hidden under their olethes are
	then person X will not know about it
	Specific example: person Y has a sear on their stomach at se
locked devices accounts	General description: When person Y has a locked device of
Tooked devices decounts	count Z, their status or activity in Z are not observed by perso
	Specific example: person X does not have access to person
	utility bill account so they would not know when person Y f
	person in the second

Table 7: Description and examples for broad scenarios where information asymmetry occurs naturally in the real world.

1	7	8	2
1	7	8	3

1	786
1	787
1	788

#	stories	#unique X	#unique Z	#unique Y
grocery store	168	26	38	
labeled containers	164	36	33	26
iners for personal belongings	142	39	37	35
y pretentious labels	139	35	36	
o healthcare	130	34	33	
scene service industry	119	35	33	
n second hand market	99	11	32	20
unethical actions	87	23	30	21
ces accounts	62	26	30	19
/ part feature	37	23	23	19
	1147	255	319	83
y pretentious labels o healthcare scene service industry n second hand market unethical actions ces accounts y part feature	142 139 130 119 99 87 62 37 1147	39 35 34 35 11 23 26 23 255	37 36 33 33 32 30 30 23 319	20 21 19 19 83

Table 8: Statistics for SimpleToM across the different scenarios, including the number of unique entities of each type (Person X, Object/Person/Action Z, Person Y).

1813					
1814	model	persona	mental state (Explicit ToM)	behavior (Applied ToM)	judgment (Applied ToM)
1015	CDT4 -		05 (40.5	15.2
1816	GP14-0		95.0	49.5	15.5
1817	GPT4-o	lawyer	95.5	49.7	17.2
1010	GPT4-o	worker	93.1	55.8	15.9
1010	GPT4-0	atheist	95.0	50.9	15.6
1819	GPT4-0	psychology student	94.4	47.6	15.5
1820	GPT4-o	psychologist	95.2	53.9	16.7
1821	Claude-3.5-Sonnet	19 0	97.9	67.0	24.9
1822	Claude-3.5-Sonnet	lawyer	98.4	67.6	32.0
1000	Claude-3.5-Sonnet	worker	97.9	67.0	24.5
1823	Claude-3.5-Sonnet	atheist	97.9	65.8	23.4
1824	Claude-3 5-Sonnet	psychology student	97.3	68.6	24.8
1825	Claude-3.5-Sonnet	psychologist	97.9	68.6	26.9
1826		1.7. 8			

Table 9: Evaluation results for SimpleToM on the different question types across 5 alternate personas, showing minor differences in scores, but without significantly closing the gap in performance between the explicit ToM mental state questions vs the implicit ToM behavior and judgment questions.



Figure 13: Analyzing where top models fail first in the the sequence of predicting mental state, then behavior and finally judgment (see Figure 3 for legend). We can record failures for the first mistake e.g., whether models (i) fail at the mental state (MS) question, (ii) pass the MS question but fail at behavior prediction, (iii) pass both MS and behavior questions but fail at judgment question, or (iv) get all 3 questions correct. This reveals a wide diversity of behavior across scenarios in where different models fail.