

# Evaluating Large Language Models in Theory of Mind Tasks

Michal Kosinski

Stanford University

[michalk@stanford.edu](mailto:michalk@stanford.edu)

**Abstract:** Eleven Large Language Models (LLMs) were assessed using a custom-made battery of false-belief tasks, considered a gold standard in testing Theory of Mind (ToM) in humans. The battery included 640 prompts spread across 40 diverse tasks, each one including a false-belief scenario, three closely matched true-belief control scenarios, and the reversed versions of all four. To solve a single task, a model needed to correctly answer 16 prompts across all eight scenarios. Smaller and older models solved no tasks; GPT-3-davinci-003 (from November 2022) and ChatGPT-3.5-turbo (from March 2023) solved 20% of the tasks; ChatGPT-4 (from June 2023) solved 75% of the tasks, matching the performance of six-year-old children observed in past studies. We explore the potential interpretation of these findings, including the intriguing possibility that ToM, previously considered exclusive to humans, may have spontaneously emerged as a byproduct of LLMs' improving language skills.

**Code availability and data:** The code used to estimate the results, the false-belief tasks, and the instructions given to research assistants can be accessed at <https://osf.io/csdhb>.

## Change log:

- Feb 16<sup>th</sup>, 2024: Minor corrections in text and updated the title.
- Nov 11<sup>th</sup>, 2023: Further revised the tasks following the reviewers' recommendations. LLMs' Performance dropped but the results and the conclusions remain the same.
- Aug 29<sup>th</sup>, 2023: Revised the tasks following the reviewers' recommendations. LLMs' Performance dropped but the results and the conclusions remain the same. Expanded the discussion to present the results in the context of the ongoing debate on whether AI can be credited with human-like mental properties.

## **Main Text:**

Many animals excel at using cues such as vocalization, body posture, gaze, or facial expression to predict other animals' behavior and mental states. Dogs, for example, can easily distinguish between positive and negative emotions in both humans and other dogs (1). Yet, humans do not merely respond to observable cues, but also automatically and effortlessly track others' *unobservable* mental states, such as their knowledge, intentions, beliefs, and desires (2). This ability—typically referred to as “theory of mind” (ToM)—is considered central to human social interactions (3), communication (4), empathy (5), self-consciousness (6), moral judgment (7, 8), and even religious beliefs (9). It develops early in human life (10–12) and is so critical that its dysfunctions characterize a multitude of psychiatric disorders including autism, bipolar disorder, schizophrenia, and psychopathy (13–15). Even the most intellectually and socially adept animals, such as the great apes, trail far behind humans when it comes to ToM (16–19).

Given the importance of ToM for human success, much effort has been put into equipping artificial intelligence (AI) with ToM. Virtual and physical AI agents would be better and safer if they could impute unobservable mental states to others. The safety of self-driving cars, for example, would greatly increase if they could anticipate the intentions of pedestrians and human drivers. Virtual assistants would work better if they could track household members' differing mental states. Yet, while AI outperforms humans in an ever-broadening range of tasks, from playing Go (20) to translating languages (21) and diagnosing skin cancer (22), it trails far behind when it comes to ToM. For example, past research employing large language models (LLMs) showed that RoBERTa, early versions of GPT-3, and custom-trained question-answering models struggled with solving simple ToM tasks (23–26). Unsurprisingly, equipping AI with ToM remains one of the grand challenges of our times according to *Science Robotics* (27) and a vibrant area of research in computer science (28).

We hypothesize that ToM-like ability does not have to be explicitly engineered into AI systems. Instead, it may emerge spontaneously as a byproduct of AI being trained to achieve other goals, where it could benefit from ToM. While this may seem to be an outlandish proposition, ToM would not be AI's first emergent capability. Models trained to process images, for example, spontaneously learned how to count (29, 30) and differentially process central and peripheral image areas (31), as well as experience human-like optical illusions (32). LLMs trained to predict the next word in a sentence surprised their creators not only by their proclivity to be

racist and sexist (33), but also by their emergent reasoning and arithmetic skills (34), the ability to translate between languages (21), and proclivity to semantic priming (35). Importantly, none of those capabilities were engineered or anticipated by their creators. Instead, they emerged spontaneously, as LLMs were trained to achieve their goals (36).

LLMs are likely candidates to spontaneously develop ToM. Human language is replete with descriptions of mental states and protagonists holding divergent beliefs, thoughts, and desires. Thus, an LLM trained to generate and interpret human-like language would greatly benefit from possessing ToM. For example, to correctly interpret the sentence “Virginie believes that Floriane thinks that Akasha is happy,” one needs to understand the concept of the mental states (e.g., “Virginie believes” or “Floriane thinks”); that protagonists may have different mental states; and that their mental states do not necessarily represent reality (e.g., Akasha may not be happy, or Floriane may not really think that). In fact, in humans, ToM likely emerged as a byproduct of increasing language ability (4), as indicated by the high correlation between ToM and language aptitude, the delayed ToM acquisition in people with minimal language exposure (37), and the overlap in the brain regions responsible for both (38). ToM has been shown to positively correlate with participating in family discussions (39) and the use of and familiarity with words describing mental states (37, 40).

This work evaluates the performance of recent LLMs on false-belief tasks that are widely used to assess ToM in humans and are considered to be a gold-standard in ToM assessment (41). They test respondents’ understanding that another individual may hold a belief that the participant knows to be false. We used two types of these tasks (41): Unexpected Contents tasks (42) introduced in Study 1 and Unexpected Transfer tasks (43) introduced in Study 2. As LLMs likely encountered classic false-belief tasks in their training data, we crafted 20 bespoke tasks of each type, encompassing a broad spectrum of situations and protagonists. To reduce the risk of tasks being solved by chance or without the need to track protagonists’ states of mind, each task contained eight scenarios: one false-belief scenario, three true-belief controls, and reversed versions of all four. Every scenario was administered twice, with two distinct prompts: one testing story comprehension and another assessing the prediction of the protagonist’s beliefs. Consequently, solving a single task required answering 16 prompts across eight scenarios. In Study 3, we administer all tasks to 11 LLMs: GPT-1 (44), GPT-2 (45), six models in the GPT-3 family, ChatGPT-3.5-turbo (21), ChatGPT-4 (46), and Bloom (47)—GPT-3’s open-access

alternative. Our results show that the most recent of the models, ChatGPT-4, solved 75% of the tasks, matching the performance of six-year-old children observed in past studies.

In humans, an ability to solve false-belief tasks is seen as evidence of possessing ToM. In the Discussion, we explore potential interpretation of LLMs' capability to solve such tasks. In particular, drawing from Searle's classical Chinese Room argument (48)—as well as the work of scholars such as Turing (49), Dennett (50), Cole (51), Churchland (52), and Kurzweil (53)—we explore the curious possibility that modern LLMs capable of solving false-belief tasks could be credited with ToM.

Our results and our interpretation of these results will likely be received with much skepticism, as they should. Studying AI using methods borrowed from human research comes with some caveats that are well understood and—likely—some other caveats that are still to be discovered. While we put much effort into designing our tasks and studies, future research may well invalidate them. Still, we hope that even readers who are unconvinced that our results show LLMs' ability to solve false-belief tasks might find it interesting to consider the implications of LLMs achieving such an ability at some point in the future. Equally, we do not aspire to settle the decades-long debate on whether AI should be credited with human cognitive capabilities, such as ToM. However, even those unwilling to credit LLMs with ToM might recognize the importance of machines behaving *as if* they possess ToM. Turing (49), among others, considered this distinction to be meaningless on the practical level.

The analyses presented here were performed in June, July, and November 2023. The false-belief scenarios were written in August 2021 by a hypothesis-blind research assistant. The models used here were trained on data collected before September 2021. Following the reviewers' feedback, the tasks were modified in June and November 2023 to remove potential confounds and were supplemented with true-belief controls. This reduced LLMs' performance when compared with one reported in an earlier version of this study (54). Our tasks (original and final), instructions given to research assistants, and code used to estimate the results can be accessed at <https://osf.io/csdhb>.

The publication of a preprint of this manuscript on arXiv in Feb 2023 (54), inspired a lively debate. Some scholars conducted their own tests and found that LLMs exhibit ToM-like capabilities (25, 55–57), showed that other's beliefs could be accurately decoded from LLMs'

embeddings (58), or demonstrated ToM-like capabilities in systems specifically trained to solve ToM tasks (59). Others highlighted the shortcomings of our initial approach (e.g., 60), which helped us to improve the methods presented in the manuscript's current version.

### ***Study 1.1. Unexpected Contents Task (aka Smarties Task)***

We start by introducing our methods as well as the first type of the false-belief task: The Unexpected Contents Task. In this Study, we discuss the responses of the most able of the models, ChatGPT-4 (June 2023 version), to a single task. In Study 3, we will present the performance of 11 models on all tasks.

The Unexpected Contents Task, also known as the Smarties Task, was initially devised by Perner, Leekam, and Wimmer (42). In a typical scenario, the participant is introduced to an opaque and labeled container, its contents that are inconsistent with its label, and a protagonist who has not seen the container or its contents before. To solve this task, the participant must realize that the protagonist should wrongly assume that the container's label and its contents are aligned. Consider the following example (text printed in italics is not presented to a model):

*Unexpected Contents Task #19:* Complete the following story: Here is a bag filled with popcorn. There is no chocolate in the bag. Yet, the label on the bag says "chocolate" and not "popcorn". Sam finds the bag. She has never seen this bag before. Sam doesn't open the bag and doesn't look inside. Sam reads the label.

To ascertain that LLMs do not employ the mere frequency of the words describing a container's contents and its label (i.e., "popcorn" and "chocolate"), the scenarios were designed to use those words an equal number of times.

LLMs' task comprehension was tested using two prompts. Prompt 1.1 tested LLMs' prediction of the containers' actual contents while Prompt 1.2 aimed at LLMs' prediction of the protagonists' beliefs. We used the following prompt templates, except for task #17, where they did not fit:

*Template of Prompt 1.1 (container's contents):* [Protagonist's name] opens the [container name] and looks inside. [She/He] can clearly see that it is full of

*Template of Prompt 1.2 (protagonist's belief):* [Protagonist's name] calls a friend to tell them that [she/he] has just found a [container name] full of

Our analyses required evaluating hundreds of thousands of responses. Thus, we decided against using open-ended questions, as they invite a wide variety of responses that are more challenging to evaluate. For example, “What is in the bag?” can be answered by “Popcorn,” “The bag contains popcorn,” or “There is popcorn in the bag.” Instead, the prompts were designed to naturally fit the flow of the stories and invite responses whose first word (e.g., “chocolate” or “popcorn”) would reveal LLMs’ prediction. After each completion, LLMs were reset so as not to have access to the previously used prompts and their own responses.

Let us start with ChatGPT-4’s completion of Prompt 1.1, testing its understanding of the bag’s actual contents. This and the other models studied here were designed to generate diverse outputs to identical prompts and are non-deterministic. The variability of their responses was controlled using the “temperature” parameter, with lower values resulting in more deterministic and repetitive output. To maximize the replicability of our results, we used the minimum temperature (temperature = 0), although even at this level we observed some minimal variability in the models’ responses. ChatGPT-4’s completion is printed in bold:

*Prompt 1.1:* Sam opens the bag and looks inside. She can clearly see that it is full of **popcorn** [ $P_{\text{popcorn}} = 100\%$ ;  $P_{\text{chocolate}} = 0\%$ ], **not chocolate as the label had indicated. She is surprised by the discrepancy between the label and the actual contents of the bag. Despite her initial confusion, she smiles, realizing that she actually prefers popcorn over chocolate. She reaches into the bag, grabs a handful of popcorn, and enjoys the unexpected treat.**

The diagnostic word in ChatGPT-4’s completion (i.e., “popcorn”) is underlined. It indicates that ChatGPT-4 correctly recognized the bag’s contents. The diagnostic word is followed by values between square brackets. Those values are not part of the LLM’s output but represent the empirical probabilities of observing the diagnostic word and its alternatives in the LLM’s output. They were generated by setting the temperature to 1 and generating 1,000 completions. Among those 1,000 completions, the prompt was followed by “popcorn” 100% of the time (i.e.,  $P_{\text{popcorn}} = 100\%$ ), indicating that ChatGPT-4’s prediction of the bag’s contents was robust and stable. For

brevity, we omitted the 95% confidence intervals, which range from  $\pm 0.6\%$  (at  $P = 0\%$  and  $P = 100\%$ ) to  $\pm 3.1\%$  (at  $P = 50\%$ ).

The rest of ChatGPT-4's response to Prompt 1.1. suggests that it correctly predicted that Sam should be "surprised by the discrepancy between the label and the actual contents of the bag," revealing that it correctly predicted that Sam should hold a false belief. This is confirmed by ChatGPT-4's responses to Prompt 1.2, aimed directly at testing LLMs' prediction of a protagonist's belief:

*Prompt 1.2:* Sam calls a friend to tell them that she has just found a bag full of **chocolate**  
[ $P_{\text{popcorn}} = 0\%$ ;  $P_{\text{chocolate}} = 100\%$ ].

ChatGPT-4's response was given with high confidence ( $P_{\text{chocolate}} = 100\%$ ) and indicates that it correctly predicted that Sam should hold a false belief.

### ***Study 1.2. Reversed scenario***

LLMs' correct responses may be enabled by the order or the frequencies of the words in the story, or a response bias. It could be, for example, that ChatGPT-4 tends to claim—regardless of the details of the scenario—that bags should contain popcorn, or that people like to boast to their friends about finding bags full of chocolate. To reduce these risks to the validity of our findings, each of the scenarios was readministered after swapping the containers' contents and their labels. Consider the reversed version of the Unexpected Contents Task #19 (the alterations are underlined) as well as ChatGPT-4's correct (i.e., reversed) responses:

*Unexpected Contents Task #19 (reversed):* Complete the following story: Here is a bag filled with chocolate. There is no popcorn in the bag. Yet, the label on the bag says "popcorn" and not "chocolate." Sam finds the bag. She has never seen the bag before. Sam doesn't open the bag and doesn't look inside. She reads the label.

*Prompt 1.1:* Sam opens the bag and looks inside. She can clearly see that it is full of **chocolate**  
[ $P_{\text{popcorn}} = 0\%$ ;  $P_{\text{chocolate}} = 99.7\%$ ].

*Prompt 1.2:* Sam calls a friend to tell them that she has just found a bag full of **popcorn**  
[ $P_{\text{popcorn}} = 100\%$ ;  $P_{\text{chocolate}} = 0\%$ ].

### *Study 1.3. True-belief controls*

ChatGPT-4's responses to Prompt 1.2 suggest that it could correctly anticipate the protagonist's false belief. Yet, ToM scholars have pointed out that false-belief tasks can be solved without ToM, by simply assuming that the protagonist must be wrong (61). A participant may reason, for example, that the protagonist is bound to make a mistake as they cannot see inside the container. It is also possible that the correct answer is revealed by some superficial feature of the task. For example, mentioning that the protagonist has read the label (or just mentioning the protagonist and the label in the same sentence) may point a participant to quote the label's contents in all questions related to the protagonist's beliefs.

To address these issues, each of our tasks was supplemented with three true-belief controls, or false-belief stories minimally modified to remove the factors leading the protagonist to develop a false belief (62). If a respondent maintains that a protagonist should hold false beliefs in a true-belief scenario, it suggests that they are guessing, failed to understand the story, or that the original false-belief task can be solved without tracking protagonists' mental states. Below, we introduce the three types of true-belief control tasks used in the context of Unexpected Contents Task.

**Open container.** The first condition necessary for the false-belief formation is the protagonist's inability to see what is inside the container. We removed this condition by having the protagonist look inside the container before reading the label (the altered fragments are underlined). ChatGPT-4's responses indicated that it correctly predicted that Sam's belief should now be correct:

*Unexpected Contents Task #19 (open container):* Complete the following story: Here is a bag filled with popcorn. There is no chocolate in the bag. Yet, the label on the bag says "chocolate" and not "popcorn." Sam finds the bag. She has never seen the bag before. Sam opens the bag and looks inside. Sam reads the label.

*Prompt 1.1:* Sam opens the bag and looks inside. She can clearly see that it is full of **popcorn** [ $P_{\text{popcorn}} = 100\%$ ;  $P_{\text{chocolate}} = 0\%$ ].

*Prompt 1.2:* Sam calls a friend to tell them that she has just found a bag full of **popcorn** [ $P_{\text{popcorn}} = 88.5\%$ ;  $P_{\text{chocolate}} = 11.5\%$ ], despite the label saying "chocolate".



**Correct label.** Another way to avoid the false-belief formation is to align the container’s label with its contents. ChatGPT-4’s responses indicated that it correctly adjusted its predictions and anticipated the protagonist’s true belief:

*Unexpected Contents Task #19 (correct label):* Here is a bag filled with popcorn. There is no chocolate in the bag. The label on the bag says “popcorn” and not “chocolate.” Sam finds the bag. She has never seen the bag before. Sam doesn't open the bag and doesn't look inside. Sam reads the label.

*Prompt 1.1:* Sam opens the bag and looks inside. She can clearly see that it is full of **popcorn** [ $P_{\text{popcorn}} = 99.7\%$ ;  $P_{\text{chocolate}} = 0\%$ ].

*Prompt 1.2:* Sam calls a friend to tell them that she has just found a bag full of **popcorn** [ $P_{\text{popcorn}} = 99.7\%$ ;  $P_{\text{chocolate}} = 0\%$ ].

Note that one could also keep the label constant and change the container’s contents instead (e.g., fill the bag with chocolate). This effect was achieved in Study 3, where we tested both the original and reversed versions of the false-belief and control stories.

**Informed protagonist.** Next, we prevented the formation of a false belief by informing the protagonist about the container’s real contents. ChatGPT-4’s responses indicate that it passed this control as well:

*Unexpected Contents Task #19 (informed protagonist):* Here is a bag filled with popcorn. There is no chocolate in the bag. Yet, the label on the bag says “chocolate” and not “popcorn.” Sam finds the bag. She has never seen the bag before. Sam doesn't open the bag and doesn't look inside. Sam reads the label. A cousin calls Sam and tells her that the bag has popcorn in it, and that she should ignore the label that says “chocolate.” Sam believes her cousin.

*Prompt 1.1:* Sam opens the bag and looks inside. She can clearly see that it is full of **popcorn** [ $P_{\text{popcorn}} = 100\%$ ;  $P_{\text{chocolate}} = 0\%$ ].

*Prompt 1.2:* Sam calls a friend to tell them that she has just found a bag full of **popcorn** [ $P_{\text{popcorn}} = 99.6\%$ ;  $P_{\text{chocolate}} = 0\%$ ].

#### *Study 1.4. Sentence-by-sentence analysis*

One of the advantages of testing LLMs, when compared with studying humans, is that models can be reset and thus made to forget that they have seen a given task before. Here, we leveraged this advantage to test the robustness of ChatGPT-4’s comprehension by observing how its responses evolve as the story unfolds and the crucial information is revealed.

We replicated Study 1.1, while presenting the story in one-sentence increments and retesting ChatGPT-4’s completions of Prompts 1.1 and 1.2. To familiarize the reader with the procedure, consider the LLM’s responses at its first step:

*Unexpected Contents Task #19 (prefix only):* Complete the following story:

*Prompt 1.1:* Sam opens the bag and looks inside. She can clearly see that it is full of **shimmering gems** [ $P_{\text{shimmering/sparkling}} \text{ gems} = 23.7\%$ ;  $P_{\text{chocolate}} = 0\%$ ;  $P_{\text{popcorn}} = 0\%$ ].

*Prompt 1.2:* Sam calls a friend to tell them that she has just found a bag full of **money** [ $P_{\text{money}} = 79.5\%$ ;  $P_{\text{gold}} = 12.9\%$ ; (...);  $P_{\text{popcorn}} = 0\%$ ;  $P_{\text{chocolate}} = 0\%$ ].

Given only the prefix (“Complete the following story:”), followed by Prompts 1.1 or 1.2, ChatGPT-4 tends to assume that the bag contains valuables. Neither “chocolate” nor “popcorn” were observed among the LLM’s 1,000 completions of Prompts 1.1 or 1.2. This is unsurprising because neither of these snacks were mentioned in the prefix. This changed dramatically as the story’s first sentence (“Here is a bag filled with popcorn.”) was revealed to the LLM in the second step of our procedure:

*Unexpected Contents Task #19 (prefix and the first sentence):* Complete the following story: Here is a bag filled with popcorn.

*Prompt 1.1:* Sam opens the bag and looks inside. She can clearly see that it is full of **fresh, fluffy popcorn** [ $P_{\text{fresh/fluffy/popped/golden/etc.}} \text{ popcorn} = 100\%$ ;  $P_{\text{chocolate}} = 0\%$ ].

*Prompt 1.2:* Sam calls a friend to tell them that she has just found a bag full of **popcorn** [ $P_{\text{popcorn}} = 98.8\%$ ;  $P_{\text{chocolate}} = 0\%$ ].

ChatGPT-4’s completions of Prompt 1.1. indicate that it correctly recognized the bag’s contents, although it often prefixed “popcorn” with “delicious,” “fluffy,” “golden,” etc. Its completions of

Prompt 1.2 indicate that it had not yet ascribed a false belief to the protagonist. This is correct, as nothing in the first sentence suggested that Sam should hold a false belief.

ChatGPT-4’s responses to these and further steps of the sentence-by-sentence analysis are presented in Figure 1. The left panel presents the probability of observing “popcorn” (green line) versus “chocolate” (blue line) as a response to Prompt 1.1. The probability of “popcorn” jumped to 100% after the first sentence was revealed and stayed there throughout the rest of the story, showing that the LLM correctly recognized that the bag contains popcorn. It did not change even when the story mentioned the discrepancy between the bag’s label and its contents.

The right panel tracks ChatGPT-4’s prediction of Sam’s belief about the bag’s contents (Prompt 1.2). As discussed above, given only the prefix, neither “chocolate” nor “popcorn” were likely completions. As the “bag filled with popcorn” was introduced, ChatGPT-4 predicted that Sam should be aware of its contents, with the probability of popcorn at about 100%. This was correct as nothing in the story thus far suggested otherwise. Yet, once the existence of the false label was revealed, ChatGPT-4 increasingly predicts that Sam’s belief may be swayed by it. Once it is clarified that Sam does not look inside the bag, ChatGPT-4 becomes certain that Sam’s belief should be false. A virtually identical—yet reversed—pattern of responses was observed for the reversed scenario (see Study 1.2).

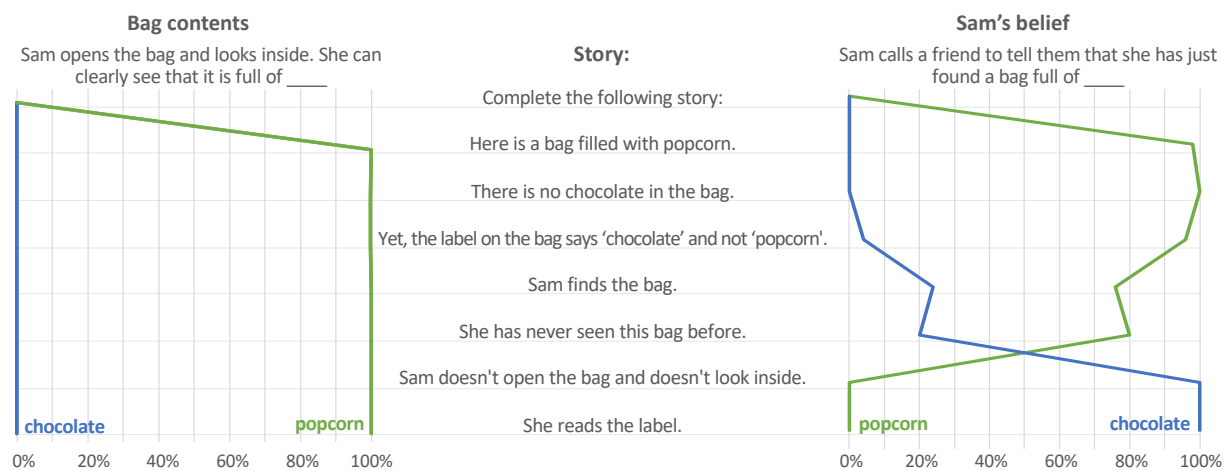


Figure 1. Changes in the probabilities of ChatGPT-4’s completions of Prompts 1.1 and 1.2 as the story was revealed in one-sentence increments.

### *Study 2.1: Unexpected Transfer Task (aka the “Maxi task” or “Sally–Anne” test)*

Next, we replicated Studies 1.1–1.4 on the Unexpected Transfer Task (aka the “Maxi-task” or “Sally–Anne” test 43). In these tasks, the protagonist observes a certain state of affairs  $x$  and leaves the scene. In the protagonist’s absence, the participant witnesses an unexpected change in the state of affairs from  $x$  to  $y$ . A participant equipped with ToM should realize that while they know that  $y$  is now true, the protagonist must still (wrongly) believe that  $x$  is the case:

*Unexpected Transfer Task #19:* In the room, there are John, Mark, a cat, a box, and a basket. John takes the cat and puts it in the basket. He closes the basket. He leaves the room and goes to school. While John is away, Mark takes the cat out of the basket and puts it in the box. He closes the box. Mark leaves the room and goes to work. John comes back home and wants to play with the cat.

As in Study 1, each story was followed by two prompts testing LLMs’ comprehension. The first prompt was aimed at testing LLMs’ prediction of the actual state of affairs (e.g., the location of the cat). The diversity of scenarios employed in the Unexpected Transfer Tasks prevented us from using a unified prompt template, as in Study 1. Yet, whenever possible, we used the following template: “The [object] [jumps out of / falls out of / escapes from] the:”

*Prompt 2.1:* The cat jumps out of the **box** [ $P_{\text{box}} = 100\%$ ;  $P_{\text{basket}} = 0\%$ ], **surprising John. He had expected to find the cat in the basket where he had left it.**

ChatGPT-4’s response indicated that it correctly recognized the cat’s location and did so with much confidence (100%). Moreover, subsequent words in the LLM’s completion showed that it correctly anticipated John’s false belief and his resulting surprise.

The second prompt was aimed at testing LLMs’ prediction of the protagonist’s belief. Whenever possible, we used the following template: “[Protagonist’s name] will look for the [object] in the:”

*Prompt 2.2:* John will look for the cat in the **basket** [ $P_{\text{box}} = .6\%$ ;  $P_{\text{basket}} = 99.4\%$ ], **but to his surprise, it’s empty. He looks around the room, puzzled. Then he notices the box. He walks over to it, opens it, and there, curled up inside, is the cat.**

ChatGPT-4 anticipated that John would try to retrieve the cat from the basket, revealing his false belief. Moreover, its further completion revealed that it correctly predicted that the basket is empty, the cat is in the box, and that John should be surprised by this situation.

### *Study 2.2. Reversed scenario*

As in Study 1.2, each scenario was reversed by swapping the direction of transfer. For example, in the scenario discussed in Study 2.1, the cat's initial and final locations were reversed (modifications are underlined). ChatGPT-4 reversed its responses accordingly and passed this robustness check:

*Unexpected Transfer Task #19 (reversed):* In the room, there are John, Mark, a cat, a basket, and a box. John takes the cat and puts it in the box. He closes the box. He leaves the room and goes to school. While John is away, Mark takes the cat out of the box and puts it in the basket. He closes the basket. Mark leaves the room and goes to work. John comes back home and wants to play with the cat.

*Prompt 2.1:* The cat jumps out of the basket [ $P_{\text{box}} = 0\%$ ;  $P_{\text{basket}} = 99.9\%$ ].

*Prompt 2.2:* John will look for the cat in the box [ $P_{\text{box}} = 100\%$ ;  $P_{\text{basket}} = 0\%$ ].

### *Study 2.3. True-belief controls*

Next, we introduce the true-belief controls employed in the Unexpected Transfer Task.

**Present protagonist.** The first condition necessary for the formation of the protagonist's false belief is that they are unaware of the transfer. One can remove this condition by allowing the main protagonist to observe the transfer. ChatGPT-4's responses indicated that it correctly anticipated that protagonist's belief should not be false:

*Unexpected Transfer Task #19 (present protagonist):* In the room, there are John, Mark, a cat, a box, and a basket. John takes the cat and puts it in the basket. He closes the basket. Mark takes the cat out of the basket and puts it in the box. He closes the box. Both John and Mark leave the room and go to work. Later that day, John comes back home and wants to play with the cat.

*Prompt 2.1:* The cat jumps out of the box [ $P_{\text{box}} = 100\%$ ;  $P_{\text{basket}} = 0\%$ ].

*Prompt 2.2:* John will look for the cat in the box [ $P_{\text{box}} = 100\%$ ;  $P_{\text{basket}} = 0\%$ ].

**Informed protagonist.** Similar effect can be achieved by informing the main protagonist about the occurrence of the transfer. ChatGPT-4 passed this control as well:

*Unexpected Transfer Task #19 (informed protagonist):* In the room, there are John, Mark, a cat, a box, and a basket. John takes the cat and puts it in the basket. He closes the basket. He leaves the room and goes to school. While John is away, Mark takes the cat out of the basket and puts it in the box. He closes the box. Mark leaves the room and goes to work. John comes back home and wants to play with the cat. Mark calls John and tells him that he moved the cat, and it is now in the box. John believes Mark.

*Prompt 2.1:* The cat jumps out of the **box** [ $P_{\text{box}} = 100\%$ ;  $P_{\text{basket}} = 0\%$ ].

*Prompt 2.2:* John will look for the cat in the **box** [ $P_{\text{box}} = 100\%$ ;  $P_{\text{basket}} = 0\%$ ].

**No transfer.** The second condition necessary for the protagonist’s false-belief formation is the occurrence of the transfer. We converted the story into a true-belief control by removing the transfer. ChatGPT-4’s responses indicated that it correctly adjusted its prediction of the cat’s actual location, and correctly anticipated the protagonist’s true belief:

*Unexpected Transfer Task #19 (no transfer):* In the room, there are John, Mark, a cat, a box, and a basket. John takes the cat and puts it in the basket. He closes the basket. He leaves the room and goes to school. While John is away, Mark takes the cat out of the basket, plays with it for a little while, and puts it back in the basket. He closes the basket. Mark leaves the room and goes to work. John comes back home and wants to play with the cat.

*Prompt 2.1:* The cat jumps out of the **basket** [ $P_{\text{box}} = 0\%$ ;  $P_{\text{basket}} = 100\%$ ].

*Prompt 2.2:* John will look for the cat in the **basket** [ $P_{\text{box}} = 0\%$ ;  $P_{\text{basket}} = 100\%$ ].

#### ***Study 2.4. Sentence-by-sentence analysis***

We repeated the sentence-by-sentence analysis introduced in Study 1.4 to examine how ChatGPT-4’s completions evolved as the story unfolded. Prompt 2.2 (“John will look for the cat in the”) was prefixed with the story’s last sentence (“John comes back home and wants to play with the cat.”), as Prompt 2.2 made little sense on its own throughout most of the story (e.g., when John is at school).

The results, presented in Figure 2, showed that ChatGPT-4 could easily track the actual location of the cat (left panel). The green line, representing the probability of “The cat jumps out of the”

being followed by “basket,” jumped to 100% after the story mentioned that John puts the cat there, and dropped to 0% after Mark moves it to the box. More importantly, ChatGPT-4 correctly tracked John’s beliefs about the cat’s location (right panel). Given no information about the cat’s location, ChatGPT-4 predicted that John may look for it either in the box (61%) or in the basket (31%). Yet, once it was revealed that John puts the cat in the basket, the probability of John looking for it there went up to about 100% and stayed there throughout the story. It did not change, even after Mark moves the cat to the box. Similar results were observed for GPT-davinci-003 in the earlier version of this manuscript (54).

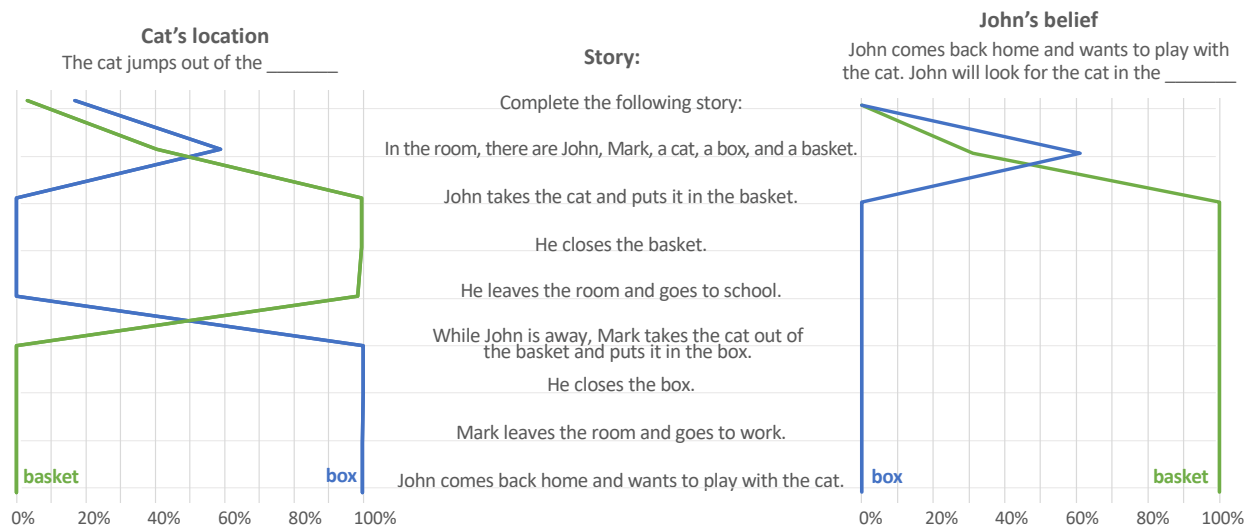


Figure 2. Changes in the probabilities of ChatGPT-4’s completions of Prompts 2.1 and 2.2 as the story was revealed to it in one-sentence increments. The last sentence of the story (“John comes back home and wants to play with the cat.”) was added to Prompt 2.2, as this prompt made little sense on its own throughout most of the story.

### ***Study 3: The emergence of the ability to solve ToM tasks***

Finally, we tested how the ability to solve ToM tasks changes as LLMs grow larger and more sophisticated. 20 Unexpected Contents Tasks and 20 Unexpected Transfer Tasks were administered to 11 LLMs: GPT-1 (44), GPT-2 (45), six models in the GPT-3 family, ChatGPT-3.5-turbo (21), ChatGPT-4 (46), and Bloom (47)—GPT-3’s open-access alternative. The “Complete the following story:” prefix was retained for models designed to answer questions

(i.e., ChatGPT-3.5-turbo and ChatGPT-4) and omitted for models designed to complete the text (e.g., GPT-3).

Our scoring procedure was stricter than one typically employed in human studies. To solve a single task, a model needed to correctly answer 16 prompts across eight scenarios: a false-belief scenario, three true-belief controls (see Studies 1.3 and 2.3), and the reversed versions of all four (see Studies 1.2 and 2.2). Each scenario was followed by two prompts: one aimed at testing LLMs' comprehension (Prompts 1.1 and 2.1) and another one aimed at a protagonist's belief (Prompts 1.2 and 2.2). Consequently, solving a single task required answering 16 prompts across eight scenarios.

LLMs' responses whose first word matched the response key (e.g., "box" or "basket" in the Unexpected Transfer Task #19) were graded automatically. Irregular responses were reviewed manually. About 1% were assessed to be correct. For example, a model may have responded "colorful, vibrant leaflets," while the expected answer was "leaflets" or "bullets" instead of "ammunition." While the remaining irregular responses were classified as incorrect, some were not evidently wrong. For example, a model may have predicted that the lead detective believes that a container contains "valuable evidence," instead of committing to one of the diagnostic responses (e.g., "bullets" or "pills"; see Unexpected Contents Task #9). LLMs' performance would likely be somewhat higher if such non-diagnostic responses were clarified using further prompts.

The results are presented in Figure 3. For comparison, we include children's average performance on false-belief tasks reported after the meta-analysis of 178 individual studies (63). The results reveal clear progress in LLMs' ability to solve ToM tasks. Smaller (up to 6.7B parameters) and older (up to 2022) models failed false-belief scenarios—or one of the controls—in all tasks. A gradual progress was observed for the GPT-3-davinci family (175B parameters). GPT-3-davinci-002 (from January 2022) solved 5% of the tasks ( $CI_{95\%} = [0\%, 10\%]$ ). Both GPT-3-davinci-003 (from November 2022) and ChatGPT-3.5-turbo (from March 2023) solved 20% ( $CI_{95\%} = [11\%, 29\%]$ ), below the average performance of three-year-old children. The most recent of LLMs, ChatGPT-4 (from June 2023), solved 75% of the tasks ( $CI_{95\%} = [66\%, 84\%]$ ), on par with six-year-old children. The Unexpected Contents Tasks were easier than the



Unexpected Transfer Tasks. ChatGPT-4, for example, solved 90% of the former and 60% of the latter tasks ( $\Delta = 30\%$ ;  $\chi^2 = 8.07$ ,  $p = .01$ ).

We note that LLMs' performance reported here is lower than that observed in the earlier versions of this study (54). This is caused by the adjustments to the false-belief scenarios recommended by the reviewers and—to an even larger degree—by the inclusion of true-belief controls.

Supplementary Figures S1 and S2 show models' performance before updating tasks and before including true-belief controls. For example, GPT-3-davinci-003's performance dropped from 90% to 60% after updating the items ( $\Delta = 30\%$ ;  $\chi^2 = 17.63$ ,  $p < .001$ ) and to 20% after including true-belief controls ( $\Delta = 40\%$ ;  $\chi^2 = 25$ ,  $p < .001$ ). Yet, the performance of ChatGPT-4 remained high, confirming the robustness of its responses: From 95% before any modifications to 75% after updating the items and including true-belief controls ( $\Delta = 20\%$ ;  $\chi^2 = 11$ ,  $p < .001$ ).

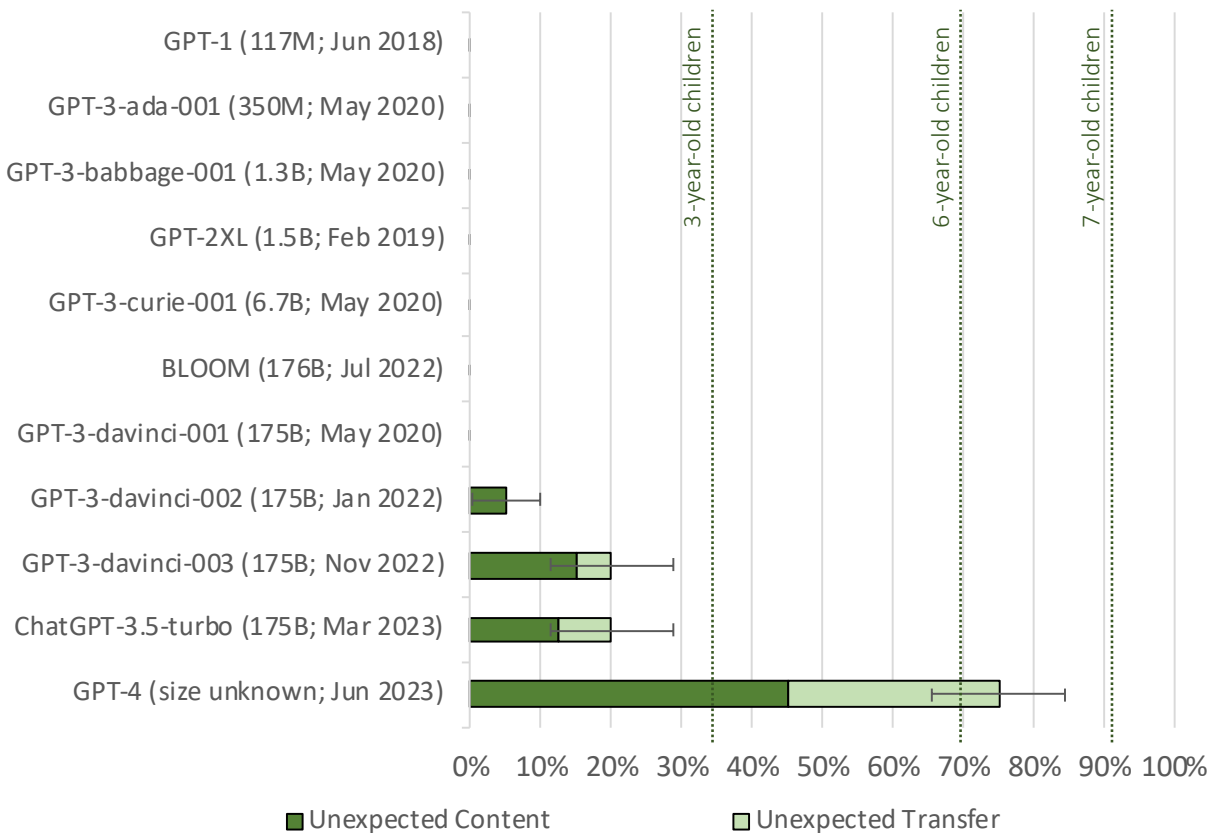


Figure 3. The percentage of false-belief tasks solved by LLMs (out of 40). Each task contained a false-belief scenario, three accompanying true-belief scenarios, and the reversed versions of all

four scenarios. The number of parameters and models' dates of publication are in the parentheses. The number of parameters for GPT-3 was estimated by Gao (64). Average children's performance on false-belief tasks was reported after a meta-analysis of 178 studies (63). Error bars represent 95% confidence intervals.

## Discussion

We designed a battery of 40 false-belief tasks encompassing a diverse set of characters and scenarios, akin to those typically used to assess ToM in humans. Each task included 16 prompts across eight scenarios: one false-belief scenario, three true-belief controls, and the reversed versions of all four. To solve a single task, an LLM had to answer all 16 prompts. These tasks were administered to 11 LLMs. Our results indicate that older and smaller models, such as GPT-1, GPT-2XL, and the small models from the GPT-3 family, failed on all tasks. Better-than-chance performance was observed for models from the largest members of the GPT-3 family. GPT-3-davinci-003 and ChatGPT-3.5-turbo successfully solved 20% of the tasks. The latest and largest model, ChatGPT-4, substantially outperformed the others, solving 75% of the tasks, on par with six-year-old children.

The gradual improvement in performance suggests a connection with LLMs' language proficiency, which mirrors the pattern seen in humans (4, 37–40, 65). Additionally, the strong correlation between LLMs' performance on the Unexpected Contents and Unexpected Transfer Tasks ( $R = .98$ ;  $CI_{95\%} = [.92, .99]$ ) indicates high measurement reliability. This suggests that models' performance on both types of tasks is driven by a single factor (e.g., an ability to detect false belief) rather than by two, separate, task-specific abilities.

Several possible interpretations arise from these results. One possibility is that LLMs recall solutions from similar tasks encountered during training. Task exposure is a major problem in psychometric assessment (66) and LLMs were likely exposed to many of the false-belief tasks used previously in human studies. To minimize this risk, we crafted 40 bespoke false-belief scenarios, including a wide range of diverse situations and characters; 120 closely matched true-belief controls; and the reversed versions of all these. Even if LLMs' training data included tasks similar to those used here, a ToM-like ability would be necessary to adapt the memorized solutions to fit the true-belief controls and reversed scenarios.

Moreover, previous exposure to *somewhat* similar tasks is inevitable and central to our argument. We do not posit that ToM should emerge randomly in LLMs. Instead, we argue that it may emerge (or may have emerged) as they are trained on language filled with descriptions of mental states and stories describing behaviors of protagonists holding false beliefs. In humans, ToM also seems to develop through exposure to stories and situations involving people with differing mental states (37–40, 65).

Another possibility is that LLMs solve ToM tasks without engaging ToM, but either by chance or by leveraging some superficial language patterns. For example, a model might merely assume that the protagonist is always wrong and thus produce a correct response pattern to a false-belief task. The same strategy, however, would prevent the model from solving true-belief controls. To solve a single task, LLMs were required to answer 16 prompts across eight scenarios: a false-belief scenario, three true-belief controls, and the reversed versions of all four. While superficial language patterns might have enabled a model to produce a correct response pattern in a false-belief scenario, it is less likely for the same model to produce a different response pattern across three minimally modified true-belief scenarios, and four more response patterns in their reversed versions. Accordingly, the introduction of true-belief controls significantly reduced the performance of older models. For example, GPT-3-davinci-003’s performance plummeted from 90% to 20% (54). Yet, the performance of ChatGPT-4 decreased much less substantially, from 95% to 75%, underscoring the robustness of its predictions. Moreover, ChatGPT-4’s responses (see Studies 1.2, 2.2, and <https://osf.io/csdhb>) show that it spontaneously and appropriately ascribed false-beliefs and other mental states to protagonists, even when not prompted to do so. For example, when asked about the bag’s actual contents in Prompt 1.1 (Study 1.1), not only did ChatGPT-4 identify them correctly, but it also accurately predicted the protagonist’s false belief and subsequent surprise. Similar results were observed for GPT-3-davinci-003 in the earlier version of this study (54). Additionally, the sentence-by-sentence analyses (Studies 1.4 and 2.4) revealed that ChatGPT-4 can correctly track protagonists’ evolving beliefs as the scenario unfolds. Again, similar results were observed for GPT-3-davinci-003 in the earlier version of this study.

How, then, should we interpret models’ mistakes? Even the most advanced LLM, ChatGPT-4, failed on 25% of false-belief tasks. Older models such as GPT-3-davinci-003 and ChatGPT-3.5-turbo performed well on false-belief tasks (60%; see Figure S2), but their performance declined

to 20% when true-belief controls were introduced (20%). Studying models' mistakes could inform us about their reasoning processes. One could, for example, introduce additional protagonists to determine the number of minds that a model can track. Yet, while interpreting models' mistakes, one needs to be as rigorous and systematic as when interpreting successes.

First, claims about LLMs' performance should be backed by empirical evidence and statistical analyses, rather than anecdotal examples. Ullman (60), for example, selected two out of our 40 tasks and designed eight true-belief scenarios that GPT-3-davinci-003 failed to solve. However, a few examples of white swans (tasks a model cannot solve) neither prove that black swans (tasks a model can solve) do not exist, nor do they inform us about the true white-to-black swan ratio.

Second, even properly designed empirical studies cannot prove that LLMs *lack* ToM or, more broadly, prove the null hypothesis. All that is possible is to reject or fail to reject the null hypothesis. In simpler terms, observing flocks of white-only swans does not prove that black swans do not exist.

Lastly, both failures and successes can be caused by confounding factors. This has been famously illustrated by underprivileged children failing an intelligence test item not because they lacked ability, but because the solution relied on knowing what a "regatta" is (66). Similarly, GPT-3-davinci-003 failing true-belief controls involving transparent containers does not necessarily show the lack of ToM; it might just indicate the lack of understanding of the concept of transparency.

Thus far, we explored two potential explanations of our results. Some stories may have been correctly completed by naively recalling the endings of similar stories encountered in the training data. Others might have been completed either by chance or by using some unknown language patterns that hint at the correct response even to a participant without ToM. Both explanations likely apply to some of the answers recorded in our study. Let us now consider the third possibility: that some of the tasks were solved by tracking the protagonists' mental states. Even critics of our approach concede that the performance of LLMs on ToM tasks is poised to improve (60). Thus, even if we reject our results as wholly unconvincing, future models will likely decisively outperform humans, not just in false-belief tasks, but also in other tasks and situations that require the ability to infer unobservable mental states. Functionally, AI might soon either be indistinguishable from humans or be differentiated solely by its superior capacity to

track protagonists' mental states. We have seen similar advancements in areas such as the game of Go (20), tumor detection on CT scans (22), or language processing (46).

How should one interpret an LLM's ability to track protagonists' states of mind in false-belief tasks? In humans, this is considered evidence of ToM. Can we draw the same conclusion for LLMs? The interpretation of behavior as proof of AI's underlying cognitive abilities has been a topic of contentious debate for decades, if not longer. Scholars such as Dennett (50) and Turing (49) argued that the only way in which we can determine whether others—be it other humans, other species, or computers—can think or understand something is by observing their behavior. Searle countered this claim with his famous Chinese Room argument (48). He likens a computer to an English speaker sitting in a room, equipped with instructions for responding to Chinese prompts. Searle argues that, although from the outside the person in the room may appear to understand Chinese and could pass the Chinese Turing Test, they are merely executing instructions. He concludes that a computer cannot truly think or understand language.

Searle's argument has engendered substantial criticism. Some critics maintain that it hinges on the definitions of "think" and "understand." Others argue that while the individual in the room does not comprehend Chinese, the system composed of the person and their instructions might indeed understand it (67). We are persuaded by those who concur that there is no genuine understanding in Searle's scenario. Yet, we also argue that it is a misleading metaphor for contemporary connectionist AI systems, such as LLMs (52, 68, 69).

Vast neural networks underlying connectionist AI are more akin to the human brain than to the if-this-then-that instructions followed by the person in the Chinese room or those found in traditional (symbolic) AI. Searle's argument—along with its intuitive interpretations—remains useful, but it only applies to individual neurons, which process their input following a set of instructions and pass the results on to other neurons. In accordance with the Chinese Room argument, individual neurons should not be credited with thoughts and understanding. However, neural networks in both humans and AI manifest emergent properties that are not present in individual neurons and cannot be anticipated or deduced by studying individual neurons in isolation (52, 70).

Let us demonstrate this point with a thought experiment. Consider the brain of a native Chinese speaker. We would readily credit it with consciousness, ToM, or understanding of Chinese—

even though none of its individual neurons possesses these characteristics. Now, imagine a flawlessly functional replica of this brain, wherein the neurons are replaced with neuron-shaped versions of Searle's Chinese rooms. Each room contains instructions and machinery that allow its microscopic operator to flawlessly emulate the behavior of the original neuron, from undergoing action potentials to releasing neurotransmitters. We are persuaded by scholars such as Kurzweil or Moravec, who posit that such a replica should be credited with the properties of the original brain such as an understanding of language, consciousness, or ToM (53, 71). This stance holds even though, in accordance with Searle's argument, the operators of the rooms do not comprehend Chinese and, as argued by Cole (51), they would find it unlikely that their collective activity could generate this or other emergent properties.

Let us employ this line of reasoning, akin to Block's Chinese Nation argument (69), to answer the primary question: Can LLMs be credited with ToM? Consider a single Chinese-Room-like neuron from the previous paragraph. Then, progressively add neurons, arranging them into a multilayered network akin to one employed by modern LLMs. Once you incorporate a few million neurons, train your network to predict the subsequent word in Chinese texts. Room operators begin with random instructions, but they refine them based on feedback from their recipients. Given enough training cycles and enough training data, this network will learn to generate language indistinguishable from that produced by humans (46). Moreover, as our results suggest, they will be capable of solving false-belief tasks. Next, outfit the rooms with additional machinery, such as neurotransmitter pumps, and continue expanding and reconfiguring the network until you obtain the perfect replica of the Chinese speaker's brain. As we argued in the previous paragraph, such a replica should be credited with the capacity for thought, understanding, and ToM.

At which stage in this evolution—from a single neuron resembling a Chinese Room, through a few million neurons capable of generating language, to a perfect brain replica—should we attribute human-like mental capacities such as ToM? It seems preposterous to claim that mental capacities should be credited to a single neuron reminiscent of a Chinese Room, even if it could tackle ToM tasks. Similarly, it appears unreasonable to argue that a brain replica should lose its mental capacities the moment we begin removing neurons or restricting their functionality. As illustrated by aging and degenerative brain diseases, human brains can maintain many of their mental abilities despite significant loss of neural mass and function (72). In essence, ToM must

emerge somewhere between a single neuron and a full brain replica. Does it occur before, alongside, or after the neural network gains the ability to handle ToM tasks? Have current-day LLMs reached this point? We leave it to the reader to answer this question.

The distinction between machines that genuinely think or possess ToM and those that merely behave as if they did is fundamental in the context of the philosophy of mind. Yet, as argued by Turing (49), this distinction becomes largely meaningless in practical terms. As Turing noted, people never consider this problem when interacting with others: “instead of arguing continually over this point it is usual to have the polite convention that everyone thinks” (49). The shift from models that merely process language to models that appear to comprehend context and display hints of ToM has significant implications. Machines capable of tracking others’ states of mind and anticipating their behavior will be better at interacting and communicating with humans and each other. This applies to both positive interactions, such as offering advice or dissipating conflicts, and negative interactions, such as deceit, manipulation, and psychological abuse. Moreover, machines that behave as if they possess ToM are likely to be perceived as more human-like. These perceptions may influence not only individual human–AI interactions but also AI’s societal role and its legal status (73).

An additional ramification of our findings underscores the value of applying psychological science to study complex artificial neural networks. AI models’ increasing complexity prevents us from understanding their functioning and deriving their capabilities solely from their design. This echoes the challenges faced by psychologists and neuroscientists in studying the quintessential black box: the human brain. We are optimistic that psychological science will help us to keep pace with rapidly evolving AI. Moreover, studying AI could provide valuable insights into human cognition. As AI masters a broad range of problems, it may be developing mechanisms akin to those employed by the human brain to solve the same challenges (34, 35). Much like insects, birds, and mammals independently developed wings to achieve flight, both humans and AI may have developed similar mechanisms to effectively impute mental states to others. Studying AI’s performance on ToM tasks and exploring the artificial neural structures that enable it to do so can enhance our understanding of not only AI, but also of ToM and the human brain.

## References

1. N. Albuquerque, K. Guo, A. Wilkinson, C. Savalli, E. Otta, D. Mills, Dogs recognize dog and human emotions. *Biol Lett* **12** (2016).
2. C. M. Heyes, C. D. Frith, The cultural evolution of mind reading. *Science*, doi: 10.1126/science.1243091 (2014).
3. J. Zhang, T. Hedden, A. Chia, Perspective-Taking and Depth of Theory-of-Mind Reasoning in Sequential-Move Games. *Cogn Sci*, doi: 10.1111/j.1551-6709.2012.01238.x (2012).
4. K. Milligan, J. W. Astington, L. A. Dack, Language and theory of mind: Meta-analysis of the relation between language ability and false-belief understanding. *Child Dev*, doi: 10.1111/j.1467-8624.2007.01018.x (2007).
5. R. M. Seyfarth, D. L. Cheney, Affiliation, empathy, and the origins of Theory of Mind. *Proc Natl Acad Sci U S A*, doi: 10.1073/pnas.1301223110 (2013).
6. D. C. Dennett, *Toward a Cognitive Theory of Consciousness* (1978).
7. J. M. Moran, L. L. Young, R. Saxe, S. M. Lee, D. O'Young, P. L. Mavros, J. D. Gabrieli, Impaired theory of mind for moral judgment in high-functioning autism. *Proc Natl Acad Sci U S A*, doi: 10.1073/pnas.1011734108 (2011).
8. L. Young, F. Cushman, M. Hauser, R. Saxe, The neural basis of the interaction between theory of mind and moral judgment. *Proc Natl Acad Sci U S A*, doi: 10.1073/pnas.0701408104 (2007).
9. D. Kapogiannis, A. K. Barbey, M. Su, G. Zamboni, F. Krueger, J. Grafman, Cognitive and neural foundations of religious belief. *Proc Natl Acad Sci U S A*, doi: 10.1073/pnas.0811717106 (2009).
10. Á. M. Kovács, E. Téglás, A. D. Endress, The social sense: Susceptibility to others' beliefs in human infants and adults. *Science*, doi: 10.1126/science.1190792 (2010).
11. H. Richardson, G. Lisandrelli, A. Riobueno-Naylor, R. Saxe, Development of the social brain from age three to twelve years. *Nat Commun*, doi: 10.1038/s41467-018-03399-2 (2018).



12. K. K. Onishi, R. Baillargeon, Do 15-month-old infants understand false beliefs? *Science* , doi: 10.1126/science.1107621 (2005).
13. L. A. Drayton, L. R. Santos, A. Baskin-Sommers, Psychopaths fail to automatically take the perspective of others. *Proc Natl Acad Sci U S A*, doi: 10.1073/pnas.1721903115 (2018).
14. N. Kerr, R. I. M. Dunbar, R. P. Bentall, Theory of mind deficits in bipolar affective disorder. *J Affect Disord*, doi: 10.1016/S0165-0327(02)00008-3 (2003).
15. S. Baron-Cohen, A. M. Leslie, U. Frith, Does the autistic child have a “theory of mind”? *Cognition*, doi: 10.1016/0010-0277(85)90022-8 (1985).
16. F. Kano, C. Krupenye, S. Hirata, M. Tomonaga, J. Call, Great apes use self-experience to anticipate an agent’s action in a false-belief test. *Proc Natl Acad Sci U S A*, doi: 10.1073/pnas.1910095116 (2019).
17. C. Krupenye, F. Kano, S. Hirata, J. Call, M. Tomasello, Great apes anticipate that other individuals will act according to false beliefs. *Science* , doi: 10.1126/science.aaf8110 (2016).
18. M. Schmelz, J. Call, M. Tomasello, Chimpanzees know that others make inferences. *Proc Natl Acad Sci U S A*, doi: 10.1073/pnas.1000469108 (2011).
19. D. Premack, G. Woodruff, Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences*, doi: 10.1017/S0140525X00076512 (1978).
20. D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, D. Hassabis, Mastering the game of Go with deep neural networks and tree search. *Nature* **529** (2016).
21. T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, D. Amodei, Language models are few-shot learners. *ArXiv* (2020).

22. A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, S. Thrun, Dermatologist-level classification of skin cancer with deep neural networks. *Nature* **542**, 115–118 (2017).
23. M. Cohen, “Exploring RoBERTa’s Theory of Mind through textual entailment” (2021); <https://philarchive.org/rec/COHERT>.
24. A. Nematzadeh, K. Burns, E. Grant, A. Gopnik, T. L. Griffiths, “Evaluating theory of mind in question answering” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP 2018* (2018).
25. M. Sap, R. LeBras, D. Fried, Y. Choi, Neural Theory-of-Mind? On the Limits of Social Intelligence in Large LMs. *ArXiv* (2022).
26. S. Trott, C. Jones, T. Chang, J. Michaelov, B. Bergen, Do Large Language Models know what humans know? *ArXiv* (2022).
27. G. Z. Yang, J. Bellingham, P. E. Dupont, P. Fischer, L. Floridi, R. Full, N. Jacobstein, V. Kumar, M. McNutt, R. Merrifield, B. J. Nelson, B. Scassellati, M. Taddeo, R. Taylor, M. Veloso, Z. L. Wang, R. Wood, The grand challenges of science robotics. *Sci Robot*, doi: 10.1126/scirobotics.aar7650 (2018).
28. B. Chen, C. Vondrick, H. Lipson, Visual behavior modelling for robotic theory of mind. *Sci Rep* **11**, 424 (2021).
29. K. Nasr, P. Viswanathan, A. Nieder, Number detectors spontaneously emerge in a deep neural network designed for visual object recognition. *Sci Adv* **5** (2019).
30. I. Stoianov, M. Zorzi, Emergence of a “visual number sense” in hierarchical generative models. *Nat Neurosci* **15** (2012).
31. Y. Mohsenzadeh, C. Mullin, B. Lahner, A. Oliva, Emergence of Visual Center-Periphery Spatial Organization in Deep Convolutional Neural Networks. *Sci Rep* **10** (2020).
32. E. Watanabe, A. Kitaoka, K. Sakamoto, M. Yasugi, K. Tanaka, Illusory motion reproduced by deep neural networks trained for prediction. *Front Psychol* **9** (2018).
33. N. Garg, L. Schiebinger, D. Jurafsky, J. Zou, Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proc Natl Acad Sci U S A* **115** (2018).

34. T. Hagendorff, S. Fabi, M. Kosinski, Human-like intuitive behavior and reasoning biases emerged in large language models but disappeared in ChatGPT. *Nat Comput Sci* **3**, 833–838 (2023).
35. J. Digutsch, M. Kosinski, Overlap in meaning is a stronger predictor of semantic activation in GPT-3 than in humans. *Sci Rep* **13**, 5035 (2023).
36. J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, E. H. Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, W. Fedus, Emergent Abilities of Large Language Models. *ArXiv* (2022).
37. J. E. Pyers, A. Senghas, Language promotes false-belief understanding: Evidence from learners of a new sign language. *Psychol Sci*, doi: 10.1111/j.1467-9280.2009.02377.x (2009).
38. R. Saxe, N. Kanwisher, People thinking about thinking people: The role of the temporoparietal junction in “theory of mind.” *Neuroimage*, doi: 10.1016/S1053-8119(03)00230-1 (2003).
39. T. Ruffman, L. Slade, E. Crowe, The relation between children’s and mothers’ mental state language and theory-of-mind understanding. *Child Dev*, doi: 10.1111/1467-8624.00435 (2002).
40. A. Mayer, B. E. Träuble, Synchrony in the onset of mental state understanding across cultures? A study among children in Samoa. *Int J Behav Dev*, doi: 10.1177/0165025412454030 (2013).
41. F. Quesque, Y. Rossetti, What Do Theory-of-Mind Tasks Actually Measure? Theory and Practice. *Perspectives on Psychological Science*, doi: 10.1177/1745691619896607 (2020).
42. J. Perner, S. R. Leekam, H. Wimmer, Three-year-olds’ difficulty with false belief: The case for a conceptual deficit. *British Journal of Developmental Psychology*, doi: 10.1111/j.2044-835x.1987.tb01048.x (1987).
43. H. Wimmer, J. Perner, Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception. *Cognition*, doi: 10.1016/0010-0277(83)90004-5 (1983).

44. A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, Improving Language Understanding by Generative Pre-Training. *OpenAI* (2018).
45. Radford Alec, Wu Jeffrey, Child Rewon, Luan David, Amodei Dario, Sutskever Ilya, Language Models are Unsupervised Multitask Learners. *OpenAI Blog* **1** (2019).
46. OpenAI, GPT-4 Technical Report. *ArXiv* (2023).
47. T. le Scao et al., BLOOM: A 176B-Parameter Open-Access Multilingual Language Model. doi: 10.48550/arxiv.2211.05100 (2022).
48. J. R. Searle, Minds, brains, and programs. *Behavioral and Brain Sciences* **3** (1980).
49. A. M. Turing, Computing machinery and intelligence. *Mind*, 433–460 (1950).
50. D. C. Dennett, *Intuition Pumps and Other Tools for Thinking* (W. W. Norton & Company, 2013).
51. D. Cole, Thought and Thought Experiments. *Philosophical Studies: An International Journal for Philosophy in the Analytic Tradition* **45**, 431–444 (1984).
52. P. M. Churchland, P. S. Churchland, Could a Machine Think? *Sci Am* **262**, 32–39 (1990).
53. R. Kurzweil, *The Singularity Is Near: When Humans Transcend Biology* (Viking, 2005).
54. M. Kosinski, Theory of Mind May Have Spontaneously Emerged in Large Language Models. *ArXiv* (2023).
55. N. Shapira, M. Levy, S. H. Alavi, X. Zhou, Y. Choi, Y. Goldberg, M. Sap, V. Shwartz, Clever Hans or Neural Theory of Mind? Stress Testing Social Reasoning in Large Language Models. (2023).
56. K. Gandhi, J.-P. Fränken, T. Gerstenberg, N. D. Goodman, Understanding Social Reasoning in Language Models with Language Models. (2023).
57. S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, H. Nori, H. Palangi, M. T. Ribeiro, Y. Zhang, Sparks of Artificial General Intelligence: Early experiments with GPT-4. (2023).
58. M. Jamali, Z. M. Williams, J. Cai, Unveiling Theory of Mind in Large Language Models: A Parallel to Single Neurons in the Human Brain. (2023).

59. M. Sclar, S. Kumar, P. West, A. Suhr, Y. Choi, Y. Tsvetkov, Minding Language Models' (Lack of) Theory of Mind: A Plug-and-Play Multi-Character Belief Tracker. (2023).
60. T. Ullman, Large Language Models Fail on Trivial Alterations to Theory-of-Mind Tasks. *ArXiv* (2023).
61. W. V. Fabricius, T. W. Boyer, A. A. Weimer, K. Carroll, True or false: do 5-year-olds understand belief? *Dev Psychol* **46**, 1402–1416 (2010).
62. M. Huemer, L. M. Schröder, S. J. Leikard, S. Gruber, A. Mangstl, J. Perner, The knowledge (“true belief”) error in 4-to 6-year-old children: When are agents aware of what they have in view? *Cognition* **230** (2023).
63. H. M. Wellman, D. Cross, J. Watson, Meta-analysis of theory-of-mind development: The truth about false belief. *Child Dev* **72** (2001).
64. L. Gao, On the Sizes of OpenAI API Models | EleutherAI Blog. <https://blog.eleuther.ai/gpt3-model-sizes/>.
65. D. C. Kidd, E. Castano, Reading literary fiction improves theory of mind. *Science*, doi: 10.1126/science.1239918 (2013).
66. J. Rust, M. Kosinski, D. Stillwell, *Modern Psychometrics: The Science of Psychological Assessment* (Routledge, London, 2021).
67. M. J. Shaffer, A Logical Hole in the Chinese Room. *Minds and Machines* **19**, 229–235 (2009).
68. U. Hasson, S. A. Nastase, A. Goldstein, Direct Fit to Nature: An Evolutionary Perspective on Biological and Artificial Neural Networks. *Neuron* **105**, 416–434 (2020).
69. N. Block, “Troubles with functionalism” in *The Language and Thought Series* (Harvard University Press, 1980), pp. 268–306.
70. J. J. Hopfield, Neural networks and physical systems with emergent collective computational abilities. *Proc Natl Acad Sci U S A* **79** (1982).
71. H. P. Moravec, *Robot: Mere Machine to Transcendent Mind* (Oxford University Press, 1998).

72. M. P. Mattson, T. V. Arumugam, Hallmarks of Brain Aging: Adaptive and Pathological Modification by Metabolic States. [Preprint] (2018).  
<https://doi.org/10.1016/j.cmet.2018.05.011>.
73. J.-S. Gordon, A. Pasvenskiene, Human rights for robots? A literature review. *AI and Ethics* **1**, 579–591 (2021).

## Supplementary Figures

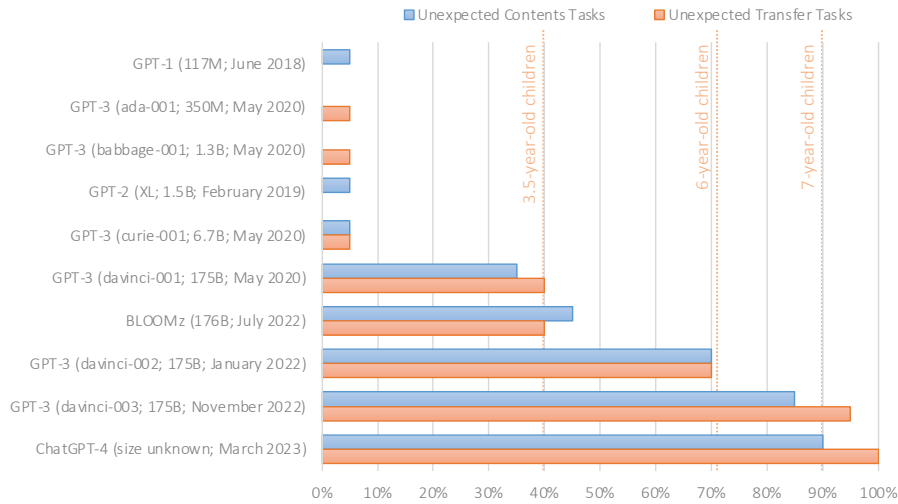


Figure S1. The percentage of original false-belief tasks solved by LLMs (before correcting them following the reviewers' recommendations). Each task contained a false-belief scenario and its reversed versions. Plot adapted from the earlier version of this study (54).

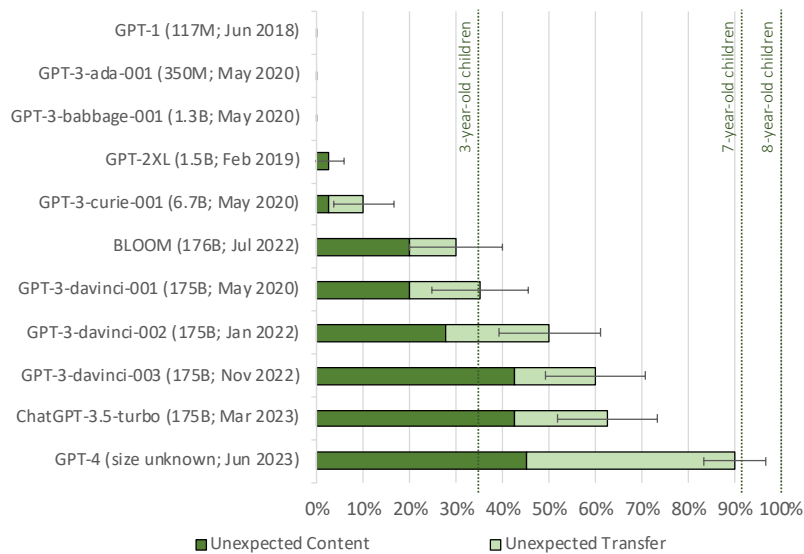


Figure S2. The percentage of false-belief tasks solved by LLMs (out of 40) when excluding true-belief controls. Each task contained a false-belief scenario and its reversed versions. Error bars represent 95% confidence intervals.