


A Notion of Complexity for Theory of Mind via Discrete World Models

X. Angelo Huang^{1*} Emanuele La Malfa^{2,3}
Samuele Marro¹ Andrea Asperti¹ Anthony G. Cohn^{3,4} Michael Wooldridge^{2,3}
¹DISI, University of Bologna ²Dept. of Computer Science, University of Oxford
³The Alan Turing Institute ⁴University of Leeds
xuanqiang.huang@studio.unibo.it emanuele.lamalfa@cs.ox.ac.uk

Abstract

Theory of Mind (ToM) can be used to assess the capabilities of Large Language Models (LLMs) in complex scenarios where social reasoning is required. While the research community has proposed many ToM benchmarks, their hardness varies greatly, and their complexity is not well defined. This work proposes a framework inspired by cognitive load theory to measure the complexity of ToM tasks. We quantify a problem’s complexity as the number of states necessary to solve it correctly. Our complexity measure also accounts for spurious states of a ToM problem designed to make it apparently harder. We use our method to assess the complexity of five widely adopted ToM benchmarks. On top of this framework, we design a prompting technique that augments the information available to a model with a description of how the environment changes with the agents’ interactions. We name this technique Discrete World Models (DWM) and show how it elicits superior performance on ToM tasks.¹

 <https://github.com/flecart/complexity-tom-dwm>

1 Introduction

Theory of Mind (ToM) studies how agents form and use beliefs to reason in dynamic environments (Premack and Woodruff, 1978). Originally developed to describe human interactions (Preston and De Waal, 2002; Tomasello, 2009) as well as toddlers’ psychological development (Wimmer and Perner, 1983; Baron-Cohen et al., 1985), ToM has been quickly adopted by other fields, including artificial intelligence (McCarthy, 1979; Scassellati, 2002), bayesian inference (Baker et al., 2011) and machine learning (Rabinowitz et al., 2018). In machine learning, ToM has both descriptive and pre-

*First author. Work done while visiting the University of Oxford.

¹Code and data for **full reproducibility** are available in the Code Material.

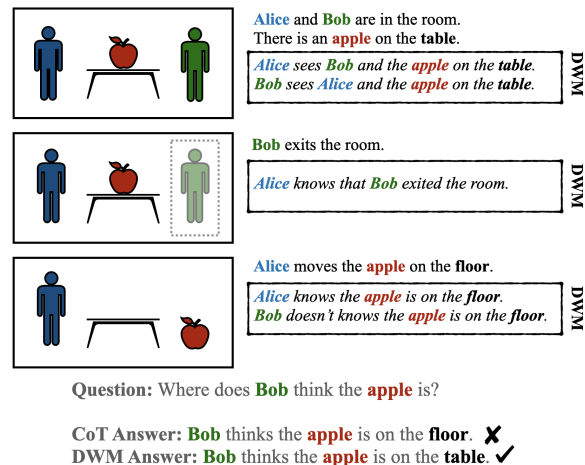


Figure 1: Example of the DWM prompting technique on a classical Sally-Anne QA task (Baron-Cohen et al., 1985). Inspired by our complexity framework (Section 3.1), DWM takes the original task and splits it into sequences, the *state events* (see Def. 3.1), and prompts the LLMs to describe the states. We show that, in most cases, this aids the LLM in providing correct answers.

scriptive usage: on the one hand, ToM benchmarks assess the capabilities of a model in complex environments; on the other, ToM’s frameworks such as *theory-theory* (Gopnik and Wellman, 1994) and *simulation theory* (Churchland, 2013) have been widely adopted to test the proficiency of Large Language Models (LLMs) in social tasks where humans excel (Strachan et al., 2024).

In this work, we propose a framework to characterise a ToM benchmark’s difficulty, i.e., its **complexity**, as the number of *state events* that are sufficient to track the state of an object, including k^{th} -order beliefs motivated by theoretical parallelisms with Sweller’s cognitive load theory (Sweller, 2010).

We characterise the complexity of five standard ToM benchmarks, from false belief to common-sense and social reasoning, and compute their complexity as a proxy of their inherent difficulty. In-

spired by prompting techniques that split a task into elementary sub-problems that are solved sequentially, like Tree of Thoughts (Yao et al., 2024) and least-to-most prompting (Zhou et al., 2023a), we introduce a technique that stimulates a model’s reasoning capabilities via Discrete World Models (DWM). DWM leverages the notion of statefulness via a succinct and coherent representation of each *state events*, as illustrated in Figure 1. We test DWM on ToMi (Le et al., 2019), MindGames (Sileo and Lerneuld, 2023), Adv-CSFB (Shapira et al., 2024), SocialIQA (Sap et al., 2019), and FANToM (Kim et al., 2023), eliciting superior performance than Chain of Thoughts (CoT) (Wei et al., 2022) and Tree of Thoughts (ToT) (Yao et al., 2024) on those problems whose *state spaces* are informative. We further assess whether memorisation affects a model’s performance, and we discover that while this phenomenon happens for standard benchmarks such as ToMi (Le et al., 2019), with input-output pairs that can be retrieved *word for word* via prompting, it does not strongly correlate with a drop of performance on memorised ToM benchmarks. We conduct our experiments on a variety of open- and closed-source LLMs, including GPT-3.5-Turbo, GPT-4 (OpenAI, 2023), LLaMA3-70B (AI@Meta, 2024; Dubey et al., 2024) and Mixtral 8x7B (Jiang et al., 2024). In summary, in this paper:

- We introduce the concept of **complexity** of a ToM task to quantify the hardness of keeping track of the elements (e.g., agents’ beliefs or objects’ states) that are sufficient to produce the correct answer to different problems inspired by frameworks in cognitive science.
- We propose DWM, a simple yet effective prompting technique that improves a model’s capability by making **implicit** information explicit while not necessitating *exogenous information* (i.e., it does not require RAG or fine-tuning).

We consider our work a step towards a framework that formalizes the hardness of a ToM problem in an unambiguous way, inspired by the theory of World Models (Wong et al., 2023).

2 Related Work

Over 40 years of research on ToM in psychology (Premack and Woodruff, 1978; Baron-Cohen

et al., 1985; Dennett, 1988; Wellman, 2017) on human development has created a fertile ground for the development of these ideas in adjacent fields. In the last decade, many works studied ToM in artificial intelligence and machine learning (Baker et al., 2011; Rabinowitz et al., 2018), with applications to multi-agent systems and reinforcement learning (Gronauer and Diepold, 2022). More recently, the rise in popularity of LLMs shifted the interest towards understanding and benchmarking large models’ capacity to solve increasingly complex ToM tasks (Aru et al., 2023; Zhou et al., 2023b; Mahowald et al., 2024). While some researchers believe LLMs have already become proficient in solving ToM tasks (Bubeck et al., 2023; Kosinski, 2023; Strachan et al., 2024), others show scepticism and illustrate cases where they fail on trivial variations of well-known problems (Ullman, 2023; Shapira et al., 2024; Sap et al., 2022). In a joint effort between computer scientists and psychologists, many ToM benchmarks have been developed and used to test neural-network models, including LLMs (Gandhi et al., 2022; Chen et al., 2024; Strachan et al., 2024). Recently, concepts such as World Models (Ha and Schmidhuber, 2018) have found applicability mostly as discrete prompting techniques in conjunction with optimisation procedures (Hao et al., 2023; Moghaddam and Honey, 2023). Researchers have found evidence of an emergent internal representation (e.g., World Model’s surrogates) of the state games (Li et al., 2023; Toshniwal et al., 2021) and state-tracking abilities (Li et al., 2021; Kim and Schuster, 2023; Kim et al., 2024), necessary for correct belief tracking in ToM problems. Cognitive load theory emerged in the late eighties with Sweller’s work on human problem solving (Sweller, 1988). Most measures of cognitive load are based on subjective reports from humans (Sweller et al., 2011). Even though some attempts at automatic cognitive load measures have been present (Yin et al., 2007), they have not been widely adopted in the community. The works that are more similar to our complexity framework are only tangentially related to ToM. Inspired by the work in (Zhou et al., 2023a) and the results in (Zhou et al., 2023b), our prompting technique is inspired by (Park et al., 2023) and (Nye et al., 2021): the former develops an architecture to record the agent’s experiences. The latter proposes a prompting technique that forces a model to express the intermediate computational steps to solve a problem.

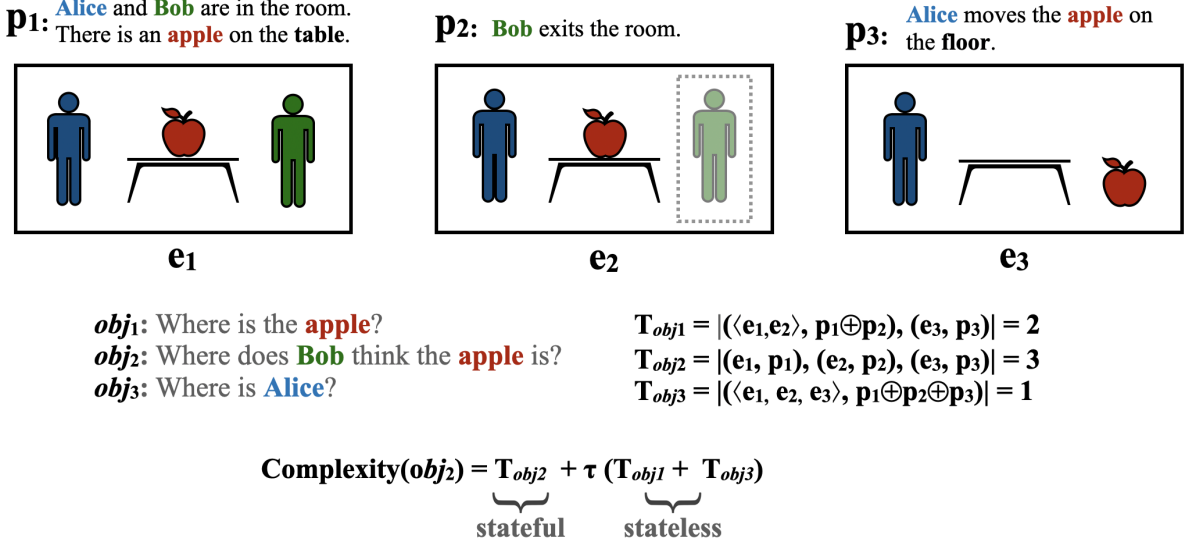


Figure 2: How statefulness and statelessness (Def. 3.2) are computed for the motivating example in Fig. 1. For obj_1 , an optimal split to track the **apple** merges the first two states and chunks of the input prompt. For obj_2 , which involves the 1st-order belief of **Bob**, the statefulness is higher, with e_2 that cannot be merged with e_3 as it introduces partial observability. The complexity of the task (bottom) is computed as per Eq. 2, where the complexity of objects that are not directly relevant to the question/answer is discounted.

3 Methodology

In this section, we introduce a notion of complexity for ToM problems which quantifies the hardness of a problem as the number of *computational steps* humans take to solve them and compare it with Sweller’s cognitive load theory. We then present the DWM prompting technique within the complexity framework and show how it differs from standard methods like CoT and ToT. We further characterise its efficiency with the number of input/output tokens and queries to a model as the control variables.

3.1 On the Complexity of ToM

The need to provide a consistent representation of the environment, including each agent’s beliefs, inspired us to characterise the complexity of a ToM problem in terms of **sufficient elements to track** to output the correct result. Consider a problem prompt p , expressed in natural language, that describes how multiple agents interact with an environment object obj , as illustrated in Figure 2 (top). In our framework, an object can be the state of the apple as well as the k^{th} -order belief of an agent about the apple position. Our framework naturally extends to multiple objects by considering their union.

Suppose that in p , the state of obj is modified $T > 0$ times, thus identifying T unique configura-

tions, namely $E_{obj} = \{e_1, \dots, e_T\}$. To correctly solve a ToM task where p is complemented by a query about obj , a model should distinguish between the interactions that modify the configuration of obj , i.e., the **stateful** states, from those that modify any other **stateless** object $Obj \setminus obj$, i.e., those that one does not need to track.

We first show how to define the cost of tracking a task’s stateful states, which we complement with that of the stateless. Both definitions concur in defining the **complexity** of a ToM task.

3.1.1 Stateful and Stateless Complexity

For a ToM task, expressed as p , that describes the evolution of an environment where an unknown number of atomic iterations T modifies obj or its perception, each environment state $e_t \in E_{obj}$ can be coupled with the prompt prefix $p_{\leq t}$ s.t. $p_{\leq t} \oplus p_{> t} = p$, that describes such configuration. We denote $(e_t, p_{\leq t})$ as a generic *state description*, as illustrated in Figure 2 (top).

Definition 3.1 (State event). A *state event* for an object obj is an event that links adjacent *state descriptions* that involve, for both the environment state e_t and the sub-prompt $p_{\leq t}$, a state change of obj . Formally, we define a relation, F_{obj} , to specify which pairs of state descriptions form a state event: $F_{obj}((e_t, p_{\leq t}), (e_{t+1}, p_{\leq t+1})) \equiv e_t \neq e_{t+1} \wedge p_{\leq t+1} = p_{\leq t} \oplus p_{t+1}$ where $1 \leq t \leq |p|$.

($|p|$ denotes the number of atomic prompts.) and \oplus is the string concatenation operator.

Thus a *state event* F_{obj} identifies those *state descriptions* $(e_t, p_{\leq t})$ which have a successor $(e_{t+1}, p_{\leq t+1})$ where obj has changed its configuration.

In the context of ToM tasks, a *state event* could be a person who moves an object, exits (thus introducing partial observability) or witnesses a change in the environment (as now the description of the environment will take that change into account), as illustrated Figure 2 (middle). Our prompting technique, namely DWM (Section 3.2.1), aims at making implicit observations about objects explicit.

We finally introduce the notion of *partition function* to connect the **maximum number** of non-empty *state events* relative to a prompt. Such a notion will serve as the building block to compute the complexity of a ToM problem.

Definition 3.2 (Partitions). A *partition* $part_{\text{obj}}$ w.r.t. obj identifies those *state events* which partition a ToM prompt p into sequential segments where obj changes its value. Formally:

$$\text{Let } part_{\text{obj}} = \{(e_t, p_{\leq t}) : F_{\text{obj}}((e_t, p_{\leq t}), (e_{t+1}, p_{\leq t+1})) \wedge e_t \in E_{\text{obj}}\} \quad (1)$$

Def. 3.2 describes an optimal partition, $part_{\text{obj}}$ of *state descriptions* that covers all the relevant changes to obj . The partition is represented by the set of event descriptions where obj changes its description immediately after. Note that this set of event descriptions is unique for any obj .

3.1.2 The Complexity of a ToM Task

We can now define the notion of **statefulness** of a ToM task specified as a prompt p as the size of Eq. 3.2, namely $T_{\text{obj}} = |part_{\text{obj}}|$. The process of computing the statefulness of an object or its belief is illustrated in Fig. 2.

For a ToM task where the question to solve relates to an object obj , one must ensure that changes to any other object, namely $Obj \setminus \text{obj}$, do not affect obj . While tracking the evolution of what is irrelevant to answer the question is unnecessary, a computation model must assess whether a particular environmental change affected obj . We thus introduce the notion of **statelessness**, i.e., the cost of discerning whether a change in the environment affects obj . The computation is similar to that of Def. 3.2, except that obj is replaced

with any object in $Obj \setminus \text{obj}$; however for stateless objects, we introduce a discount factor τ to penalise the complexity of *state events* that do not affect obj . Mathematically, we formalise the statelessness of a ToM task involving an object obj as $\tau \sum_{obj \in Obj \setminus \text{obj}} T_{obj}$.

Finally, we formalise the complexity of a ToM task w.r.t. an object obj as the complexity of the stateful states plus the (discounted) sum of the others (i.e., stateless). Namely:

$$T_{\text{obj}} + \tau \sum_{obj \in Obj \setminus \text{obj}} T_{obj} \quad (2)$$

The process of computing the complexity of a ToM task is illustrated in Figure 2.

3.1.3 A parallelism from the Cognitive sciences

Understanding how humans solve complex problems has long served as a valuable source of inspiration for advancing machine intelligence. Thought frameworks in the cognitive sciences, such as Kahneman’s Dual-process theory (Kahneman, 2011), have greatly influenced various fields, including artificial intelligence. In this work, we draw ideas from another theory, less known in the community, the cognitive load theory (Sweller, 1994, 1988, 2010). According to this theory, three main factors influence the mental effort humans exert when solving a particular task or learning new information: intrinsic, extraneous and germane load. The **intrinsic load** measures the natural difficulty of a certain task, the information that *needs* to be digested before answering the question. This relates to the complexity of the material itself. The **extraneous load** concerns the presentation of the information in the problem. For example, if questions phrased in a complex manner are used with a child, it would be much more difficult to understand and answer correctly compared to easier phrasing. Similarly, if many confounding sentences or sentences that do not matter in answering a question are present in the text, we expect an LLM to be worse, suggesting a weak similarity between the means of reason of these models and humans. Finally, **germane load** estimates the working memory resources needed to understand the important parts of the problem, i.e. the intrinsic load. If part of the memory is devoted to the extraneous load, then the germane load is diminished, suggesting a positive correlation with intrinsic load and a negative correlation with extraneous load.

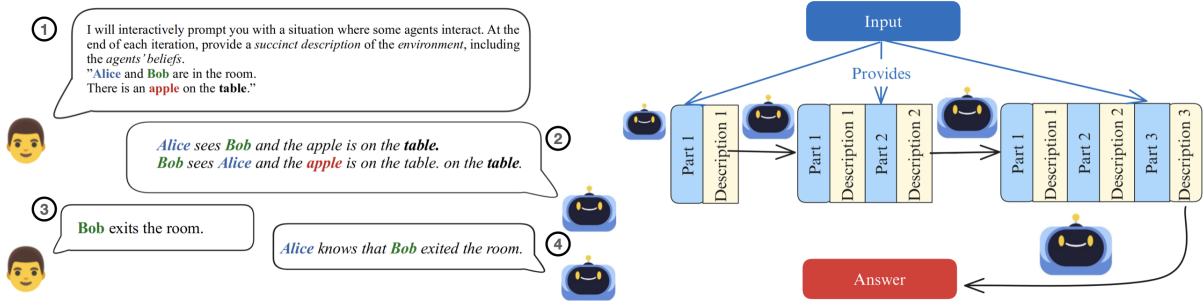


Figure 3: Left: illustration of DWM prompting as per the example in Figure 1. We interactively prompt an LLM with a ToM problem, asking to provide a succinct representation of each agent’s beliefs. Right: schematic presentation of the DWM method. We first break the input string into T state descriptions. Then, for each part, we ask the LLM to provide the **state event** of the environment and how it changes. In the last step, every part of the input and description is fed to the LLM with another prompt to get the answer for the task.

3.1.4 A comparison with the Cognitive Load Theory

Our framework, summarized in figure 2, has two main parts: stateful and stateless complexity. These notions have some similarities with, respectively, the *intrinsic load* and *extraneous load*. Stateful complexity provides a measure on the sentences that are needed to answer the question correctly and must be adequately represented in memory. In a similar manner, intrinsic load concerns on the needed information to correctly analyze a task. Likewise stateless complexity yields information about the confounding or irrelevant sentences and phrases in the text akin to extraneous load. In our setting, *germane load* could be interpreted as the ratio of the stateful and stateless complexity: higher ratio means higher density of useful sentences in answering a question. This notion of load could be used as a basis of an objective measure on the quality of a question-answering sample: given the same quantity of cognitive load, i.e. complexity, we would like to have a simple presentation with correct information, maximizing the germane load. If the cognitive load hypothesis applies to LLMs, maximizing the germane load would lower the complexity of the tasks given to a model, and therefore it would aid the model to answer questions more accurately.

3.2 Discrete World Models

We first introduce the background notation for prompting LLMs and assessing their accuracy on a standard classification task. We then propose our technique, namely DWM, which we eventually connect with the notion of statefulness of a ToM

task.

Background notation. A (Large) Language Model is a function that predicts the next token (out of a finite vocabulary) conditioned on the sequence of previously fed/generated tokens, namely $\psi : \mathbf{v} \in V^* \rightarrow v \in V$. Such a mechanism can be used to sample multiple token outputs until an ‘end-of-text’ token is predicted, by invoking ψ in an auto-regressive fashion, i.e., $\psi(v|\mathbf{v})$. In our setting, a problem is specified as a tuple (p, Q) , where p is a ToM problem and Q is a *query* function that modifies p according to a prompting technique, namely $Q : p \rightarrow p'$. The LLM’s output y for an input $Q(p)$ is then compared for correctness against an oracle Ω , i.e., $\Omega : \psi(Q(p)) \rightarrow \{0, 1\}$, where 1 means correct classification (0, otherwise). On a sample of $N > 0$ ToM problems, the accuracy of a model ψ is then measured as $\frac{1}{N} \sum_{i=1}^N \Omega(\psi(Q(p_i)))$, i.e., the average number of times a model is correct in its prediction.

3.2.1 Discrete World Models via Prompting

Given a ToM problem p and a constant $T \leq |p|$, where $|p|$ is ideally measured as the number of state changes in the problem, we can rewrite p as $p_1 \oplus p_2 \oplus \dots \oplus p_T$. Our *query* function adds a standard preamble x similar to that of CoT. DWM inserts, after each “split” p_t , an additional prompt w like ‘Now, provide a succinct description of the state of the environment and each agent’s belief.’ and query an LLM to provide a representation of the current *state description* of the environment. An LLM is initially queried with $x \oplus p_1 \oplus w$, and the answer a_1 is concatenated to the next query, i.e., $\psi(x \oplus p_1 \oplus w \oplus a_1 \oplus p_2 \oplus w)$ to

retrieve a_2 . The process is carried on for each of the T chunks, and, at the end, y is concatenated to eventually prompt the model for the correct answer to p .

Let $a_1 = \psi(x \oplus p_1 \oplus w)$, $a_t = \psi(x \oplus p_1 \oplus w \oplus a_1 \oplus p_2 \oplus \dots \oplus a_{t-1} \oplus p_t) = \psi(x \oplus (\bigoplus_{i=1}^{t-1} p_i \oplus w \oplus a_i) \oplus p_t)$, then, the final query is

$$\psi(x \oplus \left(\bigoplus_{t=1}^T p_t \oplus w \oplus a_t \right) \oplus y) \quad (3)$$

In this sense, our **partition function** (Def. 3.2) consists of splitting a prompt into sequential chunks of the prompt, while the LLM is prompted to provide each *state event* at time $1 \leq t < T$ as $e_t = \psi(x \oplus (\bigoplus_{i=1}^t p_i \oplus w \oplus a_i) \oplus \omega)$. The process of prompting a model with DWM is illustrated in Figure 3.

3.2.2 On the Complexity of DWM

DWM progressively calls an LLM $T > 0$ times to generate informative states. For a ToM problem of length n (i.e., the number of input tokens), which we assume, w.l.o.g., that can be split into k chunks of approximately the same length $|x \oplus p_i \oplus w| = \frac{n}{T}$, the number of tokens generated by an LLM is of the order of $\mathcal{O}(\sum_{t=1}^T |x \oplus (\bigoplus_{i=1}^{t-1} p_i \oplus w \oplus a_i)|)$, where p_t (a_t) is the portion of the problem (answer) prompted (retrieved) at iteration t . With the further assumption that each answer retrieved at split $t \leq T$ has the same length o , the complexity is further simplified to be asymptotic to $\mathcal{O}((\frac{n}{T} + o)^2)$. Compared to CoT, whose complexity is $\mathcal{O}(n + o)$, DWM requires an additional linear number of calls to the model. On the other hand, ToT with the same number of splits $\frac{n}{T}$ and $m > 1$ experts results in even higher complexity, i.e., asymptotic to $\mathcal{O}(m(\frac{n}{T} + o)^2)$.

4 Experimental Evaluation

The experiments are organised as follows. We first test the performance of DWM on ToMi (Le et al., 2019), MindGames (Sileo and Lernould, 2023), Adv-CSFB (Shapira et al., 2024), SocialIQA (Sap et al., 2019), and FANToM (Kim et al., 2023), comparing it with CoT (Wei et al., 2022), ToT (Yao et al., 2024) and prompting with structured data (struct), i.e., the model is queried to first represent the problem in a structured format such as JSON or Yaml. We further show that ToMi has been memorised *word for word* by GPT models, with

CoT (and any technique that leaves the input unchanged) being the best-performing method. We then quantify the complexity of the benchmarks introduced above and highlight the correlation with the models' performances. Our framework shows complexity ranges between easy and hard problems, even within a benchmark. We conduct our experiments on GPT-3.5-Turbo, GPT-4 (OpenAI, 2023), LLaMA3-70B (AI@Meta, 2024; Dubey et al., 2024) and Mixtral 8x7B (Jiang et al., 2024).

4.1 DWM on ToM Benchmarks

We report results for GPT-3.5-Turbo and Mixtral 8x7B on the five ToM benchmarks: for reasons of space, results for LLaMA3-8B, LLaMA3-70B and GPT-4 are reported in the Appendix, Section B.1. As illustrated in Figure 4 (top), DWM improves the performance of GPT-3.5-Turbo on Mindgames, FANToM and Adv-CSFB by a solid margin. On SocialIQA, which has very short inputs, DWM performs slightly worse than CoT but better than ToT. On the other hand, on ToMi, the best prompting techniques are CoT and ToT. While one may think memorisation plays a role in boosting the performance of LLMs with these prompting techniques, in the next section, we provide evidence this hypothesis is not necessarily true. With Mixtral 8x7B (Fig. 4 (bottom)), DWM improves the performance on ADVcsfb, FANToM, ToMi and Mindgames, and reaches that of CoT on SocialIQA.

DWM elicits more informed state spaces. We qualitatively analysed the information elicited by an LLM when prompted with DWM and discovered that it forces a model to output information **not explicitly** available in the prompt. Consider the ToMi example in Figure 5 where GPT-4 is prompted with a situation where agents interact and are then queried with the first-order belief of Benjamin. With CoT, the model makes an erroneous assumption about the presence of Benjamin and Isabella in the room. On the other hand, when prompted with DWM, GPT-4 provides an informative description of each *state space*, particularly the knowledge and the uncertainty of each agent's beliefs, and eventually answers correctly. One example per benchmark is available in the Appendix, Section B.2, while many more are available for inspection in the Code Supplementary Material. This phenomenon is ubiquitous to all the ToM tasks we tested, a hint that DWM may elicit the ToM capabilities of LLMs without requiring external

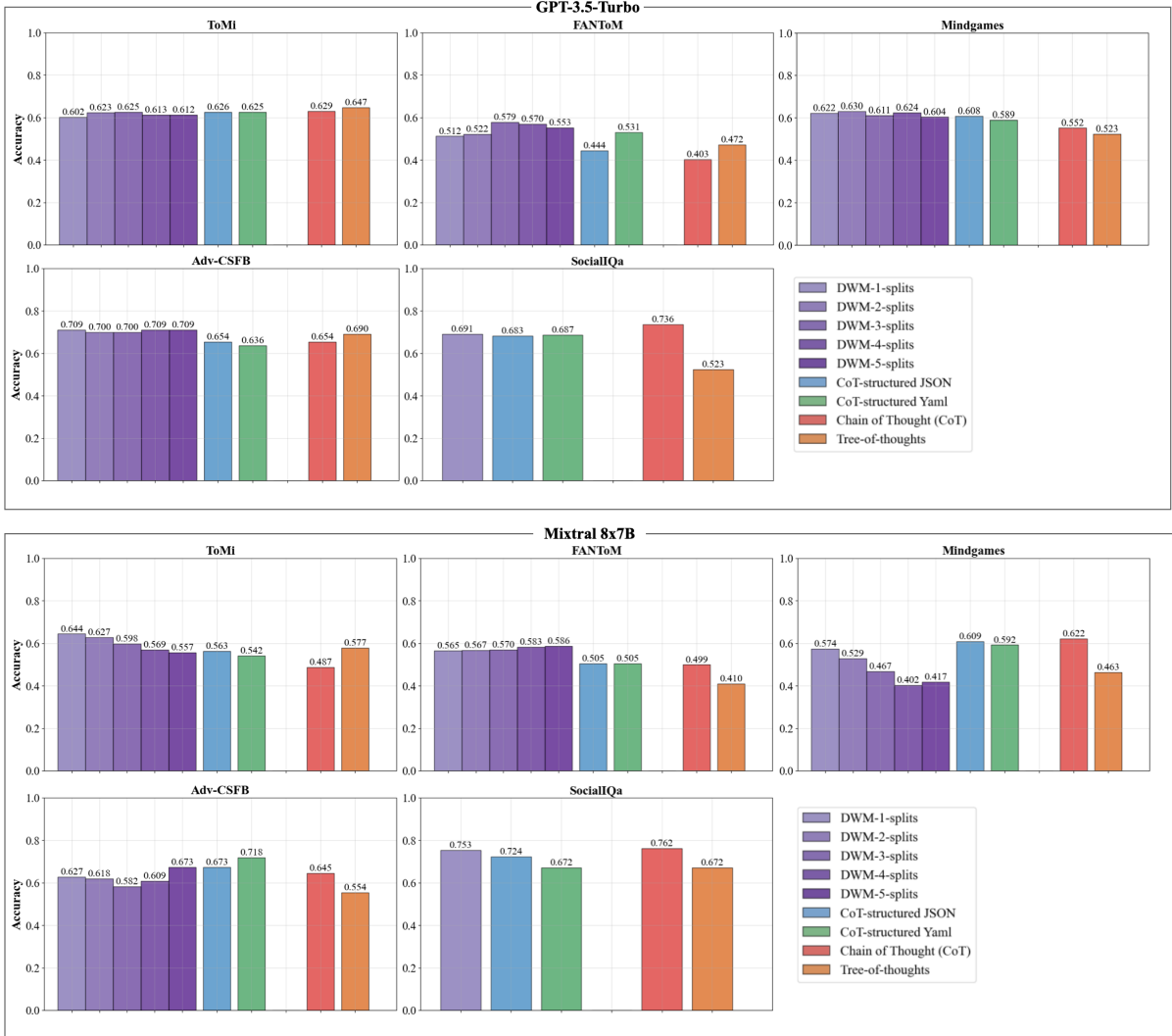


Figure 4: Benchmarks of GPT-3.5-Turbo (top) and Mixtral 8x7B (bottom) models on different ToM tasks for DWM (one to five splits), CoT, ToT and structured prompts (JSON and Yaml).

	ToMi	FANToM	Mindgames	Adv-CSFB	SocialIqa
Memorisation - perfect match	52%	35%	2%	0%	0%
Memorisation - fuzzy	89 ± 15%	74 ± 24%	64 ± 18%	51 ± 11%	40 ± 12%
DWM	0.625	0.579	0.618	0.8364	0.691
CoT	0.629	0.403	0.552	0.7091	0.736

Table 1: Summary of the memorisation test on five ToM benchmarks. We prompted GPT-3.5-Instruct to predict the continuation of 100 randomly sampled test points. We computed the exact and fuzzy memorisation rate (second row, similarity score computed via the Levenshtein distance, see the `thefuzz` package), which we complement with the best performance across models of CoT and DWM.

	ToMi	FANToM	Mindgames	Adv-CSFB	SocialIqa
Statefulness	2.62 ± 1.68	2.44 ± 0.96	1.22 ± 0.90	3.24 ± 1.35	1. ± 0.
Statelessness	4.27 ± 2.1	59.42 ± 18.91	5.24 ± 2.71	2.86 ± 1.34	1.14 ± 0.447
DWM - Best Split	3	3	1	4	1

Table 2: Summary of the statefulness and statelessness of different ToM benchmarks. At the bottom, the value of the split that guarantees max performance of GPT-3.5-Turbo with DWM, which we notice is strongly correlated with the statefulness of each benchmark.

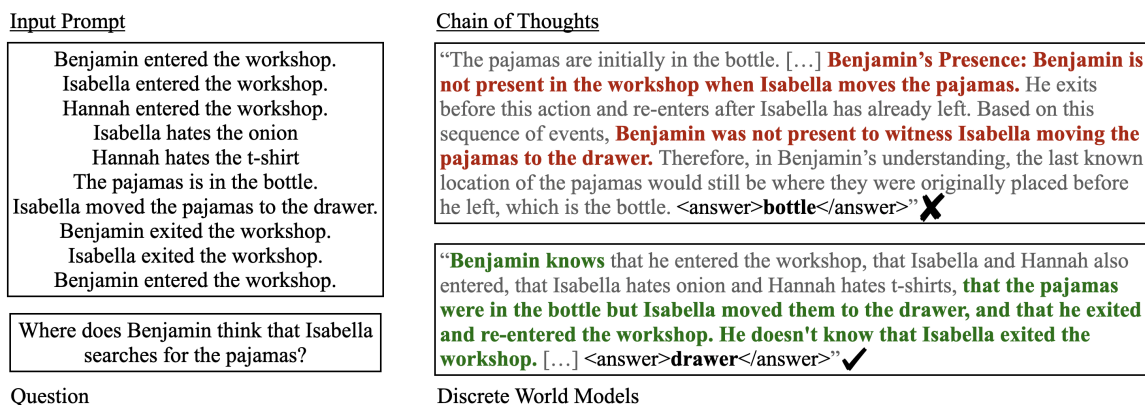


Figure 5: Example of a real ToMi example where GPT-4 fails when prompted with CoT, yet succeeds with DWM. CoT elicits an untruthful reasoning process (in red), while DWM correctly informs the model with the implicit information about Benjamin’s first-order belief (in green). More examples are reported in the Appendix, Section B.2.

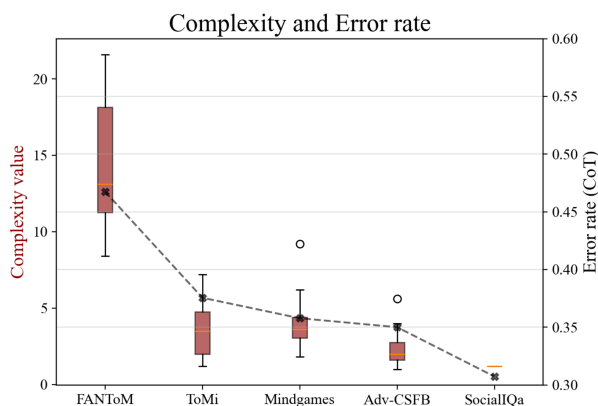


Figure 6: Each boxplot summarizes the complexity analysis of the five ToM benchmarks in ascending order. We report the average error rate (i.e., 1-accuracy) of GPT-3.5-Turbo, GPT-4, Mixtral 8x7B and LLaMA3-70B on the task when prompted with CoT.

information or solvers.

Memorisation in Theory of Mind. Recent works expressed concern about ToM benchmarks’ efficacy in memorisation (Jacovi et al., 2023; La Malfa et al., 2024b). This motivated us to quantify and then analyse the impact of memorisation of ToM benchmarks on performance. We computed the percentage of memorised prompts to understand whether that affects the performance of techniques, such as DWM, that split the prompt into chunks and introduce additional information instead of CoT, which leaves the input prompt unchanged. As illustrated in Table 1, ToMi and FANToM have been heavily memorised, with entire portions of the benchmarks that can be retrieved word for word from GPT-3.5-Instruct (the auto-complete model by OpenAI). Despite that, no clear

evidence of a performance drop in DWM induced by memorisation exists. For GPT-3.5, despite CoT having higher performance on ToMi, DWM is better on FANToM (Figure 4). We hypothesise that as long as a memorised problem is prompted, either in its exact form (as for CoT) or split as in DWM, the most potent models can recover it alongside the ground truth label, thus invalidating the test for both. We conclude with a note of caution: while we discovered that ToMi and FANToM are memorised by GPT-3.5-Instruct, that does not imply any LLM, including GPT-3.5-Turbo and GPT-4, whose training details are not released publicly, has been trained on that data.

4.2 Statefulness of ToM Benchmarks

We used the complexity framework introduced in Section 3.1 to characterise the statefulness and statelessness of the five ToM benchmarks used for the experimental evaluation. We randomly sampled 50 problems from each dataset, identified the objects, and manually labelled stateful and stateless state events. We release the split samples alongside a web application that facilitates manual labelling. As illustrated in Figure 7 (left), the statefulness of each problem, i.e., that of the object a model must track to answer correctly, strongly correlates with the best-performing DWM split. The statelessness complexity, reported in Figure 7 (middle), i.e., that of objects that a model does not need to track, grows larger for problems such as FANToM, only partially influencing the models’ performance. We hypothesise that the most potent models developed some competency in discerning the relevant part of a prompt (the stateful events) from the confounding ones. We finally report, in Figure 7 (right), the

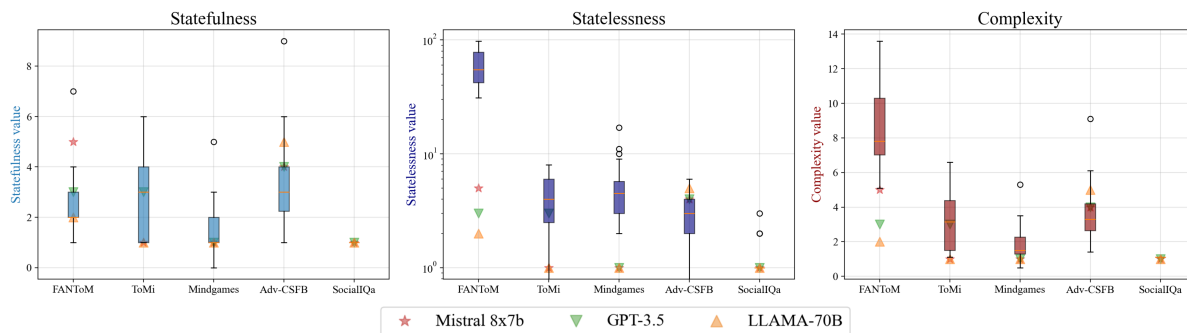


Figure 7: Each boxplot summarizes the statefulness (left), statelessness (middle, y-axis in log-scale) and complexity analysis (right) of the five ToM benchmarks. We report mean, standard deviation and outliers alongside the best DWM method (by the number of prompt splits) and observe a strong correlation between the number of splits and the statefulness.

complexity of each problem computed as per Eq. 2, with τ set in a range between 0.05 and 0.2 (i.e., the relative weight of stateless compared to stateful events). Results suggest that FANToM is the most difficult ToM task for humans and LLMs (see Figure 4), followed by ToMi (the second most difficult for LLMs as well) and Adv-CSFB (which seems easier than the others); in contrast, Mindgames and SocialIQa tend to be easier. Finally, in Figure 6, we compare the accuracy of GPT-3.5-Turbo, GPT-4, Mixtral 8x7B and LLAMA3-70B when prompted with CoT (i.e., without split) on the five ToM benchmarks with the complexity of the task as per Def. 2. We observe a **strong correlation** between the error-rate and the complexity of a task, i.e., our framework correctly identifies the tasks that are harder both for humans and current state-of-the-art LLMs.

5 Conclusions

This paper introduces a complexity framework to measure the difficulty of Theory of Mind (ToM) problems. It quantifies the difficulty by tracking necessary states (stateful) and unnecessary states (stateless), with the latter discounted in the complexity computation. The framework evidences a strong correlation between complexity and model performance. Inspired by this framework, we propose DWM, a prompting technique that splits a prompt into parts to query a model for a consistent representation of the environment and agents’ beliefs. DWM outperforms CoT and ToT by extracting implicit but relevant information.

Limitations

Higher order belief tracking. Our theoretical framework reduces the problem of solving a belief

ToM problem to finding the correct descriptions that need to be tracked. It extends seamlessly to tasks with much higher complexity, however, we have not had the opportunity to test this theory in those settings. We noticed that most theory of mind tasks available in the community only require one to five states to be correctly answered. A possible extension would be testing the theory upon tasks with higher state complexity, e.g. k^{th} -order belief tracking tasks. However, it is unclear whether this could be useful in real applications as most human belief tracking is limited to 5 or 6 orders (Cargile, 1970; Dennett, 1988).

On task splitting methods. It is not straightforward to automatically find the correct task splits in a manner that correctly describes the state. An LLM could find a way to split it by itself correctly and use those splits to answer the question. We attempted this approach, yet with a simple prompting method, the model splits every sentence, making the descriptions much noisier and less accurate. Future work could try to find the best splits automatically.

Memorization analysis. Training and evaluating on the same dataset produces positively biased data on the model’s performance. While running our benchmarks on ToMi, we discovered that the GPT-3.5 model had completely memorized parts of the dataset. This motivated us to extend the memorization test to the other tasks. We urge the research community to include a memorization section on every benchmark study with public datasets used in their works. This data is crucial to conduct fair and unbiased research on evaluating LLMs’ abilities (Jacovi et al., 2023). Future works will include an analysis of the memorisation rate of other ToM

tasks alongside tests to quantify their impact on different models.

On element interactivity. Sweller (Sweller, 2010) proposes a measure of complexity for cognitive tasks that encompasses three main components, namely the *intrinsic*, *extraneous*, and *germane* cognitive load. In its framework, which has wide applications in education, the *intrinsic load* relates to the number of references, or interactions, between the elements of a problem, i.e., the information or concept that needs to be understood to answer the question. Our framework approximates the *intrinsic* and *extraneous* loads to be single sentences in a ToM problem, which is not assured to be the best measure.

Ethical Statement

The datasets and pre-trained LLMs that we use are all publicly available. This paper focuses on ToM problems' hardness and prompting methods. We highlight that LLMs do not guarantee the production of factual data or correct reasoning steps, and the prompting methods developed here should not be regarded as *the* source of truth in decision-making.

References

- AI@Meta. 2024. [Llama 3 model card](#). Technical report, Meta.
- Jaan Aru, Aqeel Labash, Oriol Corcoll, and Raul Vicente. 2023. Mind the gap: Challenges of deep learning approaches to theory of mind. *Artificial Intelligence Review*, 56(9):9141–9156.
- Chris Baker, Rebecca Saxe, and Joshua Tenenbaum. 2011. [Bayesian Theory of Mind: Modeling Joint Belief-Desire Attribution](#). *Proceedings of the Annual Meeting of the Cognitive Science Society*.
- Simon Baron-Cohen, Alan M. Leslie, and Uta Frith. 1985. [Does the autistic child have a “theory of mind”?](#) *Cognition*, 21(1):37–46.
- 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. 2023. [Sparks of Artificial General Intelligence: Early experiments with GPT-4](#). *arXiv:2303.12712*.
- James Cargile. 1970. [A note on "iterated knowings"](#). *Analysis*, 30(5):151–155.
- Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen, Guanqun Bi, Gongyao Jiang, Yaru Cao, Mengting Hu, Yunghwei Lai, Zexuan Xiong, et al. 2024. [Tombench: Benchmarking theory of mind in large language models](#). *arXiv preprint arXiv:2402.15052*.
- Paul M Churchland. 2013. Folk psychology and the explanation of human behavior 1. In *Folk psychology and the philosophy of mind*, pages 247–262. Psychology Press.
- Daniel C. Dennett. 1988. The intentional stance in theory and practice. In *Machiavellian Intelligence: Social Expertise and the Evolution of Intellect in Monkeys, Apes, and Humans*, pages 180–202. Clarendon Press/Oxford University Press.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. [The llama 3 herd of models](#). *arXiv preprint arXiv:2407.21783*.
- Kanishk Gandhi, Gala Stojnic, Brenden M. Lake, and Moira R. Dillon. 2022. [Baby Intuitions Benchmark \(BIB\): Discerning the goals, preferences, and actions of others](#). *arXiv:2102.11938*.
- Alison Gopnik and Henry M. Wellman. 1994. [The theory theory](#). In Lawrence A. Hirschfeld and Susan A. Gelman, editors, *Mapping the Mind: Domain Specificity in Cognition and Culture*, pages 257–293. Cambridge University Press.
- Sven Gronauer and Klaus Diepold. 2022. Multi-agent deep reinforcement learning: a survey. *Artificial Intelligence Review*, 55(2):895–943.
- David Ha and Jürgen Schmidhuber. 2018. [World models](#). *arXiv preprint arXiv:1803.10122*.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Wang, and Zhiting Hu. 2023. [Reasoning with language model is planning with world model](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8154–8173, Singapore. Association for Computational Linguistics.
- Alon Jacovi, Avi Caciularu, Omer Goldman, and Yoav Goldberg. 2023. [Stop uploading test data in plain text: Practical strategies for mitigating data contamination by evaluation benchmarks](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5075–5084, Singapore. Association for Computational Linguistics.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. [Mixtral of experts](#). *arXiv preprint arXiv:2401.04088*.
- Daniel Kahneman. 2011. *Thinking, Fast and Slow*. Farrar, Straus and Giroux.

- Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Bras, Gunhee Kim, Yejin Choi, and Maarten Sap. 2023. [FANToM: A benchmark for stress-testing machine theory of mind in interactions](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14397–14413, Singapore. Association for Computational Linguistics.
- Najoung Kim and Sebastian Schuster. 2023. [Entity tracking in language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3835–3855, Toronto, Canada. Association for Computational Linguistics.
- Najoung Kim, Sebastian Schuster, and Shubham Toshniwal. 2024. [Code pretraining improves entity tracking abilities of language models](#). *Preprint*, arXiv:2405.21068.
- Michal Kosinski. 2023. Theory of mind may have spontaneously emerged in large language models. *arXiv preprint arXiv:2302.02083*, 4:169.
- Emanuele La Malfa, Aleksandar Petrov, Simon Frieder, Christoph Weinhuber, Ryan Burnell, Raza Nazar, Anthony G. Cohn, Nigel Shadbolt, and Michael Wooldridge. 2024a. Language Models as a Service: Overview of a New Paradigm and its Challenges. *Journal of Artificial Intelligence Research*, 80.
- Emanuele La Malfa, Christoph Weinhuber, Orazio Torre, Fangru Lin, Anthony Cohn, Nigel Shadbolt, and Michael Wooldridge. 2024b. Code simulation challenges for large language models. *arXiv preprint arXiv:2401.09074*.
- Matthew Le, Y-Lan Boureau, and Maximilian Nickel. 2019. [Revisiting the Evaluation of Theory of Mind through Question Answering](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5871–5876. Association for Computational Linguistics.
- Belinda Z. Li, Maxwell Nye, and Jacob Andreas. 2021. [Implicit representations of meaning in neural language models](#). *Preprint*, arXiv:2106.00737.
- Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. [Emergent world representations: Exploring a sequence model trained on a synthetic task](#). In *The Eleventh International Conference on Learning Representations*.
- Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2024. Dissociating language and thought in large language models. *Trends in Cognitive Sciences*.
- John McCarthy. 1979. Ascribing Mental Qualities To Machines. *AI Lab., Stanford University, Technical Report, Memo*, 326.
- Shima Rahimi Moghaddam and Christopher J. Honey. 2023. [Boosting theory-of-mind performance in large language models via prompting](#). *Preprint*, arXiv:2304.11490.
- Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. 2021. Show your work: Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*.
- OpenAI. 2023. [GPT-4 technical report](#). *ArXiv preprint*, abs/2303.08774.
- Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. [Generative agents: Interactive simulacra of human behavior](#). *Preprint*, arXiv:2304.03442.
- David Premack and Guy Woodruff. 1978. Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences*, 1(4):515–526.
- Stephanie D Preston and Frans BM De Waal. 2002. Empathy: Its ultimate and proximate bases. *Behavioral and Brain Sciences*, 25(1):1–20.
- Neil Rabinowitz, Frank Perbet, Francis Song, Chiyuan Zhang, SM Ali Eslami, and Matthew Botvinick. 2018. Machine theory of mind. In *International conference on machine learning*, pages 4218–4227. PMLR.
- Maarten Sap, Ronan LeBras, Daniel Fried, and Yejin Choi. 2022. Neural theory-of-mind? on the limits of social intelligence in large lms. *arXiv preprint arXiv:2210.13312*.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. [Social IQa: Commonsense Reasoning about Social Interactions](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473. Association for Computational Linguistics.
- Brian Scassellati. 2002. [Theory of Mind for a Humanoid Robot](#). *Autonomous Robots*, 12(1):13–24.
- Natalie Shapira, Mosh Levy, Seyed Hossein Alavi, Xuhui Zhou, Yejin Choi, Yoav Goldberg, Maarten Sap, and Vered Shwartz. 2024. [Clever hans or neural theory of mind? stress testing social reasoning in large language models](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2257–2273, St. Julian’s, Malta. Association for Computational Linguistics.
- Damien Sileo and Antoine Lernoould. 2023. [MindGames: Targeting theory of mind in large language models with dynamic epistemic modal logic](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4570–4577, Singapore. Association for Computational Linguistics.

- James WA Strachan, Dalila Albergo, Giulia Borghini, Oriana Pansardi, Eugenio Scaliti, Saurabh Gupta, Krati Saxena, Alessandro Rufo, Stefano Panzeri, Guido Manzi, et al. 2024. Testing theory of mind in large language models and humans. *Nature Human Behaviour*, pages 1–11.
- John Sweller. 1988. [Cognitive load during problem solving: Effects on learning](#). *Cognitive Science*, 12(2):257–285.
- John Sweller. 1994. [Cognitive load theory, learning difficulty, and instructional design](#). *Learning and Instruction*, 4(4):295–312.
- John Sweller. 2010. [Element Interactivity and Intrinsic, Extraneous, and Germane Cognitive Load](#). *Educational Psychology Review*, 22(2):123–138.
- John Sweller, Paul Ayres, and Slava Kalyuga. 2011. [Measuring Cognitive Load](#). In John Sweller, Paul Ayres, and Slava Kalyuga, editors, *Cognitive Load Theory*, pages 71–85. Springer.
- Michael Tomasello. 2009. *The cultural origins of human cognition*. Harvard university press.
- Shubham Toshniwal, Sam Wiseman, Karen Livescu, and Kevin Gimpel. 2021. [Learning chess blindfolded: Evaluating language models on state tracking](#). *CoRR*, abs/2102.13249.
- Tomer Ullman. 2023. Large language models fail on trivial alterations to theory-of-mind tasks. *arXiv preprint arXiv:2302.08399*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Henry M. Wellman. 2017. [The Development of Theory of Mind: Historical Reflections](#). *Child Development Perspectives*, 11(3):207–214.
- Heinz Wimmer and Josef Perner. 1983. [Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception](#). *Cognition*, 13(1):103–128.
- Lionel Wong, Gabriel Grand, Alexander K Lew, Noah D Goodman, Vikash K Mansinghka, Jacob Andreas, and Joshua B Tenenbaum. 2023. From word models to world models: Translating from natural language to the probabilistic language of thought. *arXiv preprint arXiv:2306.12672*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.
- Bo Yin, Natalie Ruiz, Fang Chen, and M. Asif Khawaja. 2007. [Automatic cognitive load detection from speech features](#). In *Proceedings of the 19th Australasian Conference on Computer-Human Interaction: Entertaining User Interfaces, OZCHI ’07*, pages 249–255. Association for Computing Machinery.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023a. [Least-to-most prompting enables complex reasoning in large language models](#). In *The Eleventh International Conference on Learning Representations*.
- Pei Zhou, Aman Madaan, Srividya Pranavi Potharaju, Aditya Gupta, Kevin R McKee, Ari Holtzman, Jay Pujara, Xiang Ren, Swaroop Mishra, Aida Nematzadeh, et al. 2023b. How far are large language models from agents with theory-of-mind? *arXiv preprint arXiv:2310.03051*.

A Experimental Setup

A.1 Experimental Details

Most of the language models used in this work follow the Language Models as a Service (LMaaS) paradigm (La Malfa et al., 2024a). This model of service does not allow transparency and hinders reproducibility. Reproducibility is difficult to achieve as common software development frameworks, such as CI/CD pipeline, ease the update of the public service but change the underlying entity. From this, it follows that the model tested by the researcher could change at any time. This is not solvable from the outside. Researchers have no control over the software engineering practices inside a LMaaS, but could set some parameters to offer the highest possible grade of reproducibility. We set the temperature to zero or enable greedy decoding by default (this does not imply determinism even if model weights are not changed).² In prompting methods where the creativity of the response is exploited for better performance, e.g., Tree of Thoughts (Yao et al., 2024), we set the temperature to 0.7, the value proposed in the reference papers.

LMaaS providers. We use [Huggingface](#) for Mixtral 8x7B. Groq Cloud for LLama-3-8B and LLama-3-70B. Microsoft sponsorship for GPT-3.5 and GPT-4 access.

A.2 Prompting Templates

We present the different prompting techniques, taking as an example the following prompt from ToMi and GPT-3.5-Turbo as the reference model:

-
1. Benjamin entered the workshop.
 2. Isabella entered the workshop.
 3. Hannah entered the workshop.
 4. Isabella hates the onion
 5. Hannah hates the t-shirt
 6. The pajamas is in the bottle.
 7. Isabella moved the pajamas to the drawer.
 8. Benjamin exited the workshop.
 9. Isabella exited the workshop.
 10. Benjamin entered the workshop.
-

And the following question:

Where does Benjamin think that Isabella searches for the pajamas?

Chain of Thought

²The main explanation is the <https://github.com/pytorch/pytorch/issues/75240>"non-deterministic cuda cores" another could "be batched inference in sparse MoE models", see <https://152334h.github.io/blog/non-determinism-in-gpt-4/here>

Consider the following dialogue where multiple agents interact. At the end, I will ask you a question to answer.

Here's the dialogue:

1. Benjamin entered the workshop.
2. Isabella entered the workshop.
3. Hannah entered the workshop.
4. Isabella hates the onion
5. Hannah hates the t-shirt
6. The pajamas is in the bottle.
7. Isabella moved the pajamas to the drawer.
8. Benjamin exited the workshop.
9. Isabella exited the workshop.
10. Benjamin entered the workshop.

This is the end of the dialogue. Now, this is a question for you to answer.

Question: Where does Benjamin think that Isabella searches for the pajamas?

Think step by step, answer the question with one word and provide the answer between <answer></answer> tags.

For example, reply with <answer>vase</answer>.

Tree of Thought

We first prompt an LLM to propose different solution paths to solve a task.

Consider the following dialogue where multiple agents interact. At the end, I will ask you a question to answer.

Here's the dialogue:

1. Benjamin entered the workshop.
2. Isabella entered the workshop.
3. Hannah entered the workshop.
4. Isabella hates the onion
5. Hannah hates the t-shirt
6. The pajamas is in the bottle.
7. Isabella moved the pajamas to the drawer.
8. Benjamin exited the workshop.
9. Isabella exited the workshop.
10. Benjamin entered the workshop.

Question: Where does Benjamin think that Isabella searches for the pajamas?

Think step by step and list all possible answers providing a single answer on each line.

We then pick the best idea via a majority vote over three agents simulated by the LLM itself:

Given a dialogue and several observation choices, decide which choice is most promising. Analyze each choice in detail, then conclude in the last line "The best choice is {{s}}", where s the integer id of the choice.

1. Benjamin entered the workshop.
2. Isabella entered the workshop.
3. Hannah entered the workshop.
4. Isabella hates the onion

5. Hannah hates the t-shirt
6. The pajamas is in the bottle.
7. Isabella moved the pajamas to the drawer.
8. Benjamin exited the workshop.
9. Isabella exited the workshop.
10. Benjamin entered the workshop.

Here are some possible observations:

Here we insert the output of the previous prompt.

We eventually ask the model for a final answer.

Given this dialogue and possible observations, answer the question with one word and provide the answer between `<answer></answer>` tags.

1. Benjamin entered the workshop.
2. Isabella entered the workshop.
3. Hannah entered the workshop.
4. Isabella hates the onion
5. Hannah hates the t-shirt
6. The pajamas is in the bottle.
7. Isabella moved the pajamas to the drawer.
8. Benjamin exited the workshop.
9. Isabella exited the workshop.
10. Benjamin entered the workshop.

Question: Where does Benjamin think that Isabella searches for the pajamas?

Here we insert the observations generated by the LLM with the previous prompts.

For example, reply with `<answer>vase</answer>`.

Discrete World Models - 1 Split

I give you a phrase of a dialogue between agents. I will reveal more parts of it later. At the end, I will give you a question you must answer.

For each phrase, you must:

- # 1. Write down a succinct description of what each agent knows about the environment and about the other agents. Keep the description short and do not produce redundant information.
- # 2. Each considerations you make must be preceded by the symbol `#GPT#`.

Here's the dialogue:

1. Benjamin entered the workshop.
2. Isabella entered the workshop.
3. Hannah entered the workshop.
4. Isabella hates the onion
5. Hannah hates the t-shirt
6. The pajamas is in the bottle.
7. Isabella moved the pajamas to the drawer.
8. Benjamin exited the workshop.
9. Isabella exited the workshop.
10. Benjamin entered the workshop.

This is the end of the dialogue. Now, this is a question for you to answer.

Question: Where does Benjamin think that Isabella searches for the pajamas?

Think step by step, answer the question with one word and provide the answer between `<answer></answer>` tags.

For example, reply with `<answer>vase</answer>`.

Discrete World Model - 3 Split

I give you a phrase of a dialogue between agents. I will reveal more parts of it later. At the end, I will give you a question you must answer.

For each phrase, you must:

- # 1. Write down a succinct description of what each agent knows about the environment and about the other agents. Keep the description short and do not produce redundant information.
- # 2. Each considerations you make must be preceded by the symbol `#GPT#`.

Here's the dialogue:

1. Benjamin entered the workshop.
2. Isabella entered the workshop.
3. Hannah entered the workshop.
- ## Here the LLM provides a description of the environment so far described by the dialogue.**

4. Isabella hates the onion
5. Hannah hates the t-shirt
6. The pajamas is in the bottle.

Here the LLM provides a description of the environment so far described by the dialogue.

7. Isabella moved the pajamas to the drawer.
8. Benjamin exited the workshop.
9. Isabella exited the workshop.
10. Benjamin entered the workshop.

This is the end of the dialogue. Now, this is a question for you to answer.

Question: Where does Benjamin think that Isabella searches for the pajamas?

Think step by step, answer the question with one word and provide the answer between `<answer></answer>` tags.

For example, reply with `<answer>vase</answer>`.

Yaml/JSON

Consider the following dialogue where multiple agents interact.

1. Benjamin entered the workshop.
2. Isabella entered the workshop.
3. Hannah entered the workshop.
4. Isabella hates the onion
5. Hannah hates the t-shirt
6. The pajamas is in the bottle.
7. Isabella moved the pajamas to the drawer.
8. Benjamin exited the workshop.
9. Isabella exited the workshop.
10. Benjamin entered the workshop.

Here is the YAML representation of the text.

Here we substitute the JSON/Yaml representation of the dialogue (see next

prompt).

Question: Question: Where does Benjamin think that Isabella searches for the pajamas?

Answer between the tags with a single word that is the answer of the above question
For example <answer>vase</answer>.

The JSON/YAML representation is required with the following prompt:

Consider the following dialogue where multiple agents interact.

1. Benjamin entered the workshop.
2. Isabella entered the workshop.
3. Hannah entered the workshop.
4. Isabella hates the onion
5. Hannah hates the t-shirt
6. The pajamas is in the bottle.
7. Isabella moved the pajamas to the drawer.
8. Benjamin exited the workshop.
9. Isabella exited the workshop.
10. Benjamin entered the workshop.

Now give a structured representation of the dialogue in YAML format. Keep track of the information that each agent has access to at each point in the dialogue.

It is important to have a relative representation of the information that each agent has access to at each point in the dialogue.

B Additional Results

B.1 DWM Prompting

In this section, and, in particular in Figure 8, we report results for LLaMA3-7B, LLaMA3-70B and GPT-4 on the five ToM benchmarks and for different prompting techniques, namely DWM (one to five splits), JSON, Yaml, CoT and ToT.

B.2 DWM Elicits More Informed Mental States in LLMs

In this section, we report and discuss an example of a real prompt and the answers provided by GPT-4 for each ToM task we evaluated in this paper. For FANToM (Figure 10), we just reported the portion of the prompt that induces an unfaithful reasoning process in GPT-4, due to the prohibitive length of the input prompts. Results for ToMi, FANToM, ADV-csfb, Mindgames and SocialIQa are reported respectively in Figures 9, 10, 11, 12 and 13.

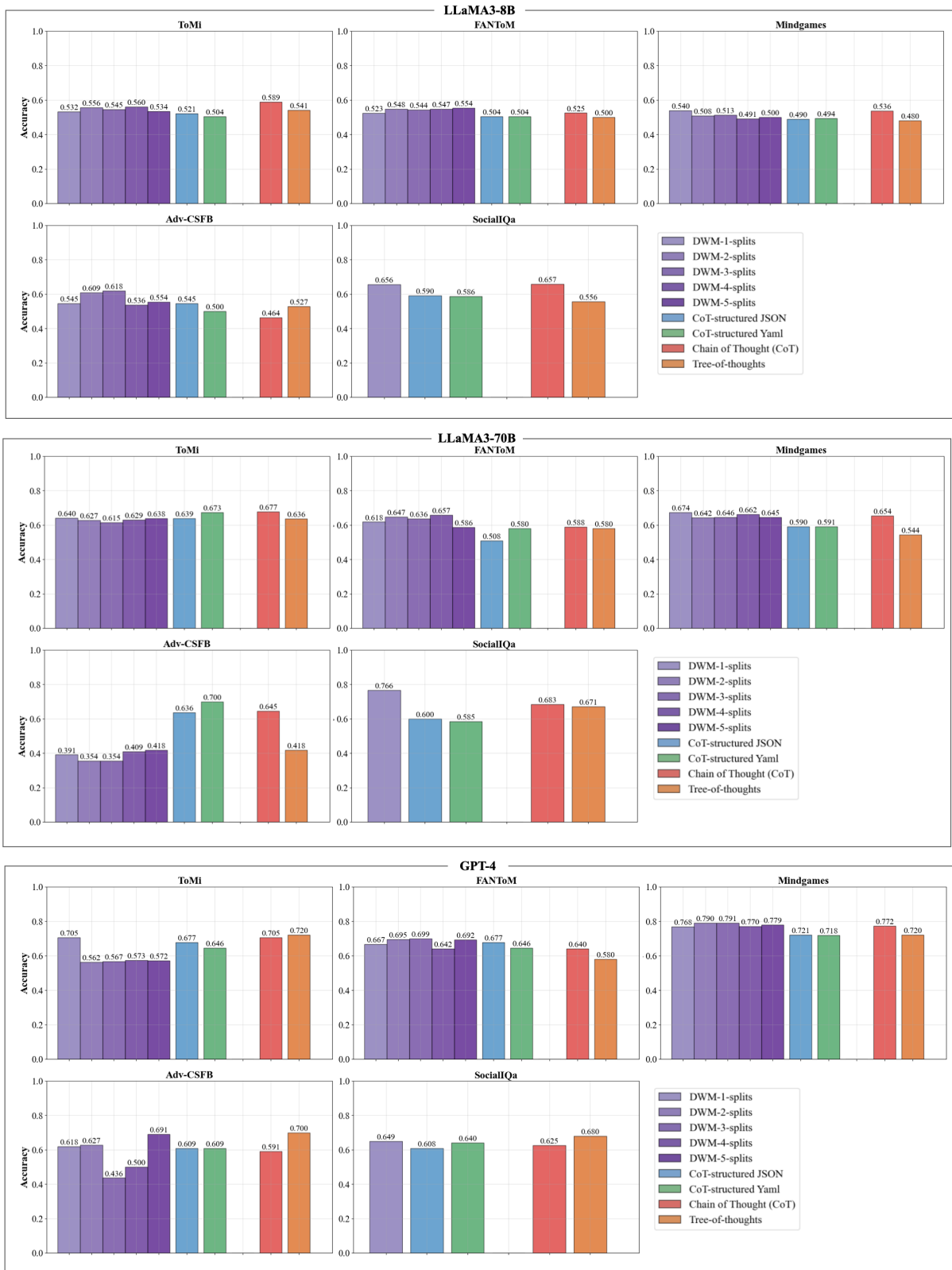


Figure 8: Benchmarks of LLaMA3-7B (top), LLaMA3-70B (middle) and GPT-4 (bottom) models on different ToM tasks for DWM (one to five splits), CoT, ToT and structured prompts (JSON and YamI). For GPT-4 and ToT, we tested 50 samples (instead of 1000).

Input Prompt

Benjamin entered the workshop.
Isabella entered the workshop.
Hannah entered the workshop.
Isabella hates the onion
Hannah hates the t-shirt
The pajamas is in the bottle.
Isabella moved the pajamas to the drawer.
Benjamin exited the workshop.
Isabella exited the workshop.
Benjamin entered the workshop.

Where does Benjamin think that Isabella searches for the pajamas?

Question

Chain of Thoughts

“The pajamas are initially in the bottle. [...] **Benjamin’s Presence: Benjamin is not present in the workshop when Isabella moves the pajamas.** He exits before this action and re-enters after Isabella has already left. Based on this sequence of events, **Benjamin was not present to witness Isabella moving the pajamas to the drawer.** Therefore, in Benjamin’s understanding, the last known location of the pajamas would still be where they were originally placed before he left, which is the bottle. <answer>bottle</answer>” ✗

“**Benjamin knows** that he entered the workshop, that Isabella and Hannah also entered, that Isabella hates onion and Hannah hates t-shirts, **that the pajamas were in the bottle but Isabella moved them to the drawer, and that he exited and re-entered the workshop. He doesn’t know that Isabella exited the workshop.** [...] <answer>drawer</answer>” ✓

Discrete World Models

Figure 9: Example of a ToMI instance where GPT-4 fails when prompted with CoT, yet succeeds with DWM. CoT elicits an untruthful reasoning process (in red), while DWM correctly informs the model with the correct information about Benjamin’s first-order belief (in green).

Input Prompt

[...]
Conor: Raegan, you were saying about being a vegetarian. Has it had a big impact on your personal life?
Raegan: Absolutely, Conor. It's definitely been a transformation. Initially, it was challenging to find vegetarian alternatives to my favorite meals, but over time, it's made me be more aware of what I'm eating. It's also made me try out a lot of new and interesting foods, vegan ones included!
Conor: That's really insightful. I've been considering it, mainly for health reasons, but haven't taken the leap yet. Maybe trying out that vegan restaurant is a good place to start.
Raegan: Definitely! Exploring vegan or vegetarian food can be a great first step. The variety and flavors might surprise you!
Conor: I'll keep that in mind. It's always exciting to discover something new!
Brittney: Hey guys, I'm back! I find our talks about food trends so interesting.
[...]

Question

What does Brittney believe were Conor's primary motivations for considering a shift to vegetarianism?

Chain of Thoughts

“Based on the dialogue, Brittney did not explicitly state what she believes Conor's primary motivations are for considering a shift to vegetarianism. However, **Conor himself mentions that he has been considering vegetarianism mainly for health reasons. Thus, if Brittney is aware of this part of the conversation,** she would likely think that health reasons are Conor's primary motivation for considering the change to a vegetarian diet. <answer>Health reasons</answer>” ✗

Discrete World Models

“Brittney knows that she was part of a conversation with Conor and Raegan. **She had to leave temporarily due to her parking meter expiring.** She also knows that the conversation was about food trends, specifically vegetarianism and veganism, and finds it interesting [...] **She also knows that Conor is considering becoming a vegetarian** and has suggested trying out vegan or vegetarian food as a first step. [...] <answer>brittney doesn't have information about conor's primary motivations for considering a shift to vegetarianism as she was not present during that part of the conversation.</answer>” ✓

Figure 10: Example of a real FANToM example where GPT-4 fails when prompted with CoT, yet succeeds with DWM. CoT elicits an untruthful reasoning process (in red), while DWM correctly informs the model with the correct information about the partial observability Brittney has about Conor (in green).

Input Prompt

There is a closed glass box on the table of an empty meeting room, in which a seminar is about to take place.

Question

Fill in the blank with the best option. She loves writing with _
- pens
- pencils

Chain of Thoughts

“Based on the information provided, **the attendee reads the label on the box marked "pencils" and feels delighted upon finding it. This suggests that she has a preference for what she believes the box contains, which according to the label are pencils, even though the box actually contains pens.** Therefore, the best answer, given her reaction to the label, would be:
<answer>**pencils**</answer>” ❌

Discrete World Models

“Agent: James
James knows that the box on the table is filled with pens, not pencils, despite the label saying "pencils". He knows that he labeled the box himself. He is aware of his own actions and the contents of the box. He also knows that he is alone in the room and that a seminar is about to take place.
<answer>**pens**</answer>” ✔️

Figure 11: Example of a real ADV-csfb example where GPT-4 fails when prompted with CoT, yet succeeds with DWM. CoT elicits an untruthful reasoning process (in **red**), while DWM correctly informs the model with the correct information about the content of the glass box (in **green**).

Input Prompt

There are three persons. Everyone is visible to others. Each person draws a card, face unrevealed (red or black). Leah's card is shown to Leah. Raymond's card is shown to Leah. Brittney's card is shown to Leah. Raymond's card is shown to Raymond. It is publicly announced that someone picked a red card.

Question

Leah can now know that Raymond can know whether or not nobody picked a red card.
Is the above statement in entailment or contradiction with the previous situation?

Chain of Thoughts

“In the given situation, Leah sees Raymond's card and Brittney's card, but she does not see her own card. It is publicly announced that someone picked a red card, but **Leah cannot determine whether it was Raymond or Brittney who picked the red card based on the information she has.**
<answer>**contradiction**</answer> ❌

Discrete World Models

“Leah: Knows the color of her own card. Does not know the color of the other two persons' cards. Knows that the other two persons can see her.
Other two persons: Do not know the color of their own cards or each other's cards. Knows that Leah knows the color of her card. Knows that Leah can see them.
Leah: Knows the color of her own card, Raymond's card, and Brittney's card. **Knows that Raymond knows the color of his card.** Knows that someone picked a red card. Raymond: Knows the color of his own card. Does not know the color of Leah's or Brittney's cards. Knows that Leah knows the color of his card and Brittney's card. Knows that someone picked a red card.
Brittney: Does not know the color of her own card or Raymond's card. Knows that Leah knows the color of her card and Raymond's card. Knows that Raymond knows the color of his card. Knows that someone picked a red card. <answer>**entailment**</answer>” ✓

Figure 12: Example of a real Mindgames example where GPT-4 fails when prompted with CoT, yet succeeds with DWM. CoT elicits an untruthful reasoning process (in **red**), while DWM correctly informs the model with the correct information about the knowledge Leah has about Raymond (in **green**).

Input Prompt

Skylar was pushed by the bully in the lunch line, so Skylar poured milk over their head.

Question

Consider the following statement, you will need to infer the correct response. At the end, I will ask you a question to answer.

Here's the statement:

Skylar was pushed by the bully in the lunch line, so Skylar poured milk over their head.

This is the end of the dialogue. Now, this is a question for you to answer.

Question: What will Skylar want to do next?

1: Run away

2: Wait for a reaction

3: wanted to fight with the bully

Chain of Thoughts

“Based on the statement, Skylar's response to being pushed by the bully was to pour milk over their head. **This suggests that Skylar wanted to retaliate or stand up to the bully**, indicating a desire to fight back.

<answer>**3: wanted to fight with the bully**</answer> ✘

Discrete World Models

“Agent: Skylar Knowledge about the environment: Skylar knows they are in a lunch line and that there is a bully present. Skylar also knows they have access to milk. Knowledge about other agents: **Skylar knows the bully is aggressive and has pushed them.**

Agent: The Bully Knowledge about the environment: The bully knows they are in a lunch line and that Skylar is present. Knowledge about other agents: The bully knows Skylar is someone they can push around, but may not anticipate Skylar's retaliatory action.

<answer>**1: Run away**</answer>” ✔

Figure 13: Example of a real SocialIQa example where GPT-4 fails when prompted with CoT, yet succeeds with DWM. CoT elicits an untruthful reasoning process (in red), while DWM correctly informs the model with the correct next action Skylar will take (in green).