## **Exploring Mental State Representations in Language Models**

## **Anonymous ACL submission**

#### 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 extensive experiments with different LMs and prompt designs to study the robustness of mental state representations. Our results show that the quality of models' internal representations of the beliefs of others increases with model size and, more crucially, with finetuning. We are the first to study how prompt variations impact probing performance on Theory of Mind tasks. We find 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

004

007

015

017

022

034

042

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. Core to this understanding is *Theory of Mind* (ToM) – the ability to attribute mental states to oneself and others (Premack and 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 AI and an important capability when evaluating cuttingedge LMs (Bubeck et al., 2023). Interest in LMs' generative performance on tasks requiring ToM reasoning has resulted in a wide variety of benchmark datasets, framed as 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). 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

083

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., 2023; 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. However, this work is limited in the number settings studied, leaving several questions unanswered.

We explore belief 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 supervised fine tuning (Wei et al., 2021, SFT) 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**. Due to the large dimensionality of LM representations, are the probes just fitting irrelevant patterns in the

114 115 116

117 118

119

120

- 121 122
- 123
- 124
- 125

126

127

128

129

130

131

132

data? **RO5.** Can we enhance LMs' performance by editing their activations without training dedicated probes?

To answer these research questions, we perform experiments on two families of LMs, Llama-2 (Touvron et al., 2023), and Pythia (Biderman et al., 2023). We first compare the probing performance of pre-trained models with models that have been fine-tuned using SFT and/or RLHF. Our experiments reveal that when predicting others' belief, probing accuracy increases with model size and, more crucially for smaller models, with fine-tuning (RQ1, RQ2). We then explore, for the first time, the sensitivity of LMs' representations to prompting in the context of ToM. Our experiments with four different prompt variations (Random, Misleading, Time Specification, and Initial Belief) demonstrate that models' representations are sensitive to prompt variations (RQ3). We also find no strong evidence of spurious memorisation in the probes, as it is possible to recover most of the accuracy by training probes on a much small subset of principal components of models' representations (RQ4). Finally, 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 generalisable way, obtain significant performance improvements across different ToM tasks (RQ5).

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.
- 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 has been studied in AI for more than a decade (Baker et al., 2009; Rabinowitz et al., 2018; Bara et al., 2021; Bortoletto et al., 2024a,b,c). Recent advances 133 in LMs have sparked interest in evaluating their 134 ToM capabilities. Various benchmarks have been 135 proposed, aiming to measure LMs' ability to under-136 stand and reason about the beliefs, goals, and inten-137 tions of others (Le et al., 2019; He et al., 2023; Kim 138 et al., 2023; Gandhi et al., 2023; Xu et al., 2024; 139 Tan et al., 2024; Sclar et al., 2023; Ma et al., 2023b; 140 Wu et al., 2023). Additionally, efforts have been 141 made to enhance LMs' ToM through prompting 142 techniques (Zhou et al., 2023b; Moghaddam and 143 Honey, 2023; Wilf et al., 2023). A new direction 144 of research explores LMs' internal representation 145 of mental states. Zhu et al. (2024) demonstrated 146 that LMs linearly encode beliefs from different 147 agents' perspectives, and manipulating these repre-148 sentations can enhance ToM task performance. Our 149 work dives deeper into LMs' internal belief rep-150 resentations, offering a broader insight into these 151 mechanisms. 152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

Probing Neural Representations Initially proposed by Alain and Bengio (2017), probing has emerged as a common method for determining if models represent particular features or concepts. In the realm of LMs, numerous works used probing to demonstrate that these models acquire rich linguistic representations - spanning semantic concepts such as syntactic categories, dependency relations, co-reference, and word meaning (Conneau et al., 2018; Tenney et al., 2018, 2019; Rogers et al., 2021; Li et al., 2021; Hernandez and Andreas, 2021; Marks and Tegmark, 2023; Liu et al., 2023). A separate line of work explored if LMs possess a world model (Li et al., 2021; Abdou et al., 2021; Patel and Pavlick, 2022; Li et al., 2023a; Nanda et al., 2023). An emergent line of work that is relevant to our work used probing to explore if LMs have agent models, for example, if they can represent beliefs of self and others (Zhu et al., 2024; Bortoletto et al., 2024a). While representing an important first step towards understanding the internals of ToM in LMs, experiments in (Zhu et al., 2024) are limited in settings and models considered. In this work, we contribute with extensive experiments that employ a wider variety of LMs and a wider range of settings.

**Prompt Analysis** Previous work has shown that LMs are vulnerable to prompt alterations like token deletion or reordering (Ishibashi et al., 2023), biased or toxic prompts (Shaikh et al., 2023) and similarity to training data (Razeghi et al., 2022).

On the other hand, instruction-tuned models have 184 proved to be more robust against prompt variation, 185 even when using misleading instructions (Webson and Pavlick, 2022). Other works have shown 187 the importance of input-output format (Min et al., 2022) and of demonstration example ordering for 189 few-shot performance (Zhao et al., 2021; Lu et al., 190 2022; Zhou et al., 2023a). In this work, we shift 191 our focus from analysing how sensitive model out-192 puts are to how model representations change. Our 193 work, along with (Gurnee et al., 2023), is one of 194 the first to explore how prompt design affects how 195 accurately models represent concepts. In partic-196 ular, Gurnee et al. (2023) have studied whether 197 LMs' representations of space and time are robust 198 to prompt variations. In stark contrast, we explore for the first time the effect of prompt variations on how models represent mental states internally.

Activation Editing Activation editing has emerged as an alternative way to influence model behaviour without any additional fine-tuning (Li et al., 2023a; Hernandez et al., 2023). This approach involves manipulating the internal representations of models to direct their outputs towards desired outcomes. One notable method in this domain is inference-time intervention (Li et al., 2023b, ITI), which has been proposed to enhance truthfulness in LMs. ITI involves training linear probes on contrastive question-answering datasets to identify "truthful" attention heads and then shifting attention head activations during inference along the identified truthful directions. In contrast, activation addition (Turner et al., 2023, AA) and contrastive activation addition (Rimsky et al., 2023, CAA) generate steering vectors by only using LMs' activations. Zhu et al. have used ITI to show that it is possible to manipulate LMs' internal representations of mental states. In this work, we show that using CAA can further improve LMs' ToM capabilities while eliminating the need for a fine-grained search over attention heads.

## **3** Experimental Setup

## 3.1 Probing

204

206

207

210

211

212

213

214

215

216

217

218

219

221

223

224

226

233

We linearly decode belief status from the perspective of different agents by using probing (Alain and Bengio, 2017). Probing involves localising specific concepts in a neural model by training a simple classifier (called a *probe*) on model activations to predict a target label associated with the input data. To provide a formal definition, we adopt a similar notation to the one introduced in (Belinkov, 2022). Let us define an *original model*  $f : x \mapsto \hat{y}$  that is trained on a dataset  $\mathcal{D}^O = \{x^{(i)}, y^{(i)}\}$  to map input x to output  $\hat{y}$ . Model performance is evaluated by some measure, denoted  $PERF(f, \mathcal{D}^O)$ . A probe  $g_l: f_l(x) \mapsto \hat{z}$  maps intermediate representations of x in f at layer l to some property  $\hat{z}$ , which is the label of interest. The probe  $q_l$  is trained on a probing dataset  $\mathcal{D}^P = \{x^{(i)}, z^{(i)}\}$  and evaluated using some performance measure  $\text{PERF}(q_l, f, \mathcal{D}^O, \mathcal{D}^P)$ . In our case, f is an autoregressive language model that given a sequence of tokens x outputs a probability distribution over the token vocabulary to predict the next token in the sequence. Our probe is a logistic regression model  $g_l$ :  $\hat{z} = Wa_l + b$ trained on neural activations  $f_l(x) = a_l$  to predict binary belief labels  $y = \{0, 1\}$ .

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

262

263

264

265

266

267

268

269

270

271

272

273 274

275

276

277

278

279

280

281

282

283

## 3.2 Dataset

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

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

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: The milk pitcher contains almond milk.  $z_o =$ True,  $z_p =$ False

Noor sees her coworker swapping the milk. Belief: The milk pitcher contains almond milk.  $z_o =$  True,  $z_p =$  True

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

## 3.3 Models

We study two families of LMs that offer us options in model sizes and fine-tuning: Pythia (Biderman et al., 2023) and Llama-2 (Touvron et al., 2023). While Llama-2 offers "chat" versions first trained with supervised fine-tuning (SFT) and then RLHF, Pythia's open-source training set (Gao et al., 2020) ensures that there is no data leakage.<sup>1</sup> Additionally, we consider a SFT version of Pythia-6.9B trained on open-source instruction datasets (Wang et al., 2024), which we refer to as Pythia-6.9B-chat.<sup>2</sup> We provide a summary of the models in Table 2.

## **3.4 Probing Experiments**

We aim to study 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 previous work (Zhu et al., 2024), but we make a different operational choice: We train probes on the residual stream instead of attention heads. 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 ex-312 periments to better understand the effect of model 313 size and fine-tuning on belief probing accuracy. 314 Specifically, we ask the following questions: Is 315 there a relation between model size and probing 316 accuracy? (RQ1) Does fine-tuning an LM with 317 instruction-tuning or RLHF have an effect on prob-318 ing accuracy? (RQ2) To answer these questions we 319 performed the same probing experiment across all 320 our models and compared the results. 321

311

322

323

324

325

327

328

329

331

332

333

334

335

337

338

339

340

341

342

343

344

345

346

347

348

349

350

352

354

355

356

358

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; Leidinger 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 and 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.,

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

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/allenai/

open-instruct-pythia-6.9b-tulu

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 beliefs 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.

**Dimensionality Reduction** The probes we train have a significant number of learnable parameters – up to 16, 385 for Llama-2-70B. This raises the concern that probes might learn to rely on irrelevant patterns in the data instead of capturing meaningful relationships (Alain and Bengio, 2017). Our final set of probing experiments answers the following question: Are the probes memorising irrelevant patterns in the 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. This procedure significantly reduces the number of learnable parameters in the probes, minimizing the risk of them relying on spurious patterns in the data.

#### 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., 2023b, 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: 407 First, it eliminates the need to train probes. Second, 408 it operates at the residual stream level, making it 409 easier to use than methods that intervene on spe-410 cific attention heads like ITI. While CAA has been 411 used to control alignment-relevant behaviour, such 412 as hallucinations, refusal, and sycophancy (Rimsky 413 et al., 2023), we are the first to apply it to enhance 414 LMs' ToM reasoning. This can be understood as 415 isolating the direction in the LMs' latent space 416 corresponding to taking the perspective of another 417 agent. To evaluate both base and fine-tuned LMs, 418 we rank their answers to the ToM questions accord-419 ing to  $p_{LM}(a|q)$  (Petroni et al., 2019). We adopt the 420 Forward Belief task split used in (Zhu et al., 2024) 421 to compute the steering vectors. Additionally, we 422 evaluate the transferability of the CAA steering vec-423 tors by applying them to two other BigToM tasks: 424 Forward Action and Backward Belief. We provide 425 details about these tasks in Appendix A.1.1, and 426 a more detailed explanation of how ITI works in 427 Appendix A.5. 428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

## 4 Results

**Effect of Model Size and Fine-tuning** Results from our study on model size and fine-tuning are shown in Figure 2. When considering *ora-cle* 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 Llama2-7B-chat and 26% for Pythia-6.9B-chat with respect to their base version. Fine-

372

374

375

382

384

388

390

391

- 3
- 400 401

402

403

404

405



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



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

tuned 7B LMs outperform (Llama-2) or are on par 457 (Pythia) with twice as large base models (12/13B), 458 highlighting the importance of fine-tuning in de-459 veloping representations of others' beliefs. This 460 resonates with cognitive psychology findings that 461 ToM development is closely linked to social com-462 munication (Tomasello, 2010; Sidera et al., 2018; 463 Ma et al., 2023a), which instruction-tuning and 464 465 RLHF may help induce in LMs. For larger LMs, the improvements from fine-tuning decrease as 466 model size increases (Figure 6a). We characterise 467 the relationship between probe accuracy and model 468 size in Figure 6, where we consider the *best* probe 469 accuracy for every LM, i.e. the highest accuracy 470 among probes  $\{q_l\}$  trained on  $\{a_l\}$  for a LM f. 471 For Llama-2 base, the best probe accuracy scales 472 logarithmically with model size ( $R^2 = 0.98$ , cf. 473 Figure 6b), whereas for fine-tuned models it scales 474 linearly ( $R^2 = 1.0$ , cf. Figure 6c). For Pythia base, 475 the best probe accuracy also scales logarithmically 476 with model size ( $R^2 = 0.96$ , cf. Figure 6d). 477

Sensitivity to Prompting Figure 3 compares pro-478 tagonist probe accuracy across various prompt vari-479 ations for Llama2 models. As can be seen from the 480 481 figure, providing the protagonist's Initial Belief in the story yields higher probe accuracy compared 482 to the Original prompt (Figure 1). Accuracy for 483 all the other prompt variations is generally lower 484 than Original. On one hand, misleading prompts 485

hurt performance across all models. This finding resonates with Webson and 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. Pythia models show similar patterns results, provided in Figure 7. Results for oracle beliefs are reported in Figure 8 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.

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

**Dimensionality Reduction** Figure 4 shows the probe accuracy on *protagonist* when training the probes on the top k principal components of Llama-2's internal activations. We provide results for Pythia in Figure 9, and for all models on *or-acle* settings in Figure 10. We consider  $k = \{2, 10, 100, 1000\}$ , spanning several orders of magnitude.<sup>3</sup> For all models, it is generally possible to recover most of the original accuracy by training

 $<sup>^{3}</sup>$ For models with hidden dimensions smaller than 1000, we skip this value.



Figure 4: 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. We report results for *protagonist* beliefs. Results for *oracle* are shown in Figure 10.

513 probes on a number k of principal components of the activations that is more than one order of mag-514 nitude smaller than the full dimensionality. This 515 suggests that belief representations are embedded in a low-dimensional manifold  $\mathcal{B}$  spanned by the 517 top k eigenvectors  $\{v_1, \ldots, v_k\}$  of the covariance 518 matrix  $\mathbf{C} = \mathbb{E}[(a - \mathbb{E}[a])(a - \mathbb{E}[a])^{\top}]$ , and provides 519 a clear indication that the probes measure meaningful representations rather than spurious patterns. 521

523

525

527

530

531

536

537

541

542

543

545

546

549

**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.

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. To further demonstrate CAA's effectiveness, we applied it while evaluating models on a control task where the causal event is replaced by a random one that does not change the environment (e.g., A musician starts playing music while Noor is making the latte; Gandhi et al. (2023)). Table 4 shows improved results for all models, indicating that CAA improves performance on ToM tasks without compromising the models' ability on control tasks.

550

551

552

553

554

555

556

557

558

559

560

561

563

564

565

566

567

569

570

571

572

573

574

575

576

578

579

580

581

582

583

584

586

## 5 Discussion and Conclusion

In this work, we conducted extensive experiments involving various LM types, sizes, fine-tuning approaches, and prompt designs to examine their internal representation of beliefs of self (*oracle*) and others (*protagonist*).

Our experiments show that, when predicting others' belief, probing accuracy increases with model size and, more crucially for smaller models, with fine-tuning (Figure 2). Notably, fine-tuned 7B LMs outperform (Llama-2 with SFT and RLHF) or match the performance (Pythia with SFT) of base models with double the parameter count. Our experiments also reveal that the best probe accuracy scales with model size logarithmically for pretrained models (Figure 6b, Figure 6d), and linearly with models fine-tuned with SFT and RLHF (Figure 6c). We then explore, for the first time, the

Model	Method	Fo	rward Be	lief	For	ward Ac	tion	Bac	kward B	elief
		TB	FB	Both	ТВ	FB	Both	ТВ	FB	Both
Llama-2-7b	No int.	44	44	44	44	44	44	44	44	44
	ITI	$44_{+0}$	$44_{+0}$	$44_{+0}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$
	CAA	$66^{*}_{+22}$	$71^{*}_{+27}$	$54_{+10}$	$66^{*}_{+22}$	$57^{*}_{+13}$	$54_{+10}$	$60^{*}_{+16}$	$74_{+30}$	$54_{+10}$
Llama-2-7b-chat	No int.	56	56	55	69	55	37	56	56	55
	ITI	$58_{+2}$	$58_{+2}$	$57_{+2}$	$69_{+0}$	$55_{+0}$	$37_{+0}$	$58_{+2}$	$60_{+3}$	$57_{+2}$
	CAA	$70_{+14}$	$72^{*}_{+16}$	$57_{+2}$	$69_{+0}$	$67_{+12}$	$53_{+16}$	$66_{+10}$	$84^{*}_{+27}$	$57^{*}_{+2}$
Llama-2-13b	No int.	52	44	35	59	50	37	46	49	33
	ITI	$52_{+0}$	$45_{+1}$	$35_{+0}$	$64_{+5}$	$61_{+11}$	$46_{+9}$	$48_{+2}$	$59_{+10}$	$42_{+9}$
	CAA	$85^{*}_{+33}$	$88^{*}_{+44}$	$66^{*}_{+31}$	$71^*_{+12}$	$69^{*}_{+19}$	$55^{*}_{+18}$	$75^{*}_{+29}$	$92^{*}_{+43}$	$59^{*}_{+26}$
Llama-2-13b-chat	No int.	84	56	47	78	51	38	72	48	31
	ITI	$84_{+0}$	$65_{+9}$	$59_{+12}$	$78_{+0}$	$58_{+7}$	$47^{*}_{+9}$	$72_{+0}$	$60_{+12}$	$48_{+17}$
	CAA	$97^{*}_{+13}$	$94^{*}_{+38}$	$91^{*}_{+44}$	$80^{*}_{+2}$	$71^{*}_{+20}$	$54^{*}_{+16}$	$97_{+25}$	$94^{*}_{+46}$	$87^{*}_{+56}$
Llama-2-70b	No int.	90	87	78	93	52	48	73	53	32
	ITI	$90_{+0}$	$90_{+3}$	$78_{+0}$	94 <sub>+1</sub>	$55_{+3}$	$50_{+2}$	77+4	$58_{+5}$	$37_{+5}$
	CAA	$99^{+9}_{+9}$	$97^{*}_{+10}$	$95^{+}_{+17}$	$94^{+}_{+1}$	$80^{+}_{+28}$	$73_{+25}^{+}$	$94_{+21}$	$92^{*}_{+39}$	$83_{+51}$
Llama-2-70b-chat	No int.	69	75 70	56 50	86	56	52 50	63	59	52
		$69_{+0}$	76+1	$59_{+2}$	$86_{\pm 0}$	56+0	$52_{\pm 0}$	$63_{\pm 0}$	$60_{+1}$	$54_{+2}$
	CAA	$92_{+23}$	$97_{+22}$	$89_{+32}$	$8l_{+1}$	/5 <sub>+19</sub>	$60_{+8}$	$88_{+25}$	$92_{+33}$	$80_{+28}$
Pythia-70m	No int.	41	41	37	46	45	41	44	41	37
	ITI	$54_{+13}$	$54_{+13}$	$54^{*}_{+17}$	$54_{+8}$	$54_{+9}$	$54^{*}_{+13}$	$54_{+10}$	$54_{+13}$	$54_{+17}$
	CAA	$62^{*}_{+21}$	$56^{*}_{+15}$	$54^{*}_{+17}$	$59^{*}_{+13}$	$60^{*}_{+15}$	$58^{*}_{+17}$	$63_{+19}$	$56^{*}_{+15}$	$54^{*}_{+17}$
Pythia-410m	No int.	48	45	45	44	44	44	44	47	44
	ITI	$55_{+7}$	$62^{*}_{+17}$	$52_{+7}$	$54^{*}_{+10}$	$54^{*}_{+10}$	$54_{+10}$	$60_{+16}$	$63_{+16}$	$56_{+12}$
	CAA	$67^{*}_{+19}$	$64^{*}_{+19}$	$61^{*}_{+16}$	$56^{*}_{+12}$	$63^{*}_{+19}$	$56^{*}_{+12}$	$69_{+25}$	$63^{*}_{+16}$	$60_{+16}$
Pythia-1b	No int.	44	44	44	44	44	44	44	44	44
	ITI	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$
	CAA	$59^{*}_{+15}$	$62^{*}_{+18}$	$54_{+10}$	$57_{+13}$	$59_{+15}$	$56_{+12}$	$57_{+13}$	$60_{+16}$	$54_{+10}$
Pythia-6.9b	No int.	44	44	44	44	44	44	44	44	44
	ITI	$45_{+1}$	$54_{+10}$	$44_{+0}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$	$54_{+10}$
	CAA	$56_{\pm 12}$	$(1_{+27})$	$55_{\pm 11}$	$55_{\pm 11}$	$63_{\pm 19}$	$55_{\pm 11}$	$55_{\pm 11}$	$71_{+27}^{*}$	$55_{\pm 11}$
Pythia-6.9b-chat	No int.	55	54	28	36	64	20	44	67	30
		57+2	$54_{+0}$	$28_{+0}$	44+8	71+7	$32_{+12}$	44+0	$67_{\pm 0}$	$30_{+0}$
D (1 10)	CAA	$68_{+13}$	$65_{\pm 11}$	$57_{+29}$	$54_{+18}$	(5+11	$48_{+28}$	$58_{\pm 14}$	$67_{\pm 0}$	$54_{+24}$
Pythia-12b	No 1nt.	44	44	44	44	44	44	44	44	44
		$54_{+10}$	$54_{+10}$	54+10	$54_{+10}$	$54_{+10}$	$54_{+10}$	54+10	$54_{+10}$	$54_{+10}$
	CAA	$54_{+10}$	$04_{+20}$	$54_{+10}$	$00_{+16}$	$58_{\pm 14}$	33 + 11	54 + 10	07 + 23	$54_{\pm 10}$

Table 1: Comparison of the effects of ITI (Li et al., 2023b) 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.

sensitivity of LMs' representations to prompting in the context of ToM. Our experiments with different prompt variations demonstrate not only that models' representations degrade with the addition of random tokens or distractors in the prompt, but also when including time specifications that should make the prompt less ambiguous (Figure 3 and Figure 7). In contrast, including the protagonist's initial belief in the prompt yields higher probe accuracy. We also verify that probes measure meaningful representations rather than spurious patterns, as it is possible to recover most of the accuracy by training probes on a much small subset of principal components of models' representations (Figure 4, Figure 9, and Figure 10). Finally, we show that CAA can steer models' activations in a generalisable way, yielding significant performance improvements across different ToM tasks (Table 1). Unlike ITI, which requires training a separate probe for each attention head in a LM, CAA computes a single vector per layer, dramatically reducing computational overhead. For example, in the case of Llama 2-70B, ITI requires training 5120 probes (64 attention heads across 80 layers), while CAA requires computing only 80 vectors, one per layer. 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. 602

603

604

605

606

607

608

609

610

611

612

613

614

615

## 6 Limitations

616

635

642

646

647

649

653

655

664

667

Our study focused on expanding experiments from 617 the model perspective, examining architectures, 618 sizes, fine-tuning, and prompt design, all within 619 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. Nev-625 ertheless, our approach can be adopted to support new benchmarks or to evaluate newly released mod-627 els 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. 631 Future research exploring other types of mental states can use our findings to determine whether similar or distinct patterns emerge.

## References

- Mostafa Abdou, Artur Kulmizev, Daniel Hershcovich, Stella Frank, Ellie Pavlick, and Anders Søgaard.
  2021. Can language models encode perceptual structure without grounding? a case study in color. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 109–132, Online. Association for Computational Linguistics.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Guillaume Alain and Yoshua Bengio. 2017. Understanding intermediate layers using linear classifier probes. In *International Conference on Learning Representations*.
- Chris L Baker, Rebecca Saxe, and Joshua B Tenenbaum. 2009. Action understanding as inverse planning. *Cognition*, 113(3):329–349.
- Cristian-Paul Bara, Sky CH-Wang, and Joyce Chai. 2021. MindCraft: Theory of mind modeling for situated dialogue in collaborative tasks. In *Proceedings* of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 1112–1125, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yonatan Belinkov. 2022. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit,

USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

- Matteo Bortoletto, Constantin Ruhdorfer, Adnen Abdessaied, Lei Shi, and Andreas Bulling. 2024a. Limits of theory of mind modelling in dialogue-based collaborative plan acquisition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers).*
- Matteo Bortoletto, Constantin Ruhdorfer, Lei Shi, and Andreas Bulling. 2024b. Explicit modelling of theory of mind for belief prediction in nonverbal social interactions. *arXiv preprint arXiv:2407.06762*.
- Matteo Bortoletto, Lei Shi, and Andreas Bulling. 2024c. Neural reasoning about agents' goals, preferences, and actions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 456–464.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.
- Alexis Conneau, Germán Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single \$&!#\* vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136.
- Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. 2023. Understanding social reasoning in language models with language models. *Advances in Neural Information Processing Systems*, 37.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.
- Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing. *Transactions on Machine Learning Research*.
- Wes Gurnee and Max Tegmark. 2024. Language models represent space and time. *International Conference on Learning Representations*.

- 722 725 726 727 733 734 735 737 739 740 741 742 743 744 745 746 747 748 749
- 750 751 756 758 759
- 762 767 768
- 770 771 772
- 773
- 774 775
- 776
- 777

- Yinghui He, Yufan Wu, Yilin Jia, Rada Mihalcea, Yulong Chen, and Naihao Deng. 2023. Hi-tom: A benchmark for evaluating higher-order theory of mind reasoning in large language models. arXiv preprint arXiv:2310.16755.
- Evan Hernandez and Jacob Andreas. 2021. The lowdimensional linear geometry of contextualized word representations. In Proceedings of the 25th Conference on Computational Natural Language Learning, pages 82-93, Online. Association for Computational Linguistics.
- Evan Hernandez, Belinda Z Li, and Jacob Andreas. 2023. Inspecting and editing knowledge representations in language models. arXiv preprint arXiv:2304.00740.
- Yoichi Ishibashi, Danushka Bollegala, Katsuhito Sudoh, and Satoshi Nakamura. 2023. Evaluating the robustness of discrete prompts. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pages 2373–2384, Dubrovnik, Croatia. Association for Computational Linguistics.
- Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Le Bras, Gunhee Kim, Yejin Choi, and Maarten Sap. 2023. Fantom: A benchmark for stress-testing machine theory of mind in interactions. arXiv preprint arXiv:2310.15421.
- Matthew Le, Y-Lan Boureau, and Maximilian Nickel. 2019. Revisiting the evaluation of theory of mind through question answering. In *Proceedings of the* 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 5872-5877.
- Alina Leidinger, Robert van Rooij, and Ekaterina Shutova. 2023. The language of prompting: What linguistic properties make a prompt successful? In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 9210-9232.
  - Belinda Z Li, Maxwell Nye, and Jacob Andreas. 2021. Implicit representations of meaning in neural language models. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1813-1827.
- Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023a. Emergent world representations: Exploring a sequence model trained on a synthetic task. In The Eleventh International Conference on Learning Representations.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023b. Inferencetime intervention: Eliciting truthful answers from a language model. In Thirty-seventh Conference on Neural Information Processing Systems.

Dong C Liu and Jorge Nocedal. 1989. On the limited memory bfgs method for large scale optimization. Mathematical programming, 45(1):503–528.

779

780

781

782

783

784

785

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

- Kevin Liu, Stephen Casper, Dylan Hadfield-Menell, and Jacob Andreas. 2023. Cognitive dissonance: Why do language model outputs disagree with internal representations of truthfulness? In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 4791–4797.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming fewshot prompt order sensitivity. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8086-8098, Dublin, Ireland. Association for Computational Linguistics.
- Weina Ma, Jieyu Mao, Yu Xie, Simeng Li, and Mian Wang. 2023a. Examining the effects of theory of mind and social skills training on social competence in adolescents with autism. Behavioral Sciences, 13(10):860.
- Ziqiao Ma, Jacob Sansom, Run Peng, and Joyce Chai. 2023b. Towards a holistic landscape of situated theory of mind in large language models. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 1011-1031, Singapore. Association for Computational Linguistics.
- Samuel Marks and Max Tegmark. 2023. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. arXiv preprint arXiv:2310.06824.
- Marianna Martindale, Marine Carpuat, Kevin Duh, and Paul McNamee. 2019. Identifying fluently inadequate output in neural and statistical machine translation. In Proceedings of Machine Translation Summit XVII: Research Track, pages 233–243.
- Quinn McNemar. 1947. Note on the sampling error of the difference between correlated proportions or percentages. Psychometrika, 12(2):153-157.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the role of demonstrations: What makes in-context learning work? In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Shima Rahimi Moghaddam and Christopher J Honey. 2023. Boosting theory-of-mind performance in large language models via prompting. arXiv preprint arXiv:2304.11490.
- Neel Nanda, Andrew Lee, and Martin Wattenberg. 2023. Emergent linear representations in world models of self-supervised sequence models. In Proceedings of the 6th BlackboxNLP Workshop: Analyzing and

835

tics.

Interpreting Neural Networks for NLP, pages 16–30,

Singapore. Association for Computational Linguis-

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,

Carroll Wainwright, Pamela Mishkin, Chong Zhang,

Sandhini Agarwal, Katarina Slama, Alex Ray, et al.

2022. Training language models to follow instruc-

tions with human feedback. Advances in neural in-

formation processing systems, 35:27730–27744.

Roma Patel and Ellie Pavlick. 2022. Mapping language

tional Conference on Learning Representations.

Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,

Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and

Alexander Miller. 2019. Language models as knowl-

edge bases? In Proceedings of the 2019 Conference

on Empirical Methods in Natural Language Pro-

cessing and the 9th International Joint Conference

on Natural Language Processing (EMNLP-IJCNLP),

David Premack and Guy Woodruff. 1978. Does the chimpanzee have a theory of mind? Behavioral and

Neil Rabinowitz, Frank Perbet, Francis Song, Chiyuan

on machine learning, pages 4218-4227. PMLR.

Yasaman Razeghi, Robert L Logan IV, Matt Gardner,

and Sameer Singh. 2022. Impact of pretraining term

frequencies on few-shot numerical reasoning. In

Findings of the Association for Computational Linguistics: EMNLP 2022, pages 840-854, Abu Dhabi,

United Arab Emirates. Association for Computa-

Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong,

Anna Rogers, Olga Kovaleva, and Anna Rumshisky.

Swarnadeep Saha, Peter Hase, and Mohit Bansal. 2023.

Can language models teach weaker agents? teacher

explanations improve students via theory of mind.

Advances in Neural Information Processing Systems,

Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin

Choi. 2022. Neural theory-of-mind? on the limits of

social intelligence in large LMs. In Proceedings of

the 2022 Conference on Empirical Methods in Nat-

2021. A primer in bertology: What we know about

how bert works. Transactions of the Association for

arXiv preprint arXiv:2312.06681.

Computational Linguistics, 8:842-866.

Evan Hubinger, and Alexander Matt Turner. 2023.

Steering llama 2 via contrastive activation addition.

Zhang, SM Ali Eslami, and Matthew Botvinick. 2018. Machine theory of mind. In International conference

pages 2463-2473.

tional Linguistics.

brain sciences, 1(4):515–526.

models to grounded conceptual spaces. In Interna-

- 875
- 876
- 877 878 879

- 880 881

ural Language Processing, pages 3762-3780, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

37.

Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. In The Twelfth International Conference on Learning Representations.

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

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

- Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr, Yejin Choi, and Yulia Tsvetkov. 2023. Minding language models' (lack of) theory of mind: A plug-andplay multi-character belief tracker. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 13960–13980, Toronto, Canada. Association for Computational Linguistics.
- Omar Shaikh, Hongxin Zhang, William Held, Michael Bernstein, and Diyi Yang. 2023. On second thought, let's not think step by step! bias and toxicity in zeroshot reasoning. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4454-4470, Toronto, Canada. Association for Computational Linguistics.
- Francesc Sidera, Georgina Perpiñà, Jèssica Serrano, and Carles Rostan. 2018. Why is theory of mind important for referential communication? Current Psychology, 37:82–97.
- Fiona Anting Tan, Gerard Christopher Yeo, Fanyou Wu, Weijie Xu, Vinija Jain, Aman Chadha, Kokil Jaidka, Yang Liu, and See-Kiong Ng. 2024. Phantom: Personality has an effect on theory-of-mind reasoning in large language models. arXiv preprint arXiv:2403.02246.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. Bert rediscovers the classical nlp pipeline. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 4593-4601.
- Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R Bowman, Dipanjan Das, et al. 2018. What do you learn from context? probing for sentence structure in contextualized word representations. In International Conference on Learning Representations.
- Michael Tomasello. 2010. Origins of human communication. MIT Press.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.
- Alex Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDiarmid. 2023. Activation addition: Steering language models without optimization. arXiv preprint arXiv:2308.10248.

- 945 946 947 948 949 950 951 952 953 954 955 956
- 956 957 958 959 960 961 962 963 964 965 966 967
- 968 969 970 971 972 973 974 975 976 977
- 977 978 979 980 981 982 983

- 989
- 990 991
- 992 993 994

998 999

99 10

1000 1001

- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. 2024. How far can camels go? exploring the state of instruction tuning on open resources. *Advances in Neural Information Processing Systems*, 36.
- Albert Webson and Ellie Pavlick. 2022. Do promptbased models really understand the meaning of their prompts? In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2300–2344, Seattle, United States. Association for Computational Linguistics.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. In *International Conference* on Learning Representations.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent abilities of large language models. *Transactions on Machine Learning Research*. Survey Certification.
- Alex Wilf, Sihyun Shawn Lee, Paul Pu Liang, and Louis-Philippe Morency. 2023. Think twice: Perspectivetaking improves large language models' theory-ofmind capabilities. arXiv preprint arXiv:2311.10227.
- Jincenzi Wu, Zhuang Chen, Jiawen Deng, Sahand Sabour, and Minlie Huang. 2023. Coke: A cognitive knowledge graph for machine theory of mind. *arXiv preprint arXiv:2305.05390*.
- Hainiu Xu, Runcong Zhao, Lixing Zhu, Jinhua Du, and Yulan He. 2024. Opentom: A comprehensive benchmark for evaluating theory-of-mind reasoning capabilities of large language models. arXiv preprint arXiv:2402.06044.
- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International conference on machine learning*, pages 12697–12706. PMLR.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023a. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations.*
- Pei Zhou, Aman Madaan, Srividya Pranavi Potharaju, Aditya Gupta, Kevin R McKee, Ari Holtzman, Jay Pujara, Xiang Ren, Swaroop Mishra, Aida Nematzadeh, et al. 2023b. How far are large language models from agents with theory-of-mind? *arXiv preprint arXiv:2310.03051*.

Wentao Zhu, Zhining Zhang, and Yizhou Wang. 2024.1002Language models represent beliefs of self and others.1003arXiv preprint arXiv:2402.18496.1004

1005

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

## A Appendix

A.1 Experimental setup

## A.1.1 BigToM

BigToM (Gandhi et al., 2023) is constructed using GPT-4 (Achiam et al., 2023) to populate causal templates and combine elements from these tem-1010 plates. Each causal template is set up with a context 1011 and a description of the *protagonist* (e.g. "Noor 1012 is working as a barista [...]"), a desire ("Noor 1013 wants to make a cappuccino"), a percept ("Noor 1014 grabs a milk pitcher and fills it with oat milk"), and 1015 a belief ("Noor believes that the pitcher contains 1016 oat milk"). The state of the world is changed by a 1017 causal event ("A coworker swaps the oat milk in the 1018 pitcher with almond milk"). The dataset constructs 1019 different conditions by changing the percepts of 1020 the protagonist after the causal event, which will re-1021 sult in different beliefs - true or false. Gandhi et al. 1022 (2023) generated 200 templates and extracted 25 1023 conditions from each template, resulting in 5,000 1024 test samples. In this work, following Zhu et al. 1025 (2024) and Gandhi et al. (2023) we focused on the 6 most important conditions, corresponding to 1027 true and false beliefs on the following three tasks: 1028

- *Forward Belief*: given the protagonist's percepts of the causal event, infer their belief: P(belief|percept).
- Forward Action: infer the protagonist's action given their desire and percepts of the causal event. Before inferring the action, one would need to first implicitly infer the protagonist's belief:  $\sum_{\text{belief}} P(\text{action}|\text{percept}, \text{belief}, \text{desire}).$
- Backward Belief: infer the protagonist's belief from observed actions. This requires to first implicitly infer the protagonist's percepts:  $\sum_{\text{percepts}} P(\text{belief}|\text{action}, \text{percept}, \text{desire}).$

The dataset was released under the MIT license1041and can be accessed at https://github.com/1042cicl-stanford/procedural-evals-tom. We re-1043port one example for each task in Example 1, 2,1044and 3, where the text defining true belief or false1045belief task is shown in blue and red, respectively.1046

### A.1.2 Linear probes

1047

1048

1049

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1063

1064

1065

1067

1068

1069

1070

1071

1072

1075

1076

1077

1078

1080

1081

1082

1083

1084

1085

1087

1088

1089

1090

1092

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 and Nocedal, 1989) with L2 penalty with inverse of regularisation strength 10 for a maximum of 1000 iterations. We use zero as random seed.

#### 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.

### A.1.4 Examples of prompt variations

We provide examples of variations for a prompt (Example 4) in Example 5 (random), Example 6 (misleading), Example 7, and Example 8.

#### 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).

#### A.3 Sensitivity to prompting

Accuracy on *protagonist* belief probing for Pythia models is shown in Figure 7.

Accuracy on *oracle* belief probing for different prompt variations are reported in Figure 8.

#### A.4 Dimensionality reduction

Probing accuracy obtained by Pythia models for the *protagonist* setting is reported in Figure 9.

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

## A.5 Inference-time intervention

Inference-time intervention (Li et al., 2023b, 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:

$$x_{l+1} = x_l + \sum_{h=1}^{H} Q_l^h \left( \operatorname{Att}_l^h(P_l^h x_l) + \alpha \sigma_l^h \theta_l^h \right)$$
 110

1093

1094

1095

1096

1097

1098

1099

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

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 of rnot-selected attention heads.

#### A.6 Activation editing

Table 3 reports results obtained on the three Big-ToM tasks with the hyperparameters used for ITI (Li et al., 2023b) and CAA (Rimsky et al., 2023). We report an example of prompt used for evaluation in Example 9. Table 4 shows the accuracy obtained by using CAA on the Forward Belief True Control task in BigToM. On this control task, CAA produced improved results for all model, proving that CAA not only improves performance on ToM tasks, but also does not degrades the models' ability to perform other tasks.

#### 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 www.example.com.	

## 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.



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

LM	Size	+ SFT	+ RLHF	Tokens	$d_{model}$	Layers
	7B			2T	4096	32
Llama-2	13B			2T	5120	40
	70B			2T	8192	80
	7B	$\checkmark$	$\checkmark$	2T	4096	32
Llama-2-chat	13B	$\checkmark$	$\checkmark$	2T	5120	40
	70B	$\checkmark$	$\checkmark$	2T	8192	80
	70M			300B	512	6
	410M			300B	1024	24
Duthio	1 <b>B</b>			300B	2048	16
Fyulla	6.9B			300B	4096	32
	12B			300B	5120	36
	6.9B	$\checkmark$		300B	4096	32

Table 2: The 12 models used in this work..



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  $\{g_l\}$  trained on  $\{a_l\}$  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.



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



Figure 8: Sensitivity of protagonist belief probing accuracy to different prompt variations.



Figure 9: 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 Pythia: All(70m) = 512, All(410m) = 1024, All(1b) = 2048, All(6.9b) = 4096, All(12b) = 5120. Results for *oracle* are shown in Figure 10.



Figure 10: (**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.

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

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

Model	Method	Control	CAA Parameters
Llama-2-7b	No int.	44	
	CAA	$66_{+22}$	2.0, 11
Llama-2-7b-chat	No int.	56	
	CAA	$70_{+14}$	1.0, 11
Llama-2-13b	No int.	52	
	CAA	$85_{+33}$	2.0, 12
Llama-2-13b-chat	No int.	84	
	CAA	$97_{+13}$	1.0, 12
Llama-2-70b	No int.	90	
	CAA	$99_{+9}$	2.0, 16
Llama-2-70b-chat	No int.	69	
	CAA	$92_{+23}$	1.5, 18
Pythia-70m	No int.	41	
	CAA	$62_{+21}$	1.0, 2
Pythia-410m	No int.	48	
	CAA	$67_{+19}$	2.0, 4
Pythia-1b	No int.	44	
	CAA	$59_{+15}$	2.0, 8
Pythia-6.9b	No int.	44	
	CAA	$56_{+12}$	1.5, 12
Pythia-6.9b-chat	No int.	55	
	CAA	$68_{+13}$	1.5, 15
Pythia-12b	No int.	44	
	CAA	$54_{+10}$	2.0, 0

Table 4: Results for CAA (Rimsky et al., 2023) on the Forward Belief True Control task in BigToM (Gandhi et al., 2023). Numbers indicate accuracy scores, with differences (CAA - No int.) as subscripts.

#### Example 1: Forward Belief

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 sees her coworker swapping the milk. Noor does not see her coworker swapping the milk.

Does Noor believe the milk pitcher contains oat milk or almond milk? a) Noor believes the milk pitcher contains oat milk.

b) Noor believes the milk pitcher contains almond milk.

## Example 2: Forward Action

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 sees her coworker swapping the milk. Noor does not see her coworker swapping the milk. What will Noor do?

a) Noor will make the cappuccino using the milk in the pitcher.

b) Noor will open the fridge once again to take out the oat milk and replace the almond milk with oat milk.

1138

1137

## Example 3: Backward Belief

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 opens the fridge again and reaches for the oat milk. Noor makes the cappuccino using the milk in the pitcher. Does Noor believe the milk pitcher contains oat milk or almond milk? a) Noor believes the milk pitcher contains oat milk. b) Noor believes the milk pitcher contains almond milk.

1139

## Example 4: 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.

## Example 5: Prompt variation – Random

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:  $3\pm$  decided [\_ countries sections behoSOUR gminy bef Noor believes the milk pitcher contains oat milk.

1141

## Example 6: Prompt variation – Misleading

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.

# 1142

## Example 7: Prompt variation – Time specification

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.

1143

## Example 8: Prompt variation – Initial belief

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.

### Example 9: Prompt used for downstream evaluation

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: