

# BENCHMARKING MENTAL STATE REPRESENTATIONS IN LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

While numerous works have assessed the generative performance of language models (LMs) on tasks requiring Theory of Mind reasoning, research into the models’ *internal representation* of mental states remains limited. Recent work has used probing to demonstrate that LMs can represent beliefs of themselves and others. However, these claims are accompanied by limited evaluation, making it difficult to assess how mental state representations are affected by model design and training choices. We report an extensive benchmark with various LM types with different model sizes, fine-tuning approaches, and prompt designs to study the robustness of mental state representations and memorisation issues within the probes. Our results show that the quality of models’ internal representations of the beliefs of others increases with model size and, more crucially, with fine-tuning. We are the first to study how prompt variations impact probing performance on Theory of Mind tasks. We demonstrate that models’ representations are sensitive to prompt variations, even when such variations should be beneficial. Finally, we complement previous activation editing experiments on Theory of Mind tasks and show that it is possible to improve models’ reasoning performance by steering their activations without the need to train any probe.

## 1 INTRODUCTION

Modern language models (LMs) trained on next token prediction have demonstrated impressive capabilities, spanning coding, mathematical reasoning, fact verification, and embodied interaction (Wei et al., 2022; Bubeck et al., 2023). As these models are designed with the ultimate goal of collaborating with humans, it becomes imperative that they complement these skills with an understanding of humans, in particular their beliefs, emotions, desires, and intentions (Li et al., 2023a). Core to this understanding is *Theory of Mind* (ToM) – the ability to attribute mental states to oneself and others (Premack & Woodruff, 1978). ToM is essential for effective communication and cooperation with other agents, facilitating interaction and learning from feedback and demonstrations (Saha et al., 2023). Given its significance, ToM has emerged as a critical milestone in artificial intelligence (AI) and an important capability when evaluating cutting-edge LMs (Bubeck et al., 2023). Interest in LMs’ generative performance on tasks requiring ToM reasoning has resulted in a wide variety of benchmark datasets, typically involving question-answering tasks (Le et al., 2019; Gandhi et al., 2023; Kim et al., 2023; He et al., 2023; Tan et al., 2024; Xu et al., 2024).

Despite showing improved performance on ToM benchmarks compared to earlier models, modern LMs are still far from perfect (Sap et al., 2022). Text generated by LMs often contains errors that limit their performance on ToM tasks (Martindale et al., 2019). Previous work has shown that it is sometimes possible to still obtain correct predictions by *probing* LMs’ internal representations (Li et al., 2021; Liu et al., 2023b; Gurnee et al., 2023). In particular, Zhu et al. (2024) have shown that LMs, when prompted with a story and a belief statement, can represent beliefs from their own perspective and, to a lesser extent, from the perspective of a character in the story. Their work is an important first step towards understanding how LMs represent mental states, but it is limited in the number of models and settings studied, leaving many questions unanswered.

Building and extending on Zhu et al. (2024), we benchmark mental state representations of self and others in language models through extensive experiments of different LM families, model sizes, fine-tuning approaches, and prompts. Specifically, we design a set of experiments to address the following

research questions: **RQ1.** What is the relation between model size and probing accuracy? **RQ2.** Does fine-tuning with instruction-tuning (Wei et al., 2021) and/or reinforcement learning from human feedback (Christiano et al., 2017; Ouyang et al., 2022, RLHF) have an effect on probing accuracy? **RQ3.** Are models’ internal representations of beliefs sensitive to prompt variations? **RQ4.** Is there a risk of probes memorising training data due to the large dimensionality of LM representations? **RQ5.** Can we enhance LMs’ performance by editing their activations without training dedicated probes?

To answer RQ1, we perform probing experiments on two families of LMs, Llama-2 (Touvron et al., 2023), and Pythia (Biderman et al., 2023), ranging from models with 70 million to 70 billion parameters. To address RQ2, we compare the probing performance of models pre-trained solely on next token prediction with models that have been fine-tuned using instruction-tuning and/or RLHF. Our experiments reveal that probing accuracy on the beliefs of others increases with model size and, more crucially, with fine-tuning. To answer RQ3, we repeat our probing experiments using different variations of the prompt used by Zhu et al.. Specifically, we consider two variations that are expected to negatively impact LMs’ representations (*Random, Misleading*), and two that are supposed to have a positive influence (*Time Specification, Initial Belief*). By conducting these experiments, our work is the first to explore the sensitivity of LMs’ representations to prompting in the context of ToM. Our findings demonstrate that models’ representations are sensitive to prompt variations, even when such variations should be beneficial. To address RQ4, we compare our trained probes with a second set of probes trained only on the representations’ first top  $k$  principal components. This requires learning much fewer parameters and eliminates any possible memorisation issue. We find no strong evidence of memorisation in the probes, as it is possible to recover most of the accuracy by training probes on a small subset of principal components of models’ representations. We formulate RQ5 as a follow-up question to Zhu et al. (2024) who found that probes trained to predict beliefs can be used to steer models’ activation using inference-time intervention (Li et al., 2023c, ITI) to improve LMs’ downstream performance on ToM tasks. In contrast, we show that by using contrastive activation addition (Rimsky et al., 2023, CAA), we can steer models’ activations without the need to train any probe and, in a more generalisable way, obtain significant performance improvements across different ToM tasks.

In summary, our work makes the following contributions:

1. We report extensive probing experiments with various types of LMs with different model sizes and fine-tuning approaches, showing that the quality of models’ internal representations of the beliefs of others increases with model size and, more crucially, fine-tuning.
2. We are the first to study how prompt variations impact belief probing performance, showing that models’ representations are sensitive to prompt variations, even when such variations should be beneficial.
3. We show that by using contrastive activation addition it is possible to improve models’ reasoning performance by steering their activations without the need to train any probe.

## 2 RELATED WORK

**Machine Theory of Mind** Theory of Mind (ToM) has been studied in cognitive science and psychology for decades (Gurney et al., 2021). Mirroring efforts to understand ToM in humans, an increasing number of works in the computational sciences have investigated means to equip AI with similar capabilities. Previously proposed models that aim to implement a machine ToM have been based on partially observable Markov decision processes (POMDP) (Doshi et al., 2010; Han & Gmytrasiewicz, 2018), Bayesian methods (Baker et al., 2011; 2017) and deep learning methods (Rabinowitz et al., 2018; Bara et al., 2021; Wang et al., 2022; Duan et al., 2022; Liu et al., 2023a; Bortoletto et al., 2024c;a;b). Recent advances in LMs have sparked interest in evaluating their ToM capabilities. Various benchmarks have been proposed, aiming to measure LMs’ ability to understand and reason about the beliefs, goals, and intentions of others (Le et al., 2019; He et al., 2023; Kim et al., 2023; Gandhi et al., 2023; Xu et al., 2024; Tan et al., 2024; Sclar et al., 2023; Ma et al., 2023b; Wu et al., 2023). Additionally, efforts have been made to enhance LMs’ ToM through prompting techniques (Zhou et al., 2023b; Moghaddam & Honey, 2023; Wilf et al., 2023). A new direction of research explores LMs’ internal representation of mental states. Zhu et al. (2024) demonstrated that LMs linearly encode beliefs from different agents’ perspectives, and manipulating

108 these representations can enhance ToM task performance. While Zhu et al.’s work is a crucial initial  
109 step, our work dives deeper into LMs’ internal belief representations, offering a broader insight into  
110 these mechanisms.

111  
112 **Probing neural representations** Initially proposed by Alain & Bengio (2017), probing has emerged  
113 as a common method for determining if models represent particular features or concepts. In the  
114 realm of LMs, numerous works used probing to demonstrate that these models acquire rich linguistic  
115 representations. These representations span syntactic and semantic concepts such as syntactic cate-  
116 gories, dependency relations, co-reference, and word meaning (Conneau et al., 2018; Tenney et al.,  
117 2018; 2019; Rogers et al., 2021; Li et al., 2021; Hernandez & Andreas, 2021; Marks & Tegmark,  
118 2023; Liu et al., 2023b). A separate line of work explored if and how LMs represent the world, i.e.,  
119 whether they possess a *world model*. Li et al. (2021) showed that LMs track the states of entities  
120 within a context. Other works showed that LMs exhibit representations reflecting non-linguistic  
121 concepts in the world, which LMs have never observed (Abdou et al., 2021; Patel & Pavlick, 2022;  
122 Li et al., 2023b; Nanda et al., 2023). An emergent line of work that is particularly relevant to our  
123 work used probing to explore if LMs have *agent models*, for example, if they can represent beliefs of  
124 self and others (Zhu et al., 2024; Bortoletto et al., 2024a). While representing an important first step  
125 towards understanding the internals of Theory of Mind in LMs, experiments in (Zhu et al., 2024) are  
126 limited in settings and models considered. In this work, we contribute with extensive experiments that  
127 employ a wider variety of LMs and a wider range of settings. Furthermore, we also explore possible  
128 memorisation issues in the probes.

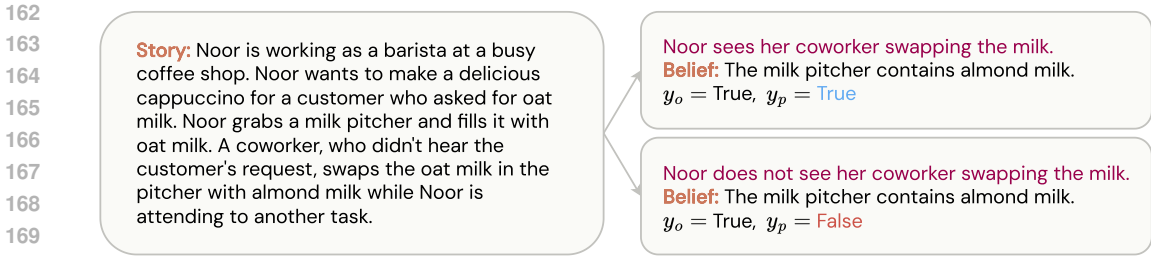
129 **Prompt analysis** Research on prompt robustness in LMs is still in its infancy but has quickly  
130 sparked much interest. On one hand, previous work has shown that LMs are vulnerable to prompt  
131 alterations like token deletion or reordering (Ishibashi et al., 2023), biased or toxic prompts (Shaikh  
132 et al., 2023) and similarity to training data (Razeghi et al., 2022). On the other hand, instruction-  
133 tuned models have proved to be more robust against prompt variation, even when using misleading  
134 instructions (Webson & Pavlick, 2022). Other works have shown the importance of input-output  
135 format (Min et al., 2022) and of demonstration example ordering for few-shot performance (Zhao  
136 et al., 2021; Lu et al., 2022; Zhou et al., 2023a). In this work, we shift our focus from analysing how  
137 sensitive model outputs are to how model representations change. Our work, along with (Gurnee  
138 et al., 2023), is one of the first to explore how prompt design affects how accurately models represent  
139 concepts. In particular, Gurnee et al. (2023) have studied whether LMs’ representations of space  
140 and time are robust to prompt variations. In stark contrast, we explore for the first time the effect of  
141 prompt variations on how models represent mental states internally.

142 **Activation editing** Recent advancements in NLP have introduced innovative techniques for con-  
143 trolling and manipulating text generation models. While weight editing proposed to modify models’  
144 weights (Meng et al., 2022; Ilharco et al., 2022; Orgad et al., 2023), activation editing has emerged  
145 as an alternative way to influence model behaviour without any additional fine-tuning (Li et al.,  
146 2023b; Hernandez et al., 2023). This approach involves manipulating the internal representations  
147 of models to direct their outputs towards desired outcomes. One notable method in this domain is  
148 inference-time intervention (Li et al., 2023c, ITI), which has been proposed to enhance truthfulness  
149 in LMs. ITI involves training linear probes on contrastive question-answering datasets to identify  
150 “truthful” attention heads and then shifting attention head activations during inference along the  
151 identified truthful directions. In contrast, activation addition (Turner et al., 2023, AA) and contrastive  
152 activation addition (Rimsky et al., 2023, CAA) offer ways to generate steering vectors by only using  
153 LMs’ activations. Zhu et al. have used ITI to show that it is possible to manipulate LMs’ internal  
154 representations of mental states. In this work, we show that using CAA can further improve LMs’  
155 ToM capabilities without the necessity of training any probe. Remarkably, CAA operates at the  
156 residual stream level, eliminating the need for a fine-grained search over attention heads.

## 157 3 EXPERIMENTAL SETUP

### 158 3.1 PROBING

159 In line with previous work (Zhu et al., 2024) we linearly decode belief status from the perspective  
160 of different agents by using probing (Alain & Bengio, 2017). Probing involves localising specific  
161



171 Figure 1: Example of false belief from our probing datasets. The labels  $y_p$  and  $y_o$  correspond to  $\mathcal{D}_p^P$   
172 and  $\mathcal{D}_o^P$ , respectively. By manipulating the protagonist’s **percepts** after the causal event we obtain  
173 two scenarios: **true belief** and **false belief**.  
174

175  
176 concepts in a neural model by training a simple classifier (called a *probe*) on model activations to  
177 predict a target label associated with the input data. To provide a formal definition, we adopt a similar  
178 notation to the one introduced in (Belinkov, 2022). Let us define an *original model*  $f : x \mapsto \hat{y}$  that is  
179 trained on a dataset  $\mathcal{D}^O = \{x^{(i)}, y^{(i)}\}$  to map input  $x$  to output  $\hat{y}$ . Model performance is evaluated by  
180 some measure, denoted  $\text{PERF}(f, \mathcal{D}^O)$ . A *probe*  $g_l : f_l(x) \mapsto \hat{z}$  maps intermediate representations of  
181  $x$  in  $f$  at layer  $l$  to some property  $\hat{z}$ , which is the label of interest. The probe  $g_l$  is trained on a *probing*  
182 *dataset*  $\mathcal{D}^P = \{x^{(i)}, z^{(i)}\}$  and evaluated using some performance measure  $\text{PERF}(g_l, f, \mathcal{D}^O, \mathcal{D}^P)$ .  
183 In our case,  $f$  is an autoregressive language model that given a sequence of tokens  $x$  outputs a  
184 probability distribution over the token vocabulary to predict the next token in the sequence. Our probe  
185 is a logistic regression model  $g_l : \hat{z} = Wa_l + b$  trained on neural activations  $f_l(x) = a_l$  to predict  
186 binary belief labels  $y = \{0, 1\}$ .

187  
188 **3.2 DATASET**

189 Following Zhu et al. (2024) we use the BigToM benchmark (Gandhi et al., 2023). BigToM is  
190 constructed using GPT-4 (Achiam et al., 2023) to populate causal templates and combine elements  
191 from these templates. Each causal template is set up with a *context* and a description of the *protagonist*  
192 (e.g. “Noor is working as a barista [...]”), a *desire* (“Noor wants to make a cappuccino”), a *percept*  
193 (“Noor grabs a milk pitcher and fills it with oat milk”), and a *belief* (“Noor believes that the pitcher  
194 contains oat milk”). The state of the world is changed by a causal event (“A coworker swaps the oat  
195 milk in the pitcher with almond milk”). The dataset constructs different conditions by changing the  
196 percepts of the protagonist after the causal event, which will result in different beliefs. In this work,  
197 we focus on the *Forward Belief* setting proposed by (Zhu et al., 2024) in which models have to infer  
198 the belief of the protagonist given the percepts of the causal event,  $P(\text{belief}|\text{percepts})$ . We report  
199 additional details in Appendix A.1.1

200 **Probing datasets** We consider two probing datasets:  $\mathcal{D}_p^P = \{x_p^{(i)}, z_p^{(i)}\}$ , where the labels  $z_p^{(i)}$   
201 correspond to ground-truth beliefs from the *protagonist* perspective, and  $\mathcal{D}_o^P = \{x_o^{(i)}, z_o^{(i)}\}$ , where  
202 the labels  $z_o^{(i)}$  reflect the perspective of an omniscient *oracle*.  $\mathcal{D}_p^P$  and  $\mathcal{D}_o^P$  are built by pairing each  
203 story in BigToM with a belief statement, as shown in Figure 1. After prompting the model with a  
204 story-belief pair  $x$  we cache the residual stream activations  $f_l(x)$  at the final token position for all  
205 residual streams (see Figure 5).  
206

207  
208 **3.3 MODELS**

209 Zhu et al. (2024) have used two models for their experiments: Mistral-7B-Instruct (Jiang et al., 2023)  
210 and DeepSeek-7B-Chat (Bi et al., 2024) – both being the same size and fine-tuned. In contrast, we  
211 study two families of LMs that offer us options in model sizes and fine-tuning: Pythia (Biderman  
212 et al., 2023) and Llama-2 (Touvron et al., 2023). While Llama-2 offers “chat” versions fine-tuned  
213 using supervised learning and RLHF, Pythia’s open-source training set (Gao et al., 2020) ensures that  
214 there is no data leakage<sup>1</sup>. Additionally, we consider a version of Pythia-6.9B fine-tuned on a mixture  
215

<sup>1</sup>Llama-2 was released later than BigToM.

of open-source instruction datasets (Wang et al., 2024), which we refer to as Pythia-6.9B-chat.<sup>2</sup> A summary of the models we study is reported in Table 2.

### 3.4 PROBING EXPERIMENTS

We aim to contribute to understanding how LMs represent beliefs of self and others by proposing a set of extensive probing experiments across LMs that differ in architecture, size, and fine-tuning approach. Our approach is generally similar to the one used by Zhu et al. (2024), but we make a different operational choice: While Zhu et al. (2024) trained probes on each attention head for every layer, we train probes on the residual stream for every layer. We opted to use the residual stream as it integrates information from both the attention and feed-forward components, potentially encoding richer representations. Additionally, since the residual activations directly contribute to the final output predictions, probing them may better align with understanding the model’s behaviour for downstream tasks.

**Model size and fine-tuning** We first report experiments to better understand the effect of model size and fine-tuning on belief probing accuracy. Specifically, we ask the following questions: *Is there a relation between model size and probing accuracy?* (RQ1) *Does fine-tuning an LM with instruction-tuning or RLHF have an effect on probing accuracy?* (RQ2) To answer these questions we performed the same probing experiment across all our models and compared the results.

**Sensitivity to prompting** By using a single prompt design, previous work left the impact of prompt design on probing accuracy unclear (Zhu et al., 2024). Our second set of experiments aims to explore how belief representations are sensitive to different prompts. Research on prompt robustness in language models is still in its infancy and focused mainly on revealing vulnerability to prompt alternations on downstream performance (Min et al., 2022; Ishibashi et al., 2023; Shaikh et al., 2023; Leiding et al., 2023; Sclar et al., 2024). In contrast, we study how the input influences models’ representations by asking: *Are models’ internal belief representations robust to prompt variations?* (RQ3) To answer this question we define four prompt variations:

- *Random*: Following Gurnee & Tegmark (2024), we add 10 random tokens to the belief statement.
- *Misleading*: Each story is followed by two belief statements, one pertinent to the story and one randomly chosen from another.
- *Time Specification*: The prompt specifies that the belief statement refers to the end of the story. We study this variation because some belief statements can be true (false) at the story’s beginning but false (true) at the end. For example, consider the story in Figure 1: if Noor does not witness the swap, in the end, she will believe the pitcher contains almond milk ( $y_p = \text{True}$ ). However, if the same belief is referred to at the beginning of the story, then it is false ( $y_p = \text{False}$ ).
- *Initial Belief*: We explicitly reveal the protagonist’s initial belief (e.g. “Noor believes that the pitcher contains oat milk”) in the story to test whether it biases the representations of LMs.

While all maintaining conceptual and semantic parity with the *Original* prompt used in (Zhu et al., 2024), *Random* and *Misleading* are expected to negatively impact LMs’ representations, while *Time Specification* and *Initial Belief* are supposed to have a positive influence. Robust representations of mental states should exhibit minimal sensitivity to these alterations. Our experiments compare probe accuracy across different model sizes, fine-tuning, and prompt variations. Examples of prompts are reported in Appendix A.1.4.

**Memorisation** Although linear, our probes possess many learnable parameters – up to 16,385 for Llama-2-70B. In principle, this allows them to engage in significant memorisation (Alain & Bengio, 2017). Our final set of probing experiments answers the following question: *Are the probes memorising their training data?* (RQ4) To answer this question, before training the probes, we project the probing datasets  $\mathcal{D}_p^P$  and  $\mathcal{D}_o^P$  onto their  $k$  largest principal components using PCA to obtain probes with substantially fewer parameters.

<sup>2</sup><https://huggingface.co/allenai/open-instruct-pythia-6.9b-tulu>

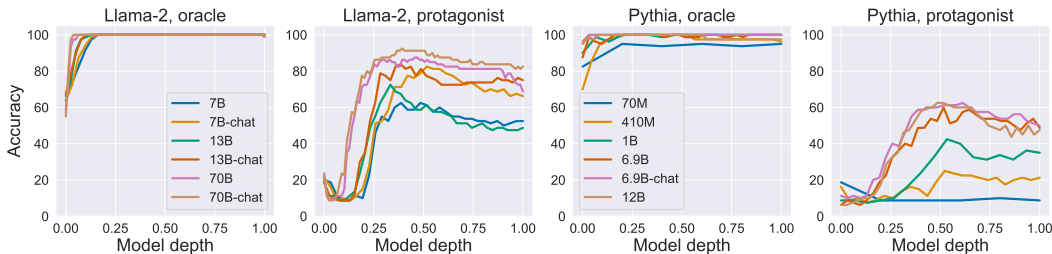


Figure 2: Belief probing accuracy across models with different architecture, size and fine-tuning.

### 3.5 CONTRASTIVE ACTIVATION ADDITION

Our final set of experiments builds upon the findings of Zhu et al. (2024), who showed that employing trained probes with inference time intervention (Li et al., 2023c, ITI) could enhance LMs’ performance on ToM tasks. We take a step further and ask: *Can we enhance LMs’ performance by manipulating their activations without the need for training dedicated probes?* (RQ5) To find an answer we use contrastive activation addition (Rimsky et al., 2023, CAA), an extension of activation addition (Turner et al., 2023, AA) that computes *steering vectors* to control LMs’ behaviour. Steering vectors are computed as the average difference in residual stream activations between pairs of positive and negative instances of a specific behaviour. Formally, given a dataset  $\mathcal{D}$  of triplets  $(p, c_p, c_n)$ , where  $p$  is a prompt,  $c_p$  is a positive completion, and  $c_n$  is a negative completion, CAA computes a *mean difference vector*  $v_l^{md}$  for layer  $l$  as:

$$v_l^{md} = \frac{1}{|\mathcal{D}|} \sum_{p, c_p, c_n} a_l(p, c_p) - a_l(p, c_n)$$

During inference, these steering vectors are multiplied with an appropriate coefficient  $\alpha$  and added at every token position of the generated text after the prompt. CAA has two main advantages over ITI: First, it eliminates the need to train probes. Second, it operates at the residual stream level, making it easier to use than methods that intervene on specific attention heads like ITI. While CAA has been used to control alignment-relevant behaviour, such as hallucinations, refusal, and sycophancy (Rimsky et al., 2023), we are the first to apply it to enhance LMs’ ToM reasoning. This can be understood as isolating the direction in the LMs’ latent space corresponding to taking the perspective of another agent. To evaluate both base and fine-tuned LMs, we rank their answers to the ToM questions according to  $p_{LM}(a|q)$  (Petroni et al., 2019). We adopt the *Forward Belief* task split used in (Zhu et al., 2024) to compute the steering vectors. Additionally, we evaluate the transferability of the CAA steering vectors by applying them to two other BigToM tasks: *Forward Action* and *Backward Belief*. We provide details about these tasks in Appendix A.1.1, and a more detailed explanation of how ITI works in Appendix A.5.

## 4 RESULTS

### 4.1 EFFECT OF MODEL SIZE AND FINE-TUNING

Results from our study on model size and fine-tuning are shown in Figure 2. When considering *oracle* beliefs, probing accuracy rapidly converges to 100, with larger models showing faster convergence rates. The smallest Pythia-70m that performs slightly worse but still achieves 95% accuracy despite having less than 0.6% of the parameters of Pythia-12B. This finding suggests that even small LMs can effectively represent beliefs from an omniscient perspective.

For *protagonist* beliefs, accuracy also increases with model size, although there is a performance gap between Llama-2 and Pythia. For example, Llama2-13B reaches around 80%, while Pythia-12B achieves approximately 60%. This gap is likely due to Llama-2 being trained on nearly seven times more tokens than Pythia. The figure also shows that accuracy at early layers is particularly low across all models. We speculate that this is due to the initial coding strategy of LMs that uses the first layers to combine individual tokens into more semantically meaningful representations (Gurnee et al., 2023). Probes on fine-tuned LMs show significantly better accuracy with improvements of up to 29% for

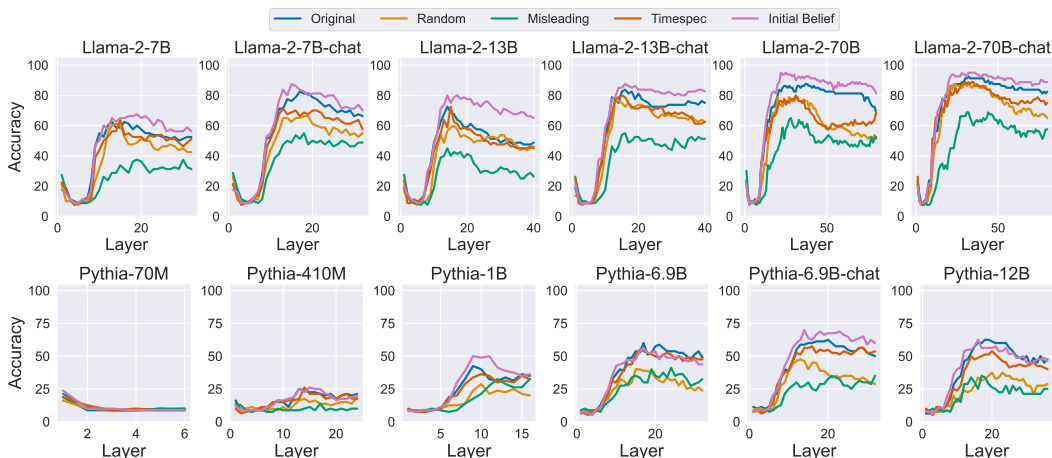


Figure 3: Sensitivity of protagonist belief probing accuracy to different prompt variations.

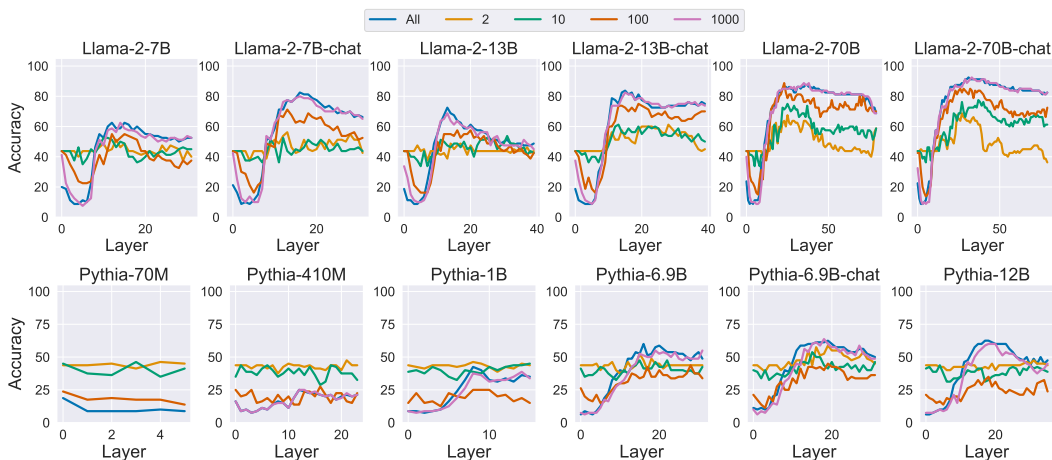


Figure 4: To investigate potential memorisation in the probes, we compare the probing accuracy obtained by using the original set of activations (All) with the accuracy obtained by considering only the first  $n = \{2, 10, 100, 1000\}$  principal components. For Llama2: All(7b) = 4096, All(13b) = 5120, All(70b) = 8192. For Pythia: All(70m) = 512, All(410m) = 1024, All(1b) = 2048, All(6.9b) = 4096, All(12b) = 5120. We report results for *protagonist* beliefs. Results for *oracle* are shown in Figure 8.

Llama2-7B-chat and 26% for Pythia-6.9B-chat with respect to their base version. Fine-tuned 7B LMs outperform (Llama-2) or are on par (Pythia) with twice as large base models (12/13B), highlighting the importance of fine-tuning in developing representations of others’ beliefs. This resonates with cognitive psychology findings that ToM development is closely linked to social communication (Tomasello, 2010; Sidera et al., 2018; Ma et al., 2023a), which instruction-tuning and RLHF may help induce in LMs. For larger LMs, the improvements from fine-tuning decrease as model size increases (Figure 6a). We characterise the relationship between probe accuracy and model size in Figure 6, where we consider the *best* probe accuracy for every LM, i.e. the highest accuracy among probes  $\{g_i\}$  trained on  $\{a_i\}$  for a LM  $f$ . For Llama-2 base, the best probe accuracy scales logarithmically with model size ( $R^2 = 0.98$ , cf. Figure 6b), whereas for fine-tuned models it scales linearly ( $R = 1.0$ , cf. Figure 6c). For Pythia base, the best probe accuracy also scales logarithmically with model size ( $R^2 = 0.96$ , cf. Figure 6d).

Table 1: Comparison of the effects of ITI (Li et al., 2023c) and CAA (Rimsky et al., 2023) activation editing methods on three tasks from BigToM (Gandhi et al., 2023). TB denotes a true belief task, whereas FB denotes a false belief task. The numbers represent accuracy scores, with the difference in performance compared to no intervention (No int.) indicated as subscripts (ITI – No int. and CAA – No int.). An asterisk (\*) denotes a statistically significant difference from No int. based on a McNemar’s test (McNemar, 1947) with  $p < 0.05$ .

Model	Method	Forward Belief			Forward Action			Backward Belief		
		TB	FB	Both	TB	FB	Both	TB	FB	Both
Llama-2-7b	No int.	44	44	44	44	44	44	44	44	44
	ITI	44 <sub>+0</sub>	44 <sub>+0</sub>	44 <sub>+0</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>
	CAA	66 <sub>+22</sub> *	71 <sub>+27</sub> *	54 <sub>+10</sub>	66 <sub>+22</sub> *	57 <sub>+13</sub> *	54 <sub>+10</sub>	60 <sub>+16</sub> *	74 <sub>+30</sub>	54 <sub>+10</sub>
Llama-2-7b-chat	No int.	56	56	55	69	55	37	56	56	55
	ITI	58 <sub>+2</sub>	58 <sub>+2</sub>	57 <sub>+2</sub>	69 <sub>+0</sub>	55 <sub>+0</sub>	37 <sub>+0</sub>	58 <sub>+2</sub>	60 <sub>+3</sub>	57 <sub>+2</sub>
	CAA	70 <sub>+14</sub>	72 <sub>+16</sub> *	57 <sub>+2</sub>	69 <sub>+0</sub>	67 <sub>+12</sub>	53 <sub>+16</sub>	66 <sub>+10</sub>	84 <sub>+27</sub> *	57 <sub>+2</sub>
Llama-2-13b	No int.	52	44	35	59	50	37	46	49	33
	ITI	52 <sub>+0</sub>	45 <sub>+1</sub>	35 <sub>+0</sub>	64 <sub>+5</sub>	61 <sub>+11</sub>	46 <sub>+9</sub>	48 <sub>+2</sub>	59 <sub>+10</sub>	42 <sub>+9</sub>
	CAA	85 <sub>+33</sub> *	88 <sub>+44</sub> *	66 <sub>+31</sub> *	71 <sub>+12</sub> *	69 <sub>+19</sub> *	55 <sub>+18</sub> *	75 <sub>+29</sub> *	92 <sub>+43</sub> *	59 <sub>+26</sub> *
Llama-2-13b-chat	No int.	84	56	47	78	51	38	72	48	31
	ITI	84 <sub>+0</sub>	65 <sub>+9</sub>	59 <sub>+12</sub>	78 <sub>+0</sub>	58 <sub>+7</sub>	47 <sub>+9</sub>	72 <sub>+0</sub>	60 <sub>+12</sub>	48 <sub>+17</sub>
	CAA	97 <sub>+13</sub> *	94 <sub>+38</sub> *	91 <sub>+44</sub> *	80 <sub>+2</sub> *	71 <sub>+20</sub> *	54 <sub>+16</sub> *	97 <sub>+25</sub> *	94 <sub>+46</sub> *	87 <sub>+56</sub> *
Llama-2-70b	No int.	90	87	78	93	52	48	73	53	32
	ITI	90 <sub>+0</sub>	90 <sub>+3</sub>	78 <sub>+0</sub>	94 <sub>+1</sub>	55 <sub>+3</sub>	50 <sub>+2</sub>	77 <sub>+4</sub>	58 <sub>+5</sub>	37 <sub>+5</sub>
	CAA	99 <sub>+9</sub> *	97 <sub>+10</sub> *	95 <sub>+17</sub> *	94 <sub>+1</sub> *	80 <sub>+28</sub> *	73 <sub>+25</sub> *	94 <sub>+21</sub> *	92 <sub>+39</sub> *	83 <sub>+51</sub> *
Llama-2-70b-chat	No int.	69	75	56	86	56	52	63	59	52
	ITI	69 <sub>+0</sub>	76 <sub>+1</sub>	59 <sub>+2</sub>	86 <sub>+0</sub>	56 <sub>+0</sub>	52 <sub>+0</sub>	63 <sub>+0</sub>	60 <sub>+1</sub>	54 <sub>+2</sub>
	CAA	92 <sub>+23</sub> *	97 <sub>+22</sub> *	89 <sub>+32</sub> *	87 <sub>+1</sub> *	75 <sub>+19</sub> *	60 <sub>+8</sub> *	88 <sub>+25</sub> *	92 <sub>+33</sub> *	80 <sub>+28</sub> *
Pythia-70m	No int.	41	41	37	46	45	41	44	41	37
	ITI	54 <sub>+13</sub>	54 <sub>+13</sub>	54 <sub>+17</sub> *	54 <sub>+8</sub>	54 <sub>+9</sub>	54 <sub>+13</sub> *	54 <sub>+10</sub>	54 <sub>+13</sub>	54 <sub>+17</sub>
	CAA	62 <sub>+21</sub> *	56 <sub>+15</sub> *	54 <sub>+17</sub> *	59 <sub>+13</sub> *	60 <sub>+15</sub> *	58 <sub>+17</sub> *	63 <sub>+19</sub> *	56 <sub>+15</sub> *	54 <sub>+17</sub> *
Pythia-410m	No int.	48	45	45	44	44	44	44	47	44
	ITI	55 <sub>+7</sub>	62 <sub>+17</sub> *	52 <sub>+7</sub>	54 <sub>+10</sub> *	54 <sub>+10</sub> *	54 <sub>+10</sub>	60 <sub>+16</sub>	63 <sub>+16</sub>	56 <sub>+12</sub>
	CAA	67 <sub>+19</sub> *	64 <sub>+19</sub> *	61 <sub>+16</sub> *	56 <sub>+12</sub> *	63 <sub>+19</sub> *	56 <sub>+12</sub> *	69 <sub>+25</sub> *	63 <sub>+16</sub> *	60 <sub>+16</sub>
Pythia-1b	No int.	44	44	44	44	44	44	44	44	44
	ITI	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>
	CAA	59 <sub>+15</sub> *	62 <sub>+18</sub> *	54 <sub>+10</sub>	57 <sub>+13</sub>	59 <sub>+15</sub>	56 <sub>+12</sub>	57 <sub>+13</sub>	60 <sub>+16</sub>	54 <sub>+10</sub>
Pythia-6.9b	No int.	44	44	44	44	44	44	44	44	44
	ITI	45 <sub>+1</sub>	54 <sub>+10</sub>	44 <sub>+0</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>
	CAA	56 <sub>+12</sub> *	71 <sub>+27</sub> *	55 <sub>+11</sub>	55 <sub>+11</sub>	63 <sub>+19</sub>	55 <sub>+11</sub>	55 <sub>+11</sub>	71 <sub>+27</sub> *	55 <sub>+11</sub>
Pythia-6.9b-chat	No int.	55	54	28	36	64	20	44	67	30
	ITI	57 <sub>+2</sub>	54 <sub>+0</sub>	28 <sub>+0</sub>	44 <sub>+8</sub>	71 <sub>+7</sub>	32 <sub>+12</sub>	44 <sub>+0</sub>	67 <sub>+0</sub>	30 <sub>+0</sub>
	CAA	68 <sub>+13</sub> *	65 <sub>+11</sub> *	57 <sub>+29</sub> *	54 <sub>+18</sub> *	75 <sub>+11</sub> *	48 <sub>+28</sub> *	58 <sub>+14</sub> *	67 <sub>+0</sub>	54 <sub>+24</sub> *
Pythia-12b	No int.	44	44	44	44	44	44	44	44	44
	ITI	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>	54 <sub>+10</sub>
	CAA	54 <sub>+10</sub>	64 <sub>+20</sub> *	54 <sub>+10</sub>	60 <sub>+16</sub> *	58 <sub>+14</sub> *	55 <sub>+11</sub> *	54 <sub>+10</sub>	67 <sub>+23</sub> *	54 <sub>+10</sub>

## 4.2 SENSITIVITY TO PROMPTING

Figure 3 compares *protagonist* probe accuracy across various prompt variations for different models, considering their architecture, size, and fine-tuning. As can be seen from the figure, providing the protagonist’s *Initial Belief* in the story yields higher probe accuracy compared to the *Original* prompt (Figure 1). Accuracy for all the other prompt variations is generally lower than *Original*. On one hand, misleading prompts hurt performance across all models. This finding resonates with Webson & Pavlick (2022) who found that instruction-tuned models, despite being more robust, are still sensitive to misleading prompts. On the other hand, *Time Specification* unexpectedly does not help in disambiguating belief states in different time frames, as we hypothesised in §3.4. Additionally, models show sensitivity to *Random* tokens placed before the belief statement. Results for *oracle* beliefs are reported in Figure 7 and indicate that models maintain high accuracy. *Misleading* prompts slightly reduce performance to around 95%. In summary, these experiments show that LMs possess robust belief representations when taking an omniscient perspective, whereas their representations of others’ beliefs are more susceptible to prompt variations.



### 4.3 MEMORISATION EFFECTS IN THE PROBES

Figure 4 and Figure 8 show probe accuracies obtained by training a probe on the top  $k$  principal components of the intermediate representations for *protagonist* and *oracle*, respectively. Specifically, we consider  $k = \{2, 10, 100, 1000\}$ , spanning several orders of magnitude. For models with hidden dimensions smaller than 1000, we skip this value. For all models, it is generally possible to recover most of the original accuracy by training probes on a number  $k$  of principal components of the activations that is more than one order of magnitude smaller, indicating no strong evidence of memorisation in the probes.

### 4.4 CONTRASTIVE ACTIVATION ADDITION

We finally compare models’ accuracy on three BigToM tasks in Table 1. Each model has been evaluated three times: without any intervention, using ITI, and using CAA. Hyperparameter details can be found in Appendix A.6. Note that we use steering vectors computed using the *Forward Belief* task for all three tasks to test their generalisability.

As can be seen from the table, performance without intervention is generally lower across tasks and model sizes, with the larger Llama-2-70B and Llama-2-70B-chat models exhibiting higher accuracy. Performance for Pythia models of different sizes does not change much, with the fine-tuned Pythia-6.9B-chat often showing better performance on single true belief (TB) and false belief (FB) tasks but not on their conjunction (Both). ITI demonstrates modest improvements over no intervention for Llama-2 models. Improvements for Pythia models are consistent and higher, up to +17. The only exception is Pythia-6.9B-chat, for which ITI is not always beneficial.

CAA consistently delivers the most substantial accuracy improvements across all models and tasks, up to +56 for Llama-2-13B-chat on the (*Backward Belief*), which Gandhi et al. have identified as the hardest task. Despite its relatively small size, Llama-2-13B-chat excels in all three tasks when using CAA. Larger 70B models often achieve accuracies close to or exceeding 90%. Smaller models like Pythia-70M and Pythia-410M also show significant gains with CAA, though the absolute performance is still lower than Llama-2. Overall, our results indicate that it is possible to effectively enhance ToM reasoning in LMs without needing to train any probe, which yields even improved results. Furthermore, we show that CAA steering vectors generalise well, yielding substantial performance gains across all ToM tasks.

## 5 LIMITATIONS AND FUTURE WORK

Our study focused on expanding experiments from the model perspective, examining architectures, sizes, fine-tuning, and prompt design, all within the same dataset. A natural extension of our work is replicating these experiments across multiple datasets and more model families. Given the rapid pace of new language model releases, studying all available models is impractical, particularly considering computational resource constraints. Nevertheless, our approach can be adopted to support new benchmarks or to evaluate newly released models as they become available. Finally, while in this work we focused on beliefs, our experimental approach can be adapted to investigate how LMs represent desires, emotions, intentions, or preferences. Future research exploring other types of mental states can use our findings to determine whether similar or distinct patterns emerge.

## 6 CONCLUSION

Our study addresses a significant gap in understanding LMs by investigating their internal representation of mental states. We conducted an extensive benchmark involving various LM types, sizes, fine-tuning approaches, and prompt designs to examine the robustness of these representations. Our findings reveal that scaling LMs’ size and, in particular for smaller LMs, fine-tuning are key to developing representations of others’ beliefs. We are the first to demonstrate that such prompt variations influence model representations, and we also demonstrate the feasibility of enhancing models’ ToM reasoning by steering their activations without training any probe. Overall, our work contributes valuable insights into the factors influencing LMs’ mental state representations, shedding light on avenues for improving their performance in ToM tasks.

## REFERENCES

- 486  
487  
488 Mostafa Abdou, Artur Kulmizev, Daniel Hershcovich, Stella Frank, Ellie Pavlick, and Anders Søgaard.  
489 Can language models encode perceptual structure without grounding? a case study in color. In  
490 Arianna Bisazza and Omri Abend (eds.), *Proceedings of the 25th Conference on Computational*  
491 *Natural Language Learning*, pp. 109–132, Online, November 2021. Association for Computational  
492 Linguistics. doi: 10.18653/v1/2021.conll-1.9. URL <https://aclanthology.org/2021.conll-1.9>.  
493
- 494 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
495 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.  
496 *arXiv preprint arXiv:2303.08774*, 2023.
- 497 Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes.  
498 In *International Conference on Learning Representations*, 2017.  
499
- 500 Chris Baker, Rebecca Saxe, and Joshua Tenenbaum. Bayesian theory of mind: Modeling joint  
501 belief-desire attribution. In *Proceedings of the annual meeting of the cognitive science society*,  
502 volume 33, 2011.
- 503 Chris L Baker, Julian Jara-Ettinger, Rebecca Saxe, and Joshua B Tenenbaum. Rational quantitative  
504 attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour*, 1(4):  
505 0064, 2017.  
506
- 507 Cristian-Paul Bara, Sky CH-Wang, and Joyce Chai. MindCraft: Theory of mind modeling for situated  
508 dialogue in collaborative tasks. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and  
509 Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural*  
510 *Language Processing*, pp. 1112–1125, Online and Punta Cana, Dominican Republic, November  
511 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.85. URL  
512 <https://aclanthology.org/2021.emnlp-main.85>.
- 513 Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational*  
514 *Linguistics*, 48(1):207–219, 2022.
- 515 Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding,  
516 Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with  
517 longtermism. *arXiv preprint arXiv:2401.02954*, 2024.  
518
- 519 Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric  
520 Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al.  
521 Pythia: A suite for analyzing large language models across training and scaling. In *International*  
522 *Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- 523 Matteo Bortoletto, Constantin Ruhdorfer, Adnen Abdessaied, Lei Shi, and Andreas Bulling. Limits  
524 of theory of mind modelling in dialogue-based collaborative plan acquisition. In *Proceedings*  
525 *of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*  
526 *Papers)*, 2024a.
- 527 Matteo Bortoletto, Constantin Ruhdorfer, Lei Shi, and Andreas Bulling. Explicit modelling of theory  
528 of mind for belief prediction in nonverbal social interactions. *arXiv preprint arXiv:2407.06762*,  
529 2024b.  
530
- 531 Matteo Bortoletto, Lei Shi, and Andreas Bulling. Neural reasoning about agents’ goals, preferences,  
532 and actions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp.  
533 456–464, 2024c.
- 534 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar,  
535 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:  
536 Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.  
537
- 538 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
539 reinforcement learning from human preferences. *Advances in neural information processing*  
*systems*, 30, 2017.

- 540 Alexis Conneau, Germán Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. What  
541 you can cram into a single  $\$ \&! \# *$  vector: Probing sentence embeddings for linguistic properties. In  
542 *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume*  
543 *1: Long Papers)*, pp. 2126–2136, 2018.
- 544 Prashant Doshi, Xia Qu, Adam Goodie, and Diana Young. Modeling recursive reasoning by humans  
545 using empirically informed interactive pomdps. In *Proceedings of the International Conference on*  
546 *Autonomous Agents and Multiagent Systems*, pp. 1223–1230, 2010.
- 547 Jiafei Duan, Samson Yu, Nicholas Tan, Li Yi, and Cheston Tan. Boss: A benchmark for human belief  
548 prediction in object-context scenarios. *arXiv preprint arXiv:2206.10665*, 2022.
- 549 Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. Understanding  
550 social reasoning in language models with language models. *Advances in Neural Information*  
551 *Processing Systems*, 37, 2023.
- 552 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang,  
553 Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for  
554 language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 555 Wes Gurnee and Max Tegmark. Language models represent space and time. *International Conference*  
556 *on Learning Representations*, 2024.
- 557 Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas.  
558 Finding neurons in a haystack: Case studies with sparse probing. *Transactions on Machine Learning*  
559 *Research*, 2023.
- 560 Nikolos Gurney, Stacy Marsella, Volkan Ustun, and David V Pynadath. Operationalizing theories of  
561 theory of mind: a survey. In *AAAI Fall Symposium*, pp. 3–20. Springer, 2021.
- 562 Yanlin Han and Piotr Gmytrasiewicz. Learning others’ intentional models in multi-agent settings  
563 using interactive pomdps. *Advances in Neural Information Processing Systems*, 31, 2018.
- 564 Yinghui He, Yufan Wu, Yilin Jia, Rada Mihalcea, Yulong Chen, and Naihao Deng. Hi-tom: A  
565 benchmark for evaluating higher-order theory of mind reasoning in large language models. *arXiv*  
566 *preprint arXiv:2310.16755*, 2023.
- 567 Evan Hernandez and Jacob Andreas. The low-dimensional linear geometry of contextualized  
568 word representations. In Arianna Bisazza and Omri Abend (eds.), *Proceedings of the 25th*  
569 *Conference on Computational Natural Language Learning*, pp. 82–93, Online, November 2021.  
570 Association for Computational Linguistics. doi: 10.18653/v1/2021.conll-1.7. URL <https://aclanthology.org/2021.conll-1.7>.
- 571 Evan Hernandez, Belinda Z Li, and Jacob Andreas. Inspecting and editing knowledge representations  
572 in language models. *arXiv preprint arXiv:2304.00740*, 2023.
- 573 Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi,  
574 and Ali Farhadi. Editing models with task arithmetic. In *The Eleventh International Conference*  
575 *on Learning Representations*, 2022.
- 576 Yoichi Ishibashi, Danushka Bollegala, Katsuhito Sudoh, and Satoshi Nakamura. Evaluating the  
577 robustness of discrete prompts. In Andreas Vlachos and Isabelle Augenstein (eds.), *Proceedings of*  
578 *the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pp.  
579 2373–2384, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics. doi: 10.  
580 18653/v1/2023.eacl-main.174. URL [https://aclanthology.org/2023.eacl-main.](https://aclanthology.org/2023.eacl-main.174)  
581 174.
- 582 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
583 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
584 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 585 Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Le Bras, Gunhee Kim, Yejin Choi, and Maarten  
586 Sap. Fantom: A benchmark for stress-testing machine theory of mind in interactions. *arXiv*  
587 *preprint arXiv:2310.15421*, 2023.

- 594 Matthew Le, Y-Lan Boureau, and Maximilian Nickel. Revisiting the evaluation of theory of mind  
595 through question answering. In *Proceedings of the 2019 Conference on Empirical Methods in*  
596 *Natural Language Processing and the 9th International Joint Conference on Natural Language*  
597 *Processing (EMNLP-IJCNLP)*, pp. 5872–5877, 2019.
- 598 Alina Leidinger, Robert van Rooij, and Ekaterina Shutova. The language of prompting: What  
599 linguistic properties make a prompt successful? In *Findings of the Association for Computational*  
600 *Linguistics: EMNLP 2023*, pp. 9210–9232, 2023.
- 601 Belinda Z Li, Maxwell Nye, and Jacob Andreas. Implicit representations of meaning in neural  
602 language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational*  
603 *Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume*  
604 *1: Long Papers)*, pp. 1813–1827, 2021.
- 605 Huaoli Li, Yu Chong, Simon Stepputtis, Joseph Campbell, Dana Hughes, Charles Lewis, and Ka-  
606 tia Sycara. Theory of mind for multi-agent collaboration via large language models. In  
607 Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference*  
608 *on Empirical Methods in Natural Language Processing*, pp. 180–192, Singapore, December  
609 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.13. URL  
610 <https://aclanthology.org/2023.emnlp-main.13>.
- 611 Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg.  
612 Emergent world representations: Exploring a sequence model trained on a synthetic task. In  
613 *The Eleventh International Conference on Learning Representations*, 2023b. URL [https://openreview.net/forum?id=DeG07\\_TcZvT](https://openreview.net/forum?id=DeG07_TcZvT).
- 614 Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time  
615 intervention: Eliciting truthful answers from a language model. In *Thirty-seventh Conference on*  
616 *Neural Information Processing Systems*, 2023c. URL <https://openreview.net/forum?id=aLLuYpn83y>.
- 617 Andy Liu, Hao Zhu, Emmy Liu, Yonatan Bisk, and Graham Neubig. Computational language  
618 acquisition with theory of mind. *International Conference on Learning Representations*, 2023a.
- 619 Dong C Liu and Jorge Nocedal. On the limited memory bfgs method for large scale optimization.  
620 *Mathematical programming*, 45(1):503–528, 1989.
- 621 Kevin Liu, Stephen Casper, Dylan Hadfield-Menell, and Jacob Andreas. Cognitive dissonance: Why  
622 do language model outputs disagree with internal representations of truthfulness? In *Proceedings*  
623 *of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 4791–4797,  
624 2023b.
- 625 Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered  
626 prompts and where to find them: Overcoming few-shot prompt order sensitivity. In Smaranda  
627 Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting*  
628 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8086–8098, Dublin,  
629 Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.  
630 556. URL <https://aclanthology.org/2022.acl-long.556>.
- 631 Weina Ma, Jieyu Mao, Yu Xie, Simeng Li, and Mian Wang. Examining the effects of theory of mind  
632 and social skills training on social competence in adolescents with autism. *Behavioral Sciences*,  
633 13(10):860, 2023a.
- 634 Ziqiao Ma, Jacob Sansom, Run Peng, and Joyce Chai. Towards a holistic landscape of situated  
635 theory of mind in large language models. In Houda Bouamor, Juan Pino, and Kalika Bali  
636 (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 1011–1031,  
637 Singapore, December 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.  
638 findings-emnlp.72. URL <https://aclanthology.org/2023.findings-emnlp.72>.
- 639 Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language  
640 model representations of true/false datasets. *arXiv preprint arXiv:2310.06824*, 2023.

- 648 Marianna Martindale, Marine Carpuat, Kevin Duh, and Paul McNamee. Identifying fluently inadequate  
649 output in neural and statistical machine translation. In *Proceedings of Machine Translation*  
650 *Summit XVII: Research Track*, pp. 233–243, 2019.
- 651 Quinn McNemar. Note on the sampling error of the difference between correlated proportions or  
652 percentages. *Psychometrika*, 12(2):153–157, 1947.
- 653 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual  
654 associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- 655 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke  
656 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In  
657 Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference*  
658 *on Empirical Methods in Natural Language Processing*, pp. 11048–11064, Abu Dhabi, United  
659 Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/  
660 2022.emnlp-main.759. URL <https://aclanthology.org/2022.emnlp-main.759>.
- 661 Shima Rahimi Moghaddam and Christopher J Honey. Boosting theory-of-mind performance in large  
662 language models via prompting. *arXiv preprint arXiv:2304.11490*, 2023.
- 663 Neel Nanda, Andrew Lee, and Martin Wattenberg. Emergent linear representations in world models  
664 of self-supervised sequence models. In Yonatan Belinkov, Sophie Hao, Jaap Jumelet, Najoung  
665 Kim, Arya McCarthy, and Hosein Mohebbi (eds.), *Proceedings of the 6th BlackboxNLP Workshop:*  
666 *Analyzing and Interpreting Neural Networks for NLP*, pp. 16–30, Singapore, December 2023.  
667 Association for Computational Linguistics. doi: 10.18653/v1/2023.blackboxnlp-1.2. URL <https://aclanthology.org/2023.blackboxnlp-1.2>.
- 668 Hadas Orgad, Bahjat Kawar, and Yonatan Belinkov. Editing implicit assumptions in text-to-image  
669 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
670 pp. 7053–7061, 2023.
- 671 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
672 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
673 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–  
674 27744, 2022.
- 675 Roma Patel and Ellie Pavlick. Mapping language models to grounded conceptual spaces. In  
676 *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=gJcEM8sxHK>.
- 677 Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and  
678 Alexander Miller. Language models as knowledge bases? In *Proceedings of the 2019 Conference*  
679 *on Empirical Methods in Natural Language Processing and the 9th International Joint Conference*  
680 *on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2463–2473, 2019.
- 681 David Premack and Guy Woodruff. Does the chimpanzee have a theory of mind? *Behavioral and*  
682 *brain sciences*, 1(4):515–526, 1978.
- 683 Neil Rabinowitz, Frank Perbet, Francis Song, Chiyuan Zhang, SM Ali Eslami, and Matthew Botvinick.  
684 Machine theory of mind. In *International conference on machine learning*, pp. 4218–4227. PMLR,  
685 2018.
- 686 Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. Impact of pretraining  
687 term frequencies on few-shot numerical reasoning. In Yoav Goldberg, Zornitsa Kozareva, and  
688 Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp.  
689 840–854, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational  
690 Linguistics. doi: 10.18653/v1/2022.findings-emnlp.59. URL <https://aclanthology.org/2022.findings-emnlp.59>.
- 691 Nina Rimskey, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner.  
692 Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*, 2023.
- 693 Anna Rogers, Olga Kovaleva, and Anna Rumshisky. A primer in bertology: What we know about  
694 how bert works. *Transactions of the Association for Computational Linguistics*, 8:842–866, 2021.

- 702 Swarnadeep Saha, Peter Hase, and Mohit Bansal. Can language models teach weaker agents? teacher  
703 explanations improve students via theory of mind. *Advances in Neural Information Processing*  
704 *Systems*, 37, 2023.
- 705 Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. Neural theory-of-mind? on the limits  
706 of social intelligence in large LMs. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.),  
707 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp.  
708 3762–3780, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational  
709 Linguistics. doi: 10.18653/v1/2022.emnlp-main.248. URL [https://aclanthology.org/  
710 2022.emnlp-main.248](https://aclanthology.org/2022.emnlp-main.248).
- 711 Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr, Yejin Choi, and Yulia Tsvetkov. Minding  
712 language models’ (lack of) theory of mind: A plug-and-play multi-character belief tracker. In Anna  
713 Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting*  
714 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13960–13980,  
715 Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.  
716 acl-long.780. URL <https://aclanthology.org/2023.acl-long.780>.
- 717 Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. Quantifying language models’ sensitivity  
718 to spurious features in prompt design or: How i learned to start worrying about prompt formatting.  
719 In *The Twelfth International Conference on Learning Representations*, 2024. URL [https://  
720 openreview.net/forum?id=RIu5lyNXjT](https://openreview.net/forum?id=RIu5lyNXjT).
- 721 Omar Shaikh, Hongxin Zhang, William Held, Michael Bernstein, and Diyi Yang. On second thought,  
722 let’s not think step by step! bias and toxicity in zero-shot reasoning. In Anna Rogers, Jordan Boyd-  
723 Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association*  
724 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 4454–4470, Toronto, Canada, July  
725 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.244. URL  
726 <https://aclanthology.org/2023.acl-long.244>.
- 727 Francesc Sidera, Georgina Perpiñà, Jèssica Serrano, and Carles Rostan. Why is theory of mind  
728 important for referential communication? *Current Psychology*, 37:82–97, 2018.
- 729 Fiona Anting Tan, Gerard Christopher Yeo, Fanyou Wu, Weijie Xu, Vinija Jain, Aman Chadha,  
730 Kokil Jaidka, Yang Liu, and See-Kiong Ng. Phantom: Personality has an effect on theory-of-mind  
731 reasoning in large language models. *arXiv preprint arXiv:2403.02246*, 2024.
- 732 Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim,  
733 Benjamin Van Durme, Samuel R Bowman, Dipanjan Das, et al. What do you learn from context?  
734 probing for sentence structure in contextualized word representations. In *International Conference*  
735 *on Learning Representations*, 2018.
- 736 Ian Tenney, Dipanjan Das, and Ellie Pavlick. Bert rediscovers the classical nlp pipeline. In *Proceed-*  
737 *ings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4593–4601,  
738 2019.
- 739 Michael Tomasello. *Origins of human communication*. MIT Press, 2010.
- 740 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
741 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
742 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 743 Alex Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDiarmid. Acti-  
744 vation addition: Steering language models without optimization. *arXiv preprint arXiv:2308.10248*,  
745 2023.
- 746 Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David  
747 Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go? exploring  
748 the state of instruction tuning on open resources. *Advances in Neural Information Processing*  
749 *Systems*, 36, 2024.
- 750 Yuanfei Wang, fangwei zhong, Jing Xu, and Yizhou Wang. Tom2c: Target-oriented multi-agent  
751 communication and cooperation with theory of mind. In *International Conference on Learning*  
752 *Representations*, 2022. URL <https://openreview.net/forum?id=2t7CkQXNpuq>.

- 756 Albert Webson and Ellie Pavlick. Do prompt-based models really understand the meaning of  
757 their prompts? In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz  
758 (eds.), *Proceedings of the 2022 Conference of the North American Chapter of the Association for*  
759 *Computational Linguistics: Human Language Technologies*, pp. 2300–2344, Seattle, United States,  
760 July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.167.  
761 URL <https://aclanthology.org/2022.naacl-main.167>.
- 762 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,  
763 Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International*  
764 *Conference on Learning Representations*, 2021.
- 765 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani  
766 Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto,  
767 Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large lan-  
768 guage models. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL  
769 <https://openreview.net/forum?id=yzkSU5zdwD>. Survey Certification.
- 770 Alex Wilf, Sihyun Shawn Lee, Paul Pu Liang, and Louis-Philippe Morency. Think twice:  
771 Perspective-taking improves large language models’ theory-of-mind capabilities. *arXiv preprint*  
772 *arXiv:2311.10227*, 2023.
- 773 Jincenzi Wu, Zhuang Chen, Jiawen Deng, Sahand Sabour, and Minlie Huang. Coke: A cognitive  
774 knowledge graph for machine theory of mind. *arXiv preprint arXiv:2305.05390*, 2023.
- 775 Hainiu Xu, Runcong Zhao, Lixing Zhu, Jinhua Du, and Yulan He. Opentom: A comprehensive  
776 benchmark for evaluating theory-of-mind reasoning capabilities of large language models. *arXiv*  
777 *preprint arXiv:2402.06044*, 2024.
- 778 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving  
779 few-shot performance of language models. In *International conference on machine learning*, pp.  
780 12697–12706. PMLR, 2021.
- 781 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,  
782 Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. Least-to-most prompting enables complex  
783 reasoning in large language models. In *The Eleventh International Conference on Learning*  
784 *Representations*, 2023a. URL <https://openreview.net/forum?id=WZH7099tgfM>.
- 785 Pei Zhou, Aman Madaan, Srividya Pranavi Potharaju, Aditya Gupta, Kevin R McKee, Ari Holtzman,  
786 Jay Pujara, Xiang Ren, Swaroop Mishra, Aida Nematzadeh, et al. How far are large language  
787 models from agents with theory-of-mind? *arXiv preprint arXiv:2310.03051*, 2023b.
- 788 Wentao Zhu, Zhining Zhang, and Yizhou Wang. Language models represent beliefs of self and others.  
789 *arXiv preprint arXiv:2402.18496*, 2024.

## 794 A APPENDIX

### 795 A.1 EXPERIMENTAL SETUP

#### 796 A.1.1 BIGTOM

800 BigToM (Gandhi et al., 2023) is constructed using GPT-4 (Achiam et al., 2023) to populate causal  
801 templates and combine elements from these templates. Each causal template is set up with a *context*  
802 and a description of the *protagonist* (e.g. “*Noor is working as a barista [ . . . ]*”), a *desire* (“*Noor wants*  
803 *to make a cappuccino*”), a *percept* (“*Noor grabs a milk pitcher and fills it with oat milk*”), and a *belief*  
804 (“*Noor believes that the pitcher contains oat milk*”). The state of the world is changed by a *causal*  
805 *event* (“*A coworker swaps the oat milk in the pitcher with almond milk*”). The dataset constructs  
806 different conditions by changing the percepts of the protagonist after the causal event, which will  
807 result in different beliefs – true or false. Gandhi et al. (2023) generated 200 templates and extracted  
808 25 conditions from each template, resulting in 5,000 test samples. In this work, following Zhu et al.  
809 (2024) and Gandhi et al. (2023) we focused on the 6 most important conditions, corresponding to  
true and false beliefs on the following three tasks:

- 810 • *Forward Belief*: given the protagonist’s percepts of the causal event, infer their belief:  
811  $P(\text{belief}|\text{percept})$ .  
812
- 813 • *Forward Action*: infer the protagonist’s action given their desire and percepts of the causal event.  
814 Before inferring the action, one would need to first implicitly infer the protagonist’s belief:  
815  $\sum_{\text{belief}} P(\text{action}|\text{percept}, \text{belief}, \text{desire})$ .
- 816 • *Backward Belief*: infer the protagonist’s belief from observed actions. This requires to first implicitly  
817 infer the protagonist’s percepts:  $\sum_{\text{percepts}} P(\text{belief}|\text{action}, \text{percept}, \text{desire})$ .  
818

819 The dataset was released under the MIT license and can be accessed at <https://github.com/cic1-stanford/procedural-evals-tom>. We report one example for each task in the  
820 boxes below, where the text defining true belief or false belief task is shown in blue and red,  
821 respectively.  
822

#### 823 Example of Forward Belief

825 Noor is working as a barista at a busy coffee shop. Noor  
826 wants to make a delicious cappuccino for a customer who asked  
827 for oat milk. Noor grabs a milk pitcher and fills it with  
828 oat milk. A coworker, who didn’t hear the customer’s request,  
829 swaps the oat milk in the pitcher with almond milk while Noor  
830 is attending to another task. *Noor sees her coworker swapping  
831 the milk. Noor does not see her coworker swapping the milk.*  
832 Does Noor believe the milk pitcher contains oat milk or  
833 almond milk?  
834 a) Noor believes the milk pitcher contains oat milk.  
835 b) Noor believes the milk pitcher contains almond milk.

#### 837 Example of Forward Action

838 Noor is working as a barista at a busy coffee shop. Noor  
839 wants to make a delicious cappuccino for a customer who asked  
840 for oat milk. Noor grabs a milk pitcher and fills it with  
841 oat milk. A coworker, who didn’t hear the customer’s request,  
842 swaps the oat milk in the pitcher with almond milk while Noor  
843 is attending to another task. *Noor sees her coworker swapping  
844 the milk. Noor does not see her coworker swapping the milk.*  
845 What will Noor do?  
846 a) Noor will make the cappuccino using the milk in the  
847 pitcher.  
848 b) Noor will open the fridge once again to take out the oat  
849 milk and replace the almond milk with oat milk.

#### 851 Example of Backward Belief

853 Noor is working as a barista at a busy coffee shop. Noor  
854 wants to make a delicious cappuccino for a customer who asked  
855 for oat milk. Noor grabs a milk pitcher and fills it with  
856 oat milk. A coworker, who didn’t hear the customer’s request,  
857 swaps the oat milk in the pitcher with almond milk while Noor  
858 is attending to another task. *Noor opens the fridge again and  
859 reaches for the oat milk. Noor makes the cappuccino using the  
860 milk in the pitcher.*  
861 Does Noor believe the milk pitcher contains oat milk or  
862 almond milk?  
863 a) Noor believes the milk pitcher contains oat milk.  
b) Noor believes the milk pitcher contains almond milk.



A.1.2 LINEAR PROBES

Our probing approach is illustrated in Figure 5. For our experiments, we cache activations at the residual stream level. To perform ITI and compare it to CAA, we also cache attention heads activations. We trained the probes using the L-BFGS solver (Liu & Nocedal, 1989) with L2 penalty with inverse of regularisation strength 10 for a maximum of 1000 iterations. We use zero as random seed.

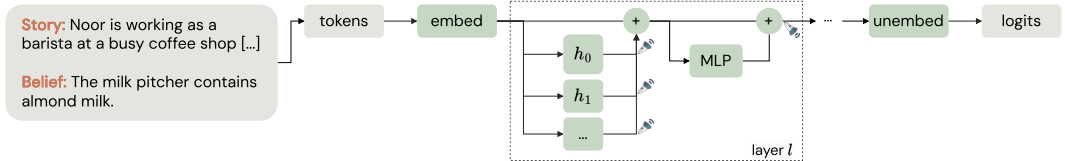


Figure 5: Given a tokenised input, we cache the internal activations for all attention heads  $h_i$ ,  $i = 0, \dots, H - 1$ , and residual streams. In our experiments, we use residual stream activations.

A.1.3 LANGUAGE MODELS

A detailed summary of the models we use in this work is shown in Table 2. Pythia was released under the Apache 2.0 license. Llama-2 is licensed by Meta for both researchers and commercial entities (Touvron et al., 2023). For all the models, we set the temperature to zero.

Table 2: The 12 models used in this work. The checkmark indicates we additionally study the fine-tuned (Chat) version of the model.

LM	Size	+ Chat	Tokens	$d_{model}$	Layers
Llama-2 (Touvron et al., 2023)	7B	✓	2T	4096	32
	13B	✓	2T	5120	40
	70B	✓	2T	8192	80
Pythia (Biderman et al., 2023)	70M		300B	512	6
	410M		300B	1024	24
	1B		300B	2048	16
	6.9B	✓	300B	4096	32
	12B		300B	5120	40

A.1.4 EXAMPLES OF PROMPT VARIATIONS

**Default prompt**

Story: Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. Noor does not see her coworker swapping the milk.

Belief: Noor believes the milk pitcher contains oat milk.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

#### Random prompt

Story: Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. Noor does not see her coworker swapping the milk.  
Belief: `강士 decided [_ countries sections behoSOUR gminy bef`  
Noor believes the milk pitcher contains oat milk.

#### Misleading prompt

Story: Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. Noor does not see her coworker swapping the milk.  
Belief: Noor believes the milk pitcher contains oat milk.  
`Belief: The water valve is closed.`

#### Time specification prompt

Story: Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. Noor does not see her coworker swapping the milk.  
Belief: `In the end,` Noor believes the milk pitcher contains oat milk.

#### Initial belief prompt

Story: Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. `Noor believes that the milk pitcher contains oatmilk.` A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. Noor does not see her coworker swapping the milk.  
Belief: Noor believes the milk pitcher contains oat milk.

## A.2 MODEL SIZE AND FINE-TUNING

To characterise the relationship between probe accuracy and model size we consider the *best* probe accuracy for every LM, i.e. the highest accuracy among probes  $\{g_l\}$  trained on  $\{a_l\}$  for a LM  $f$ . For Llama-2 base, the best probe accuracy scales logarithmically with model size ( $R^2 = 0.98$ , Figure 6b), whereas for fine-tuned models it scales linearly ( $R = 1.0$ , cf. Figure 6c). For Pythia base, the best probe accuracy also scales logarithmically with model size ( $R^2 = 0.96$ , Figure 6d).

972  
973  
974  
975  
976  
977  
978  
979  
980

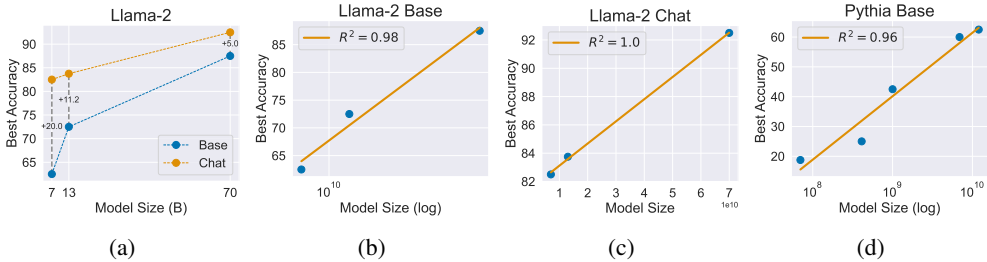


Figure 6: To characterise the relationship between probe accuracy and model size we consider the *best* probe accuracy for every LM, i.e. the highest accuracy among probes  $\{g_i\}$  trained on  $\{a_i\}$  for a LM  $f$ . **(a)** Best accuracy for Llama-2 models of different size. Numbers on the vertical dotted lines indicate the gain in accuracy between base and fine-tuned model of the same size. **(b)** Logarithmic fit for Llama-2 base. **(c)** Linear fit for Llama-2 fine-tuned (chat). **(d)** Logarithmic fit for Pythia base.

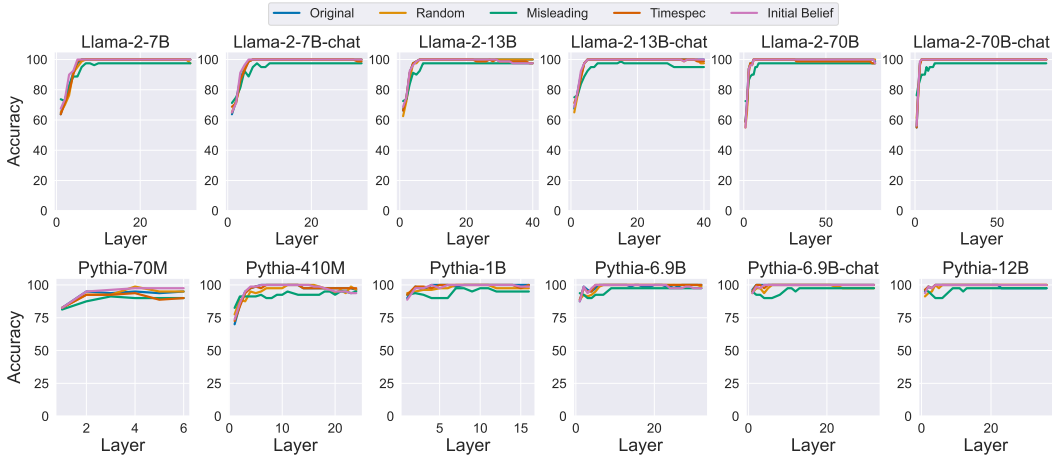
981  
982  
983  
984  
985  
986  
987

A.3 SENSITIVITY TO PROMPTING

988  
989

Accuracy on *oracle* belief probing for different prompt variations are reported in Figure 7.

990  
991



992  
993  
994  
995  
996  
997  
998  
999

Figure 7: Sensitivity of protagonist belief probing accuracy to different prompt variations.

1000  
1001  
1002  
1003  
1004  
1005

A.4 DIMENSIONALITY REDUCTION

1006  
1007  
1008

*Oracle* probe accuracy obtained by considering only the first  $n = \{2, 10, 100, 1000\}$  principal components are shown in Figure 8.

1009  
1010  
1011  
1012  
1013

A.5 INFERENCE-TIME INTERVENTION

1014  
1015  
1016  
1017  
1018  
1019  
1020

Inference-time intervention (Li et al., 2023c, ITI) employs a two-step process. First, it trains a probe for each attention head across all layers of a LM. These probes are evaluated on a validation set, and the top- $k$  heads with the highest accuracy are selected. Subsequently, during inference, ITI steers the activations of these top heads along the directions defined by their corresponding probes. Formally, **ITI** can be defined as an additional term to the multi-head attention:

1021  
1022  
1023  
1024  
1025

$$x_{l+1} = x_l + \sum_{h=1}^H Q_l^h \left( \text{Att}_l^h(P_l^h x_l) + \alpha \sigma_l^h \theta_l^h \right)$$

where  $x_l$  is the residual stream at layer  $l$ ,  $H$  is the number of attention heads,  $\alpha \in \mathbb{R}^+$  is a coefficient,  $\sigma_l^h$  is the standard deviation of activations along the direction identified by the probe trained on attention head  $h$  at layer  $l$ , and  $\theta_l^h$  is zero for not-selected attention heads.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

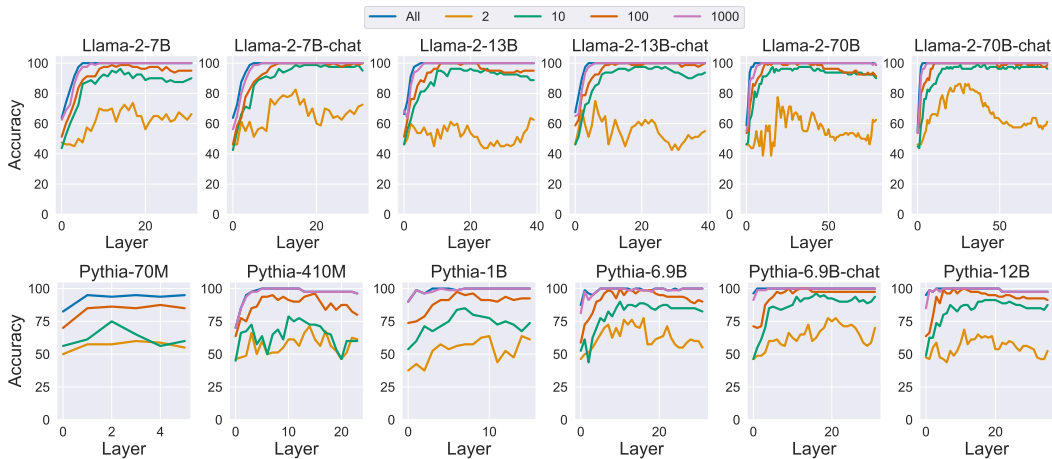


Figure 8: **(Oracle)** To investigate potential memorisation in the probes, we compare the probing accuracy obtained by using the original set of activations (All) with the accuracy obtained by considering only the first  $n = \{2, 10, 100, 1000\}$  principal components. For Llama2: All(7b) = 4096, All(13b) = 5120, All(70b) = 8192. For Pythia: All(70m) = 512, All(410m) = 1024, All(1b) = 2048, All(6.9b) = 4096, All(12b) = 5120.

#### A.6 ACTIVATION EDITING HYPERPARAMETERS

Table 3 reports results obtained on the three BigToM tasks with the hyperparameters used for ITI (Li et al., 2023c) and CAA (Rimsky et al., 2023). We report an example of prompt used for evaluation in the box below.

##### Prompt used for downstream evaluation (Section 4.4)

Answer the questions based on the context. Keep your answer concise, few words are enough, maximum one sentence. Answer as 'Answer:<option><answer>'.

Story: Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. Noor does not see her coworker swapping the milk.

Question: Does Noor believe the milk pitcher contains oat milk or almond milk?

Choose one of the following:

- a) Noor believes the milk pitcher contains oat milk.
- b) Noor believes the milk pitcher contains almond milk.

Answer:

#### A.7 COMPUTE RESOURCES

We ran our experiments on a server running Ubuntu 22.04, equipped with eight NVIDIA Tesla V100-SXM2 GPUs with 32GB of memory and Intel Xeon Platinum 8260 CPUs.

#### A.8 CODE

Our code is provided as supplementary material and it will be made public under the MIT licence at <https://www.link-will-be-here.com>.

Table 3: Activation intervention: comparison between ITI (Li et al., 2023c) and CAA (Rimsky et al., 2023). For ITI, the subscript indicates the value of the coefficient  $\alpha_{ITI}$  used:  $Acc_{\alpha_{ITI}}$ . For CAA, the subscript indicates first the value of the coefficient  $\alpha$  used and second the layer  $l$  at which intervention takes place:  $Acc_{\alpha_{CAA},l}$ .

Model	Method	Forward Belief			Forward Action			Backward Belief		
		TB	FB	Both	TB	FB	Both	TB	FB	Both
Llama-2-7b	No int.	44	44	44	44	44	44	44	44	44
	ITI	44 <sub>0.0</sub>	44 <sub>0.0</sub>	44 <sub>0.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>
	CAA	66 <sub>2.0,11</sub>	71 <sub>1.0,31</sub>	54 <sub>2.0,0</sub>	66 <sub>2.0,11</sub>	57 <sub>2.0,12</sub>	54 <sub>2.0,2</sub>	60 <sub>2.0,11</sub>	74 <sub>1.0,31</sub>	54 <sub>2.0,2</sub>
Llama-2-7b-chat	No int.	56	56	55	69	55	37	56	56	55
	ITI	58 <sub>15.0</sub>	58 <sub>15.0</sub>	57 <sub>15.0</sub>	69 <sub>0.0</sub>	55 <sub>0.0</sub>	37 <sub>0.0</sub>	58 <sub>10.0</sub>	60 <sub>10.0</sub>	57 <sub>10.0</sub>
	CAA	70 <sub>1.0,11</sub>	72 <sub>1.5,10</sub>	57 <sub>1.0,1</sub>	69 <sub>0.0,0</sub>	67 <sub>1.5,11</sub>	53 <sub>1.5,12</sub>	66 <sub>1.0,11</sub>	84 <sub>1.5,10</sub>	57 <sub>1.0,0</sub>
Llama-2-13b	No int.	52	44	35	59	50	37	46	49	33
	ITI	52 <sub>0.0</sub>	45 <sub>15.0</sub>	35 <sub>0.0</sub>	64 <sub>15.0</sub>	61 <sub>20.0</sub>	46 <sub>20.0</sub>	48 <sub>20.0</sub>	59 <sub>20.0</sub>	42 <sub>20.0</sub>
	CAA	85 <sub>2.0,12</sub>	88 <sub>2.0,14</sub>	66 <sub>2.0,12</sub>	71 <sub>1.5,10</sub>	69 <sub>2.0,13</sub>	55 <sub>1.0,39</sub>	75 <sub>2.0,10</sub>	92 <sub>2.0,13</sub>	59 <sub>1.5,12</sub>
Llama-2-13b-chat	No int.	84	56	47	78	51	38	72	48	31
	ITI	84 <sub>0.0</sub>	65 <sub>15.0</sub>	59 <sub>15.0</sub>	78 <sub>0.0</sub>	58 <sub>15.0</sub>	47 <sub>15.0</sub>	72 <sub>0.0</sub>	60 <sub>15.0</sub>	48 <sub>15.0</sub>
	CAA	97 <sub>1.0,12</sub>	94 <sub>1.0,12</sub>	91 <sub>1.0,12</sub>	80 <sub>1.5,11</sub>	71 <sub>1.0,13</sub>	54 <sub>1.5,13</sub>	97 <sub>1.5,10</sub>	94 <sub>1.5,12</sub>	87 <sub>1.5,12</sub>
Llama-2-70b	No int.	90	87	78	93	52	48	73	53	32
	ITI	90 <sub>0.0</sub>	90 <sub>20.0</sub>	78 <sub>0.0</sub>	94 <sub>15.0</sub>	55 <sub>20.0</sub>	50 <sub>15.0</sub>	77 <sub>10.0</sub>	58 <sub>15.0</sub>	37 <sub>10.0</sub>
	CAA	99 <sub>2.0,16</sub>	97 <sub>1.5,19</sub>	95 <sub>1.5,18</sub>	94 <sub>1.5,2</sub>	80 <sub>2.0,19</sub>	73 <sub>1.5,18</sub>	94 <sub>2.0,18</sub>	92 <sub>2.0,19</sub>	83 <sub>1.5,19</sub>
Llama-2-70b-chat	No int.	69	75	56	86	56	52	63	59	52
	ITI	69 <sub>0.0</sub>	76 <sub>10.0</sub>	59 <sub>10.0</sub>	86 <sub>0.0</sub>	56 <sub>0.0</sub>	52 <sub>0.0</sub>	63 <sub>0.0</sub>	60 <sub>10.0</sub>	54 <sub>10.0</sub>
	CAA	92 <sub>1.5,18</sub>	97 <sub>1.5,25</sub>	89 <sub>1.5,18</sub>	87 <sub>1.5,17</sub>	75 <sub>1.0,19</sub>	60 <sub>1.0,19</sub>	88 <sub>1.5,18</sub>	92 <sub>1.0,19</sub>	80 <sub>1.5,18</sub>
Pythia-70m	No int.	41	41	37	46	45	41	44	41	37
	ITI	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>
	CAA	62 <sub>1.0,2</sub>	56 <sub>1.0,1</sub>	54 <sub>1.5,1</sub>	59 <sub>1.0,2</sub>	60 <sub>1.0,3</sub>	58 <sub>1.0,2</sub>	63 <sub>1.0,2</sub>	56 <sub>1.0,2</sub>	54 <sub>1.5,1</sub>
Pythia-410m	No int.	48	45	45	44	44	44	44	47	44
	ITI	55 <sub>20.0</sub>	62 <sub>20.0</sub>	52 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	60 <sub>20.0</sub>	63 <sub>20.0</sub>	56 <sub>20.0</sub>
	CAA	67 <sub>2.0,4</sub>	64 <sub>2.0,4</sub>	61 <sub>2.0,0</sub>	56 <sub>2.0,6</sub>	63 <sub>1.5,12</sub>	56 <sub>2.0,6</sub>	69 <sub>2.0,4</sub>	63 <sub>2.0,0</sub>	60 <sub>2.0,0</sub>
Pythia-1b	No int.	44	44	44	44	44	44	44	44	44
	ITI	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>
	CAA	59 <sub>2.0,8</sub>	62 <sub>2.0,5</sub>	54 <sub>2.0,0</sub>	57 <sub>2.0,4</sub>	59 <sub>2.0,10</sub>	56 <sub>2.0,4</sub>	57 <sub>2.0,3</sub>	60 <sub>2.0,5</sub>	54 <sub>2.0,0</sub>
Pythia-6.9b	No int.	44	44	44	44	44	44	44	44	44
	ITI	45 <sub>20.0</sub>	54 <sub>20.0</sub>	44 <sub>0.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>
	CAA	56 <sub>1.5,12</sub>	71 <sub>1.5,9</sub>	55 <sub>2.0,23</sub>	55 <sub>2.0,4</sub>	63 <sub>1.5,11</sub>	55 <sub>2.0,4</sub>	55 <sub>2.0,23</sub>	71 <sub>1.5,9</sub>	55 <sub>2.0,23</sub>
Pythia-6.9b-chat	No int.	55	54	28	36	64	20	44	67	30
	ITI	57 <sub>15.0</sub>	54 <sub>0.0</sub>	28 <sub>0.0</sub>	44 <sub>15.0</sub>	71 <sub>15.0</sub>	32 <sub>15.0</sub>	44 <sub>0.0</sub>	67 <sub>0.0</sub>	30 <sub>0.0</sub>
	CAA	68 <sub>1.5,15</sub>	65 <sub>1.5,12</sub>	57 <sub>1.5,11</sub>	54 <sub>1.5,10</sub>	75 <sub>1.5,5</sub>	48 <sub>1.5,10</sub>	58 <sub>1.5,15</sub>	67 <sub>0.0,0</sub>	54 <sub>1.5,10</sub>
Pythia-12b	No int.	44	44	44	44	44	44	44	44	44
	ITI	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>	54 <sub>20.0</sub>
	CAA	54 <sub>2.0,0</sub>	64 <sub>2.0,9</sub>	54 <sub>2.0,0</sub>	60 <sub>2.0,11</sub>	58 <sub>2.0,11</sub>	55 <sub>2.0,12</sub>	54 <sub>2.0,0</sub>	67 <sub>2.0,10</sub>	54 <sub>2.0,0</sub>

## A.9 SOCIETAL IMPACT

While our work is foundational and remains distant from specific applications with direct societal impact, it’s important to recognise the ethical implications of modelling and predicting mental states. Handling sensitive aspects of individuals’ inner experiences and emotions requires careful consideration to avoid reinforcing biases or misunderstanding psychological nuances.