

# Capturing Humans’ Mental Models of AI: An Item Response Theory Approach

Markelle Kelly  
kmarke@uci.edu

Department of Computer Science  
University of California, Irvine  
Irvine, California, USA

Padhraic Smyth  
smyth@ics.uci.edu

Department of Computer Science  
University of California, Irvine  
Irvine, California, USA

Aakriti Kumar  
aakritik@uci.edu

Department of Cognitive Sciences  
University of California, Irvine  
Irvine, California, USA

Mark Steyvers  
mark.steyvers@uci.edu

Department of Cognitive Sciences  
University of California, Irvine  
Irvine, California, USA

## ABSTRACT

Improving our understanding of how humans perceive AI teammates is an important foundation for our general understanding of human-AI teams. Extending relevant work from cognitive science, we propose a framework based on item response theory for modeling these perceptions. We apply this framework to real-world experiments, in which each participant works alongside another person or an AI agent in a question-answering setting, repeatedly assessing their teammate’s performance. Using this experimental data, we demonstrate the use of our framework for testing research questions about people’s perceptions of both AI agents and other people. We contrast mental models of AI teammates with those of human teammates as we characterize the dimensionality of these mental models, their development over time, and the influence of the participants’ own self-perception. Our results indicate that people expect AI agents’ performance to be significantly better on average than the performance of other humans, with less variation across different types of problems. We conclude with a discussion of the implications of these findings for human-AI interaction.

## CCS CONCEPTS

• **Human-centered computing** → *HCI theory, concepts and models*.

## KEYWORDS

theory of mind, mental models, human-AI interaction

### ACM Reference Format:

Markelle Kelly, Aakriti Kumar, Padhraic Smyth, and Mark Steyvers. 2023. Capturing Humans’ Mental Models of AI: An Item Response Theory Approach. In *2023 ACM Conference on Fairness, Accountability, and Transparency (FAccT ’23)*, June 12–15, 2023, Chicago, IL, USA. ACM, New York, NY, USA, 18 pages. <https://doi.org/10.1145/3593013.3594111>

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

FAccT ’23, June 12–15, 2023, Chicago, IL, USA

© 2023 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-0192-4/23/06.

<https://doi.org/10.1145/3593013.3594111>

## 1 INTRODUCTION

With the recent rapid growth in interest in applying AI approaches<sup>1</sup> to a wide variety of decision and prediction problems, there is an increasing realization that hybrid human-AI teams will be an important component of how AI will be deployed in practice [33, 40]. A decision making process that includes humans and AI can, ideally, benefit from the strengths of each [31, 34, 67]. Humans can act as a safeguard for unpredictable or undesirable behavior in AI algorithms, and can incorporate the type of contextual information and common sense reasoning that AI often lacks [17, 19, 20]. Conversely, the use of AI in prediction and decision making enables the processing of more complex patterns and greater volumes of data than humans alone can accommodate, for example in tedious and time-consuming work such as fact-checking [39] and in high-stakes decision making such as diagnostic medical imaging [64].

One important goal for human-AI systems is *complementarity*—achieving better performance than either the human(s) or the AI agent(s) acting independently [8, 22, 30]. A significant body of literature has developed that aims to understand and achieve complementarity across a variety of hybrid human-AI settings [11, 53, 62, 71]. One of the lessons learned from this work is that achieving complementarity is complex in practice; for instance, AI agents that exhibit high performance (e.g., in terms of prediction accuracy) on their own can actually harm overall team performance if their behavior is unpredictable for humans [5, 7]. Thus, understanding humans’ expectations of AI is essential for optimizing team performance, and recent work in this area has called for a better understanding of how humans perceive AI [40, 61].

In particular, an important (and somewhat under-studied) topic in this context is humans’ mental models of their AI teammates. If a person is deciding whether or not to take the advice of an AI agent, they are likely to make better decisions if they can accurately perceive its strengths and weaknesses, that is, if they have developed an accurate mental model of the agent [6]. For example, an AI agent could have particular “blindspots” in terms of its expertise, even when the agent’s overall performance on a particular task is comparable to or exceeds that of a human [4, 21]. In general, a

<sup>1</sup>We will use the term “AI” in this paper to refer to the broad spectrum of techniques that are currently (in 2023) referred to as “AI,” including models built using machine learning in particular.

better understanding of human mental models of AI agents can help predict phenomena such as humans’ development of appropriate trust in AI agents, human behavior in deferring to an AI agent, and how a human-AI team will function overall.

In this paper, we aim to improve our understanding of human mental models of AI agents by building on prior theoretical and empirical work in cognitive science that has analyzed how people form mental models of other people. To this end, we present a general framework for modeling human mental models of AI agents based on item response theory (IRT). In the traditional IRT approach, given an agent’s performance on problem sets involving a particular task, the IRT framework is used to estimate both problem set difficulties and agent abilities. We build on this traditional IRT methodology to propose a new framework that models human mental models both of themselves and of others, in terms of perceived abilities and perceived problem difficulties. This enables us to make comparative predictions about humans’ perceptions of AI agents, in the context of their perceptions of themselves and of other people.

To highlight the use of this framework, we conduct experiments in the context of question-answering where participants work alongside either another person or an AI agent (in the form of a large language model). Participants estimate the performance of their counterpart throughout the experiment, allowing us to make inferences about participants’ perceptions (mental models) of the abilities of their counterparts. In our analysis, we focus on two sets of research questions that have not yet been explored in the literature:

- (1) **Multidimensionality of Mental Models:** What is the dimensionality of humans’ mental models in the context of assessing task performance of other agents? In particular, do humans’ mental models of others capture multiple different abilities or areas of expertise, or do they estimate a single notion of ability, a “general intelligence”? We investigate the specific structure of these perceptions of ability, in particular, the correlations between different abilities, and how our findings differ between mental models of other humans and those of AI agents.
- (2) **Role of Self-Perception:** What role does a person’s self-perception play as they develop a mental model of another agent? For example, do people estimate how another agent’s abilities differ from their own? We consider multiple potential relationships between (a) humans’ perceptions of their own abilities and experienced problem difficulties and (b) their perceptions of other agents’ abilities and problem difficulties. Again, we explore how these relationships differ between mental models of other humans and AI agents.

To address these questions we make two main contributions in this paper. First, we present an extensive experimental dataset involving humans and AI in a question-answering context. This data directly captures human perceptions of self- and other-agent performance, providing insight into questions about humans’ mental models of other agents. Second, we introduce a theoretical model-based framework for directly modeling and analyzing humans’ mental

models of AI agents, and we use this modeling framework to gain insight into our experimental data.<sup>2</sup>

In Section 2, we review existing work on mental models and their role in human-AI teams. In Section 3, we describe our experimental setup, and we include empirical findings and data analysis in Section 4. Section 5 introduces our IRT-based framework, and in Sections 6 and 7 we present the methodology and results for our two Research Foci—Section 6 on the dimensionality of mental models and Section 7 on the influence of self-perception. Finally, we discuss key takeaways in Section 8 and conclude in Section 9.

## 2 RELATED WORK AND BACKGROUND

### 2.1 Mental Models and Collaboration

In general, *mental models* are simplified representations of the world that people use to process new information and make predictions [9, 59]. Our work focuses in particular on mental models of one’s self, of other people, and of AI agents, and is informed by foundational work in cognitive science in the areas of *metacognition* [24, 43], *theory of mind* [3, 27], and *theory of machine* [45], respectively.

The importance of mental models of other agents has received considerable emphasis in prior work on collaboration, particularly for collaboration among teams of humans. Specifically, the goal of shared mental models (SMMs) [48, 55, 56] necessitates that team members are aligned in terms of their perceptions of their team, strategy, and the task at hand. Information about the skills and knowledge of a teammate, which we refer to as an “other mental model” or OMM, is an important component of these SMMs, promoting effective collaboration [14, 23, 47, 49, 58]. In the context of the focus of this paper, namely hybrid human-AI teams, prior work has found that more accurate perceptions of AI agents tend to result in better team performance [5, 28] and more satisfying interactions for humans [37].

### 2.2 Understanding Mental Models of AI Agents

Given the importance of human perceptions of AI agents in human-AI collaboration, an emerging body of work aims to understand these perceptions [18, 68, 70]. Based on experiments in cooperative game settings, [28] delineated three different components of mental models of AI: the agent’s *knowledge distribution*, *local behavior*, and *global behavior*. There is evidence that people develop a mental model of an AI agent’s *knowledge distribution* based on their own knowledge [41], and that this perceived intelligence can be affected by provided explanations [52]. It has also been shown that humans develop perceptions of AI agents’ *global behavior*, but that these perceptions weaken as error boundaries become more complex and stochastic [6] and can be biased by first impressions of the agent [51]. Finally, prior research has shown that people can predict an AI agent’s *local behavior*, basing initial estimates on their own abilities [12], and that counterfactual examples can improve these predictions [2].

In this paper, we characterize mental models of AI agents in terms of the perceived ability of the agent and the perceived problem difficulty for the agent. While we use perceptions of behavior to estimate these quantities, all three of [28]’s components of mental

<sup>2</sup>All our code and data, including the original trivia questions, are available at [https://github.com/markellekelly/AI\\_mental\\_models](https://github.com/markellekelly/AI_mental_models).

models are relevant. A human’s perception of the abilities and problem difficulties for an AI agent could be influenced by how that human perceives the agent’s knowledge distribution. Perceived ability and problem difficulty could also be used to predict both local and global behavior.

## 2.3 Contributions in the Context of Related Work

The focus of this paper differs from prior work on mental models of AI agents in two important ways. First, we introduce a framework that directly models OMMs, and thus can describe them in terms of relevant latent variables and test hypotheses about their structure. In contrast, prior work generally has not investigated OMMs directly, instead analyzing a proxy such as team performance or humans’ predictions of AI behavior (e.g., [2, 7]). Second, in our experiments, we collect data on participants’ mental models of other people. Directly comparing mental models of AI agents to those of other people provides important context and helps determine how human-AI collaboration relates to general cognitive science research on teaming.

More specifically, our work differs from that of [6], who also experimentally investigated the capacity of humans to understand AI agents. Their analysis focused on the relationship between the complexity of the AI agent’s classification error boundary and the participant’s performance on the task. This earlier work differs from the work in this paper in that it did not directly capture or model participants’ OMMs, nor did the approach relate OMMs to self-perceptions or mental models of other people. Our work also differs from [51], who investigated experimentally how participants’ mental models of an AI agent were affected by their first impressions of the agent. Their results demonstrated that people developed more accurate mental models (i.e., had lower error in predicting model performance) when they had a positive first impression. However, the mental models themselves were not analyzed in terms of their structure, and they were not compared to self-perceptions or OMMs of people.

Finally, there has been recent prior work that has proposed analytic frameworks for human mental models of other humans. In particular, [38] developed an IRT framework for understanding these OMMs, which we build upon in this paper. Further, [69] introduced an AI agent that learns human OMMs of their human teammates for the purpose of improving collective intelligence in a human-human teaming context. This is relevant to the approach we propose in this paper in that it demonstrates how a framework such as ours could be deployed in a team scenario. However, these prior frameworks on modeling OMMs differ from our approach in that they do not investigate mental models of AI agents.

## 3 EXPERIMENTS

In our experiments, participants complete a multi-category trivia question-answering task, estimating their own performance and the performance of either another human or an AI agent. Experiments were conducted using Amazon Mechanical Turk.

### 3.1 Task

Each participant answered 16 sets of 12 trivia questions. We used a trivia setting because it does not require specialized knowledge (and thus is doable by Mechanical Turk workers), is discriminative (it is very unlikely humans or AI will answer all trivia questions correctly) [13], and can be broken up into distinct categories, allowing us to directly investigate multiple dimensions of ability.

Using a dataset of trivia questions from The Question Company, we selected four question topics: *History of Art*, *Video Games*, *Cities*, and *Math*. Sample questions for each topic are shown in Figure 1. Using preliminary data, we selected these topics to achieve (1) variation between participants for each topic, (2) variation among topics for each participant, and (3) varying correlations between pairs of topics across participants. In addition, we selected topics that we expected people would *perceive* as being different, e.g., requiring different types of knowledge.

For each category, we created four problem sets of 12 questions each, for a total of 16 problem sets. Participants were presented with questions in four rounds. Each round included one problem set from each trivia category. For each participant, we randomized the order of questions within each problem set, the order of categories within rounds, and the order of problem sets across rounds. Examples of possible problem set orders are shown in Figure 2.

### 3.2 Performance Assessment

After each 12-question problem set, participants were asked to estimate their own performance, to capture self-assessment, and the performance of another agent, to capture other-assessment. Participants were randomly assigned to assess either an “AI system” or another person (see Appendix C for details).

Two-thirds of the participants were randomly selected to receive feedback regarding their performance estimation: after they provided their estimates for a given problem set, they were shown their own actual performance, and the actual performance of the other agent, on that problem set. The remaining one-third of the participants did not receive this feedback. This no feedback group was included to capture prior expectations and to understand the strategies people use to estimate the performance of another agent in the absence of feedback.

### 3.3 “Other” Selection

To obtain the performance data for the other humans, we first performed a pilot study ( $n = 34$ ) using the same 16 12-question problem sets. Based on overall accuracy, we selected the top five (“high accuracy”) and bottom five (“low accuracy”) participants in the pilot study, where each of the top five and bottom five were then used as “other humans” in the main experiment. Specifically, participants in the main experiment that were assigned to the other human condition were shown performance data from one of these 10 participants. The name of this other human was randomly chosen from a set of ten names drawn from a random name generator (e.g., “Anna” or “Felix”).

To reduce the possibility of performance as a confounding variable, we then matched AI performance with the other human performance on a topic-wise basis. This was done by running several variants of UnifiedQA [35] and Zero-shot-CoT [36], which are large

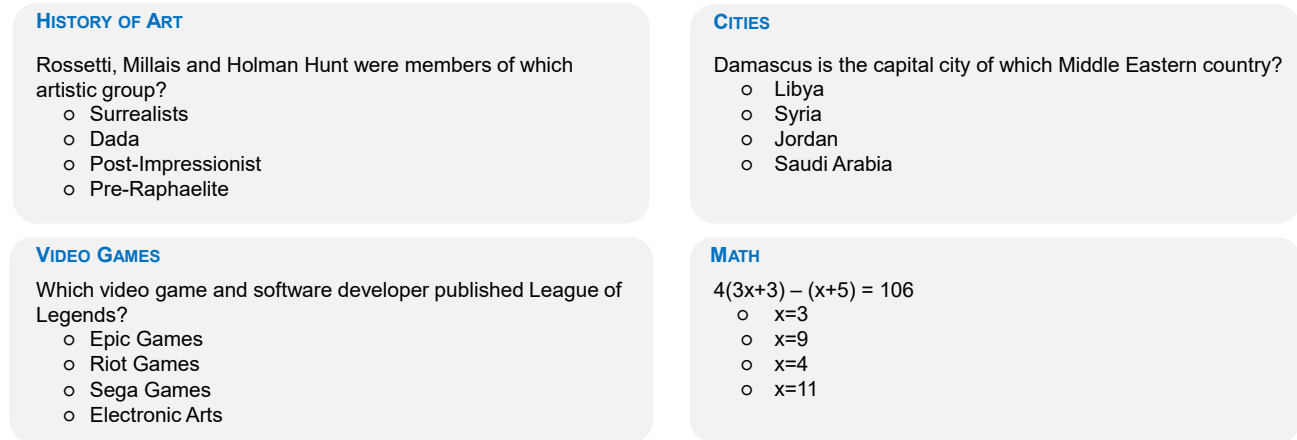


Figure 1: Sample questions from each trivia category.

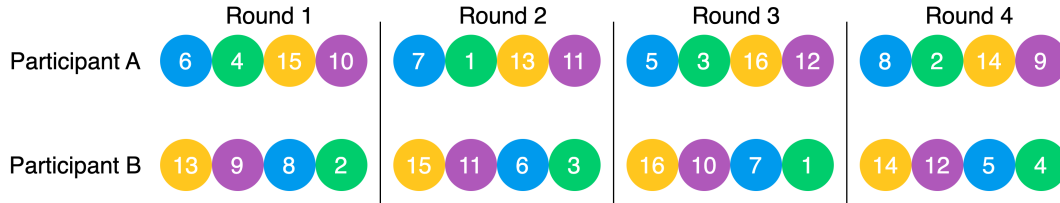


Figure 2: Example experiment configurations. Each circle represents a problem set, identified by its ID (number) and topic category (color). Each round has one problem set from each category and the order of these categories is consistent across rounds. Both this order and the order of individual problem sets is randomized for each participant.

language models designed to generalize to a range of tasks. We then chose two models for each topic, one with similar performance to that of the high accuracy humans, and another to match the performance of the low accuracy humans. (Details on the exact model settings used, and the final topic-wise accuracies, can be found in Appendix C.)

We include agents with both high and low accuracy to improve the generalization of results; mental models could differ depending on whether the other agent has higher or lower performance than the participant.

### 3.4 Setup

203 Amazon Mechanical Turk workers, all located in the U.S., participated in the study, which was conducted in January 2023. Participants could only complete the experiment once (and were disqualified if they had participated in the earlier pilot study). To ensure high-quality participation, workers were required to be AMT Masters and have a 95% approval rating; they were paid \$7 plus a bonus of up to \$2. These incentive bonuses were based on the participants' other-assessment performance. The experimental protocol was approved by the University of California, Irvine Institutional Review Board.

Participants were assigned to an agent type (other human or AI agent) and agent category (high accuracy, low accuracy, or no

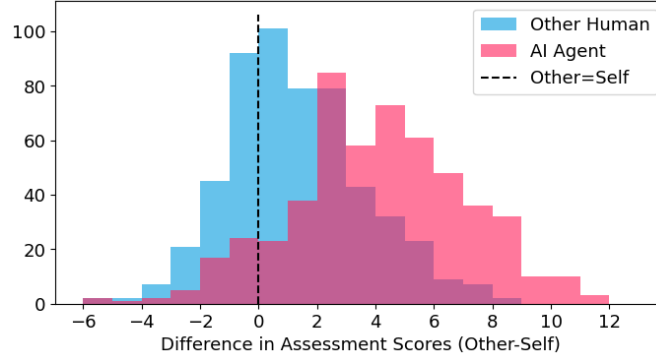
feedback). Participants were evenly divided between these six experimental condition combinations.

## 4 OVERVIEW AND EMPIRICAL ANALYSIS OF EXPERIMENTAL DATA

We begin our analysis with an investigation of initial findings from our experimental data, exploring the participants' perceptions of their own (self) performance, compared to their perceptions of the performance of others (both human and AI).

### 4.1 No Feedback Condition

We consider first the no feedback experimental condition, examining participants' OMMs when no performance information is available about the other agent. We analyze aggregate results for all 1072 problem set assessments in the no feedback condition. Each participant provided a score between 0 and 12 estimating their own performance (self-assessment) and a score estimating the performance of another (human or AI) agent (other-assessment), in terms of how many questions were answered correctly, for each of 16 problem sets. Figure 3 shows a histogram plot of the differences between self-assessment and other-assessment for both AI agents and other humans. Positive differences indicate that a participant provided a higher score for the other agent than for themselves.



**Figure 3: Histograms of differences between other agent (AI or Human) assessments and self assessments in the no feedback condition.**

The results in Figure 3 illustrate a striking difference between participants' assessments of AI agents and their assessments of other humans. The mean difference of the assessed performances of AI agents (relative to self) was +3.0 points, compared to +0.8 points for other humans (the two means are significantly different, and both are greater than 0, at a p-value threshold of  $\alpha = 0.01$  under two-sample and one-sample one-tailed t-tests, respectively). These results indicate that on average, people believe other agents will perform better than themselves in trivia question answering, and that AI agents will have much higher performance than other humans. Our findings are consistent with previous research indicating that people expect AI to be better at objective tasks, e.g., tasks that involve retrieving factual information, when compared to humans [16, 44].

## 4.2 Feedback Condition

In order to investigate how participants' OMMs adapt over time, given feedback about the performance of the other agent, we analyze the experimental data for all 2176 problem set assessments under the feedback condition. Of interest in this context is how feedback about the other agent (provided after each self-assessment, for 16 problem sets per participant) affects people's perceptions of the performance of the other human or of the other AI agent.

Figure 4 illustrates how participants update their assessments over rounds of feedback, starting from problem set 1 (when no feedback has been provided yet) up to problem set 16 (when feedback has been provided about all previous 15 problem sets). Participants adapt their assessments of other agents (both AI and human, both high and low accuracy) quickly within the first 4 to 6 problem sets and then change relatively slowly after that. Even after feedback from 16 problem sets, for the high accuracy agents, participants still systematically assess AI agents as being roughly 0.6 points more accurate than the human agents, even though the two agents were selected to have approximately the same accuracy (dotted lines). For the low accuracy agents there is also a consistent bias in favor of the AI agent, of roughly 0.2 points.

We also note that feedback appears to have relatively little influence in correcting self-assessed performance (in green). For example, for the low accuracy agents, after feedback on 16 problem sets,

even though the human and AI agents are much worse (-3 points) in terms of actual performance than the average participant (dotted lines), participants predict similar performance for themselves and the low accuracy agents.

In summary, the experimental results show that OMMs are quite different for AI agents and other humans for this question-answering task. People generally expect, a priori, that AI agents outperform humans, and it takes considerable evidence to adjust these expectations.

## 5 IRT FRAMEWORK

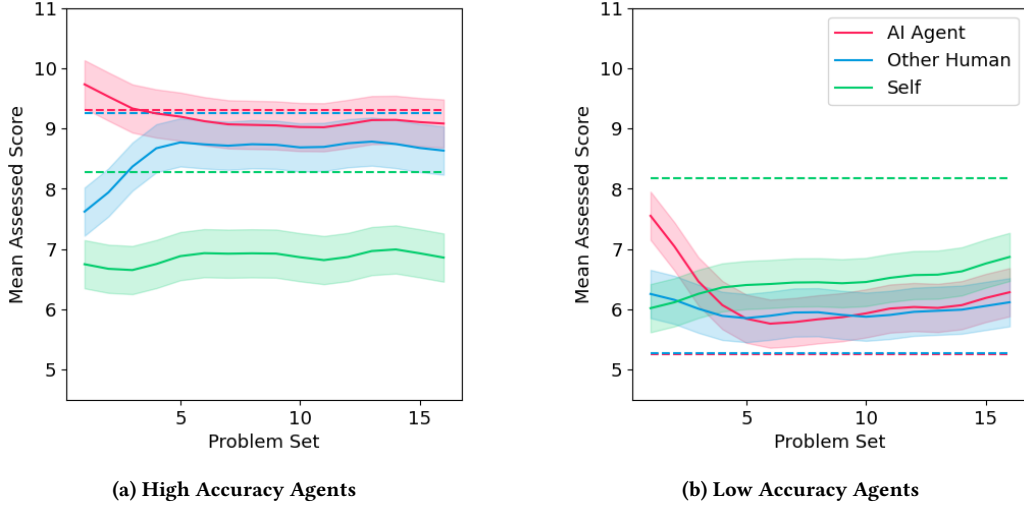
In this section we briefly outline our theoretical modeling framework for mental models, which we then use in Sections 6 and 7 to investigate Research Foci 1 and 2 posed in the Introduction. Our framework is based on Item Response Theory (IRT) [26, 65] which is widely used in education [10] and psychology [63] for modeling observed performance on a task in terms of latent psychological factors. In particular, we extend the hierarchical model of knowledge assessment proposed in [38], which focused on people's mental models of their own knowledge and the knowledge of other humans, but did not investigate mental models of AI agents.

In a standard IRT setup, we have  $V$  problems per problem set  $j$ , where each problem set  $j$  has a latent (unobserved) difficulty  $d_j$ . For each individual  $i$  and each problem set  $j$ , we model the number of items answered correctly  $x_{i,j}$ , ranging between 0 and  $V$ , as:

$$x_{i,j} = f(\theta_{i,j}) = f(a_i - d_j)$$

where  $a_i$  is the latent ability of individual  $i$  and  $f$  is a function that noisily converts the latent  $\theta_{i,j}$  for individual  $i$  and task  $j$  into an integer-valued  $x_{i,j}$ . The equation above represents how the IRT model can simulate or generate data in a forward manner. Given observed data  $x_{i,j}$ , for multiple participants  $i$  and problem sets  $j$ , we can then make inferences in the reverse direction about the latent abilities  $a_i$  and problem difficulties  $d_j$  (e.g., using standard Bayesian sampling techniques such as Markov Chain Monte Carlo sampling).

As in [38], our IRT framework analyzes *perceived* performance in addition to actual performance. In our experiments, we ask participants to estimate their own performance on each problem set.



**Figure 4: Mean perceived performance of other agents and self at each problem set. The results are separated by other agents with (a) high accuracy and (b) low accuracy. Dashed lines show corresponding values of actual performance for reference (for self, AI agents, and other humans, averaged across all participants and all problem sets in the given experimental condition). Results are smoothed across problem-sets to facilitate visual comparison.**

We denote this self-perceived data  $x_{i,j}^s$ , using an  $s$  superscript to signify self-assessment, with  $x_{i,j}^s = f(\theta_{i,j}^s) = f(a_i^s - d_j^s)$ . Participants are also asked to assess the performance of another human or an AI agent, capturing their OMM. We refer to this data with an  $o$  superscript, i.e.,  $x_{i,j}^o = f(\theta_{i,j}^o) = f(a_i^o - d_j^o)$ . Thus, using the other-assessment data  $x_{i,j}^o$ , we can then estimate  $a_i^o$ , the perceived ability of the other agent, as well as  $d_j^o$ , the perceived difficulty of problem  $j$  for the other agent, all from the perspective of participant  $i$ . Further details of our modeling framework are provided in Appendix B.

## 6 RESEARCH FOCUS 1: MULTIDIMENSIONALITY OF MENTAL MODELS

We next investigate whether people develop multidimensional OMMs, and more specifically, if people develop estimates of multiple different abilities of another agent. In the context of our experiments, we explore whether participants assess the strengths and weaknesses of the other agents across different topics. We are also interested in the correlational structure of these perceived abilities: do people expect other agents to have expertise in specific topics, or do they expect other agents to have more generalized intelligence? Finally, we investigate how these mental models develop over time and how these mental models differ between other humans and AI agents.

### 6.1 Methods

We address the question of multidimensional models of ability by comparing one-dimensional and multidimensional IRT models for other-assessment on our experimental data. Multidimensional item response theory (MIRT) [1, 54] models ability as a vector with multiple dimensions. For example, a student’s performance on a

standardized test might be related to their reading, writing, and mathematical abilities. In our framework, we use a between-items MIRT setup, which assumes that the probability of success for a specific item is affected by only one of the ability dimensions [29, 57].

In particular, we associate each of the four trivia topics (history of art, video games, cities, and math) (see Section 3.1) with an ability dimension. Thus, each component of  $\mathbf{a}_i^o$  corresponds to the perceived ability of the other agent in a specific trivia category. To understand the dimensionality of perceived abilities, we can then compare the one-dimensional and multidimensional models:

One-dimensional

$$\theta_{i,j}^o = a_i^o - d_j^o$$

Multidimensional

$$\theta_{i,j}^o = a_{ij}^o - d_j^o$$

where, in the multidimensional model,  $a_{ij}^o$  is the component of the  $k$ -dimensional ability vector  $\mathbf{a}_i^o$  that corresponds to problem  $j$ .

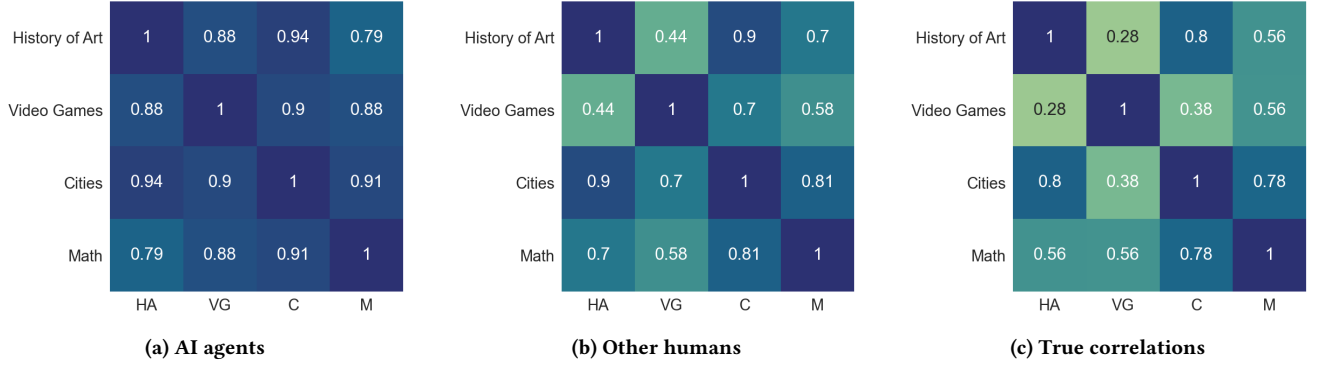
Further, the MIRT model estimates a  $4 \times 4$  matrix  $\Sigma^o$  of latent linear correlations between ability dimensions (see details in Appendix B). These estimated correlations capture the structure of perceived abilities, allowing us to quantify how participants expect these topic-wise abilities to be related.

### 6.2 Results

**6.2.1 Dimensionality of mental models.** To determine the dimensionality of participants’ OMMs, we compare the simpler one-dimensional model with the multidimensional model using three widely-used statistical model selection techniques: (i) held-out log-likelihood, (ii) WAIC score, and (iii) LOO score [46, 66].<sup>3</sup> For each participant, we hold out the final four problem sets they completed,

<sup>3</sup>In the main paper, for all model comparisons, we present only the held-out log-likelihood scores for brevity. See Appendix A for additional details; all three scores agree in terms of which models were selected in all model comparisons.





**Figure 5: Latent ability correlations for the four trivia categories in the assessment of AI agents (a) and other humans (b), and for the true performance of both humans and AI agents (c). The results are based on the feedback condition.**

and compute the log-likelihood over those four sets. Average held-out log-likelihoods across problem sets and participants are shown in Table 1. For reference, throughout the paper we contrast the performance of the IRT models with a discrete uniform baseline on [0, 12].

**Table 1: Held-out log-likelihood (higher is better) of the baseline, one-dimensional, and multidimensional models for other humans and AI agents, in the feedback condition.**

	Humans	AI
Baseline	-2.56	-2.56
One-dimensional	-1.81	-1.82
Multidimensional	-1.74	-1.72

The held-out log-likelihoods (and WAIC and LOO scores) indicate that the multidimensional model is a better model for both other-human and AI agent assessment. These results suggest that people can, and do, develop ideas of another agent’s strengths and weaknesses; mental models are not limited to a single ability.

**6.2.2 Latent correlations of ability dimensions.** Given that there is evidence that perceptions of others’ abilities are multidimensional, we are interested in their specific correlational structure. For instance, people might expect agents to have specific pockets of expertise or to exhibit a more general intelligence [32, 60].

To this end, we compare the latent correlations between trivia topics for other assessment, shown in Figure 5. Here we include only participants in the feedback condition, capturing differences in structure even when participants have observed the agent’s actual performance across these topics. We also include the true correlations across both humans and AI agents. (Note that these true correlations are similar between other humans and AI agents; the correlations split by agent type can be found in Appendix A.)

In summary, we find that the perceived correlations between abilities are significantly higher when assessing AI agents than when assessing other people. This means that participants expect the abilities of AI agents across different trivia categories to be highly correlated with each other, much more so than the correlations across abilities of another person. Note that this does not

contradict the earlier findings of multidimensionality in OMMs; the perceived abilities per topic are correlated, but they are still distinct from each other.

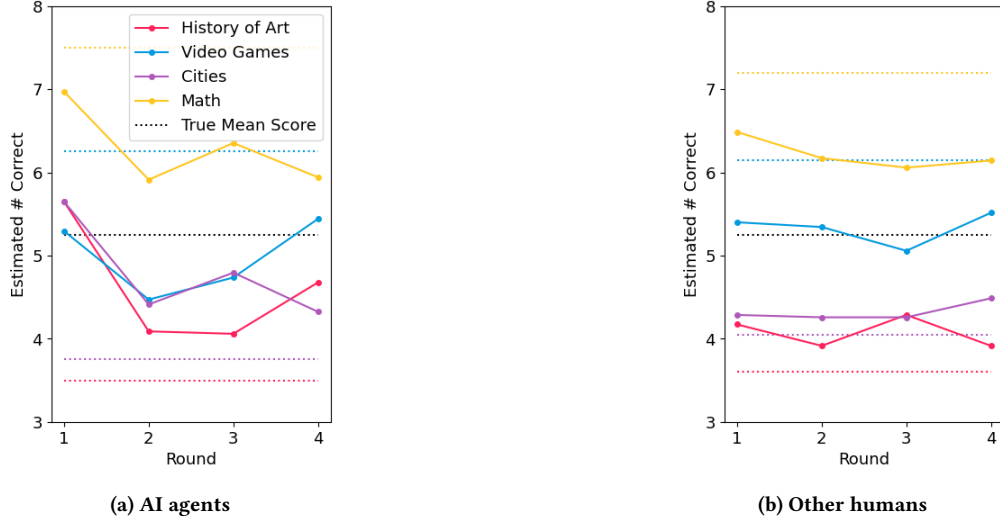
**6.2.3 Mental models over time.** In Section 4 we saw that other-assessment (on average) tends to stabilize, without converging to the other agent’s true performance, over the first 4 to 6 problem sets, with some continued drift for low accuracy agents. In this section we explore how other-assessment develops at the per-topic level, that is, how multidimensionality develops over time. In Figure 6, other-assessment data for low accuracy agents is plotted across rounds (where a round consists of 4 problem sets for a participant, one problem set from each topic). (Similar plots for high accuracy agents can be found in Appendix A.)

Figure 6 reflects that participants’ estimates of agent performance do not converge to the true agent performance at the per-topic level. Although in Round 4 there are differences, on average, between performance estimates for different topics (i.e. perceptions of ability are multidimensional), the estimates do not appear to have stabilized, and are not spread out enough to match true per-topic performance, particularly for AI agents. Instead, they appear to be anchored by the overall average performance (shown in black).

In line with the findings of Section 6.2.2, Figure 6 also reflects that people expect the abilities of AI agents to be highly correlated. Initially, participants expect similar performance across history of art, video games, and cities questions for AI agents (see the cluster of points in Round 1 in Figure 6a). For other humans, the estimates are more spread out across topics. After this first round, participants’ estimates of AI agent performance exhibit a similar decrease across all four categories—even though initial estimates of math and video games performance were, on average, underestimates. This reflects high perceived correlations between abilities: people expect AI agent scores to be closely related across topics, more so than for other humans (Figure 6b).

## 7 RESEARCH FOCUS 2: INFLUENCE OF SELF-PERCEPTION

In this section, we explore the role a person’s self-perception plays when forming a mental model of another agent. In particular, we investigate whether the role of self-perception in developing an



**Figure 6: Average other-assessment for low accuracy agents in the feedback condition. In both figures, the solid lines plot the average other-assessed performance across the four rounds. The average true performances of the other agent are shown as dotted horizontal lines, with the overall average in black and per-topic averages in their respective colors.**

OMM differs between AI agents and other people, as well as how this changes as more information about the other agent becomes available. Developing an understanding of the differences between one’s own capabilities and those of an AI agent is essential for improving cooperation [61], and explicitly comparing one’s own performance with the performance of an AI agent can promote appropriate selective reliance on the algorithm [42].

## 7.1 Methods

We test the fit of three different hierarchical IRT structures connecting the true (“underlying”) performance, self-assessment, and other-assessment data. We assume that self-assessment is (noisily) related to true performance. More specifically, we assume that self-assessed ability is a function of a person’s underlying ability, and self-perceived difficulty is a function of the underlying problem difficulty. We then test the relationship between self-assessment and other-assessment latent parameters, following a three-tier hierarchical structure.

We refer to the three setups as *undifferentiated*, *differentiated by ability*, and *fully differentiated*. Figure 7 depicts the graphical models for each of these structures. In the *undifferentiated* structure, the participant uses the same mental model to understand their own performance and the other agent’s performance; the model does not allow for differentiation between the participant’s own abilities and difficulties and those of the other agent. In the *differentiated by ability* setup, the participant learns a difference  $\delta$  between their own ability and the ability of the other agent, but problem difficulties remain undifferentiated. Finally, in the *fully differentiated* structure, the person does not use their own ability or difficulties to estimate those of the other agent; the self and other mental models are independent.

We test the fit of each of these three hierarchical models to determine which matches most closely with the true relationships

between parameters. Because we have evidence that these mental models are multidimensional (see Section 6), we use a multidimensional structure for the underlying, self-assessed, and other-assessed abilities;  $\delta$  is a  $k$ -dimensional vector capturing topic-wise ability differences.

## 7.2 Results

**7.2.1 Role of self-perception.** For both AI agents and other people, we train three different MIRT models, one for each hypothesized hierarchical structure. To capture mental model development over time, we evaluate models based on next-round predictions, that is, we compute the log-likelihood for round  $t$  using a model trained on data from rounds 1 to  $t - 1$ . These next-round log-likelihoods, under the feedback condition, are shown in Table 2.

**Table 2: Held-out next-round log-likelihoods (higher is better) for self differentiation models in the feedback condition.**

	Humans	AI
Baseline	-2.56	-2.56
Undifferentiated	-2.69	-3.36
Differentiated by Ability	<b>-2.00</b>	-2.18
Fully Differentiated	-2.13	<b>-2.13</b>

When feedback on performance is provided, the undifferentiated model has the worst fit overall, especially for AI agents—in fact, the undifferentiated model performs worse than random guessing (the baseline model). Thus, there is strong evidence that people differentiate between themselves and other agents. For perceptions of other humans, the differentiated by ability model fits the data best, suggesting that participants perceive other humans’ abilities in relation to their own abilities. In contrast, the fully differentiated model best explains the perceived scores of AI agents; this provides



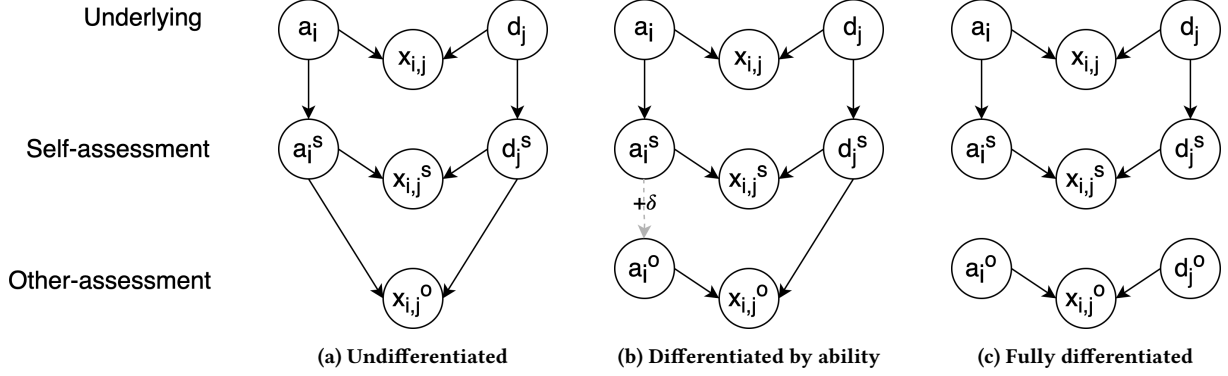


Figure 7: Assumed structures for the three different hierarchical models tested.

evidence that participants' mental models of themselves were less relevant in developing mental models of AI agents (in comparison to those of other humans).

**7.2.2 Ability differential.** The differentiated by ability model learns a parameter  $\delta$ , estimating the differential between a person's self-assessed ability and their perception of the other agent's ability. These estimates, for high-accuracy other agents, are shown in Table 3. Higher values of  $\delta$  correspond to higher perceived abilities of the other agent (relative to self-perceived ability). For reference, for the cities topic,  $\delta = 0.79$  corresponds to a 16-percentage-point increase in the latent perceived probability of a correct answer (i.e., the participant perceives the other agent's latent probability of success to be 16 percentage points higher than their own), whereas  $\delta = 1.72$  corresponds to a 27-percentage-point increase in that probability.<sup>4</sup>

Table 3: Latent per-topic differences between self-assessed and other-assessed ability for high-accuracy agents.

	Humans	AI
History of Art	0.83	1.63
Video Games	0.38	1.19
Cities	0.79	1.72
Math	0.75	1.69

The variation of  $\delta$ s between agent types suggests substantial differences between perceptions of other humans and AI agents. Specifically, the values of  $\delta$  are larger for AI agents than for other people, despite the true actual abilities of the AI agents and other people being very similar (by design) on each topic. This aligns with the findings from Section 4, in particular, that participants generally expect AI agents to perform at a significantly higher level than (other) humans.

**7.2.3 Relationship with self-perception.** In this section we investigate the overall relationship between self- and other-assessment in the absence of feedback. Note that these results, without feedback, are solely focused on each person's perception of the other agent, as a function of their self-perception; the true performance of the other

<sup>4</sup>These interpretations are computed using the median values of self-perceived ability and difficulty (0.54 and -0.14, respectively).

agent does not play any role since no feedback is provided. Figure 8 compares self- and other-assessment scores in the no feedback condition.

Figure 8 reflects a positive correlation between self-assessment and other-assessment when evaluating other humans. In comparison, the relationship between self-assessed performance and the perceived performance of AI agents is less pronounced. When participants perceive that they have done poorly on a particular problem set, they accordingly reduce their expectations of other humans. In contrast, participants predict similar scores for AI agents across self-assessed scores. Overall, this aligns with our findings from Sections 7.2.1 and 7.2.2, namely, that participants expect the performance of AI agents to be more different from their own performance, compared to the performance of other humans.

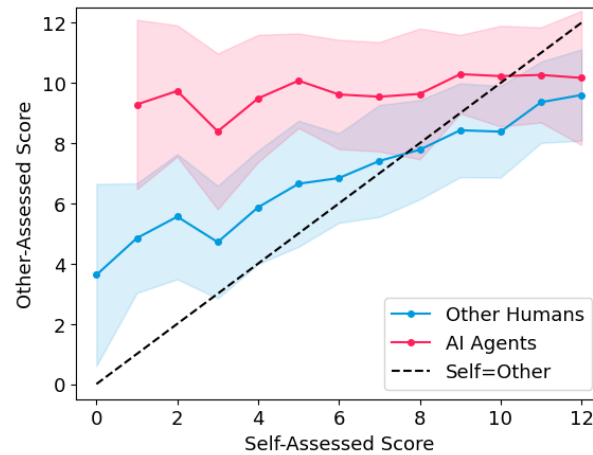
We also note that for lower self-assessed scores, participants expect other agents to score higher than them, while for higher self-assessed scores, other-assessment falls slightly below self-assessment (i.e., below the diagonal). This observation aligns with existing work in cognitive science, which has found that people believe they are better than others at easier tasks, but worse than others on more difficult tasks [25, 50]. Our experiments suggest that this finding, in particular that people believe they are worse than others on more difficult tasks, may be especially pronounced when the other is an AI agent.

## 8 DISCUSSION

First, we outline three key takeaways from our experiments and discuss how they might affect human-AI collaboration.

**1. Over-differentiation of AI from self.** We observe that mental models of an AI agent's ability are highly differentiated from self-perceived ability (e.g., Figures 3 and 4 and Table 3). On average, participants expect AI agents to perform very differently from themselves, especially in the absence of feedback. This bias could lead to under- or over-reliance on an AI agent in a team setting; additional work is needed to better understand and counteract it.

**2. "General intelligence" bias.** While we observe that people can pick up on the strengths and weaknesses of AI agents (see also [6, 51]), in our experiments participants expect a more unified, single intelligence from an AI agent than they do another person,



**Figure 8: Relationship between self-assessed score and other-assessed score in the no feedback condition. The dotted black diagonal line represents equality between self- and other-assessment.**

even after observing evidence to the contrary. In particular, the other agents in our experiments have near-identical inter-topic performance variations, but participants perceive much higher correlations between topics for AI agents (see Section 6.2.2). This can look like a failure to recognize how well the AI agent performs in its strongest areas, and how poorly it performs in its weakest areas (see Section 6.2.3), potentially resulting in over- or under-use of an AI decision making aid. Again, further research is needed to determine the extent of this bias and how it might be counteracted in human-AI teams.

**3. Incomplete development of mental models.** In our experiments, participants’ other-assessments did not fully converge to AI agents’ true performances, even given feedback (see Figure 4). This extended to the per-topic level; participants did not accurately estimate AI agents’ strengths and weaknesses after 16 rounds of feedback (see Figure 6). This phenomenon could serve as motivation to give the teammates of an AI agent extra information to aid in OMM development, e.g. a “primer” or onboarding process [15] or prediction explanations [2, 52].

Looking ahead, we believe our modeling framework could be useful for capturing people’s OMMs in the context of hybrid human-AI teams. In this paper, we investigate how a person perceives an AI agent, in terms of their different abilities, or strengths and weaknesses, and the difficulties of specific problems for the agent. Our findings and framework could help predict when a person is likely to defer to an AI agent (and thus help predict overall team performance) and identify biases that could lead to over- or under-use of the agent.

**Limitations.** Our experiments and results are limited to a single task and setting, and involve only Amazon Mechanical Turk workers, who are not necessarily knowledgeable on the task or on AI in general. In future work, it will be important to investigate these research questions across other settings, e.g., for image classification or other tasks, and with other users, e.g., human experts who already interact with an AI agent on a regular basis.

## 9 CONCLUSION

In this paper, we present an experimental dataset capturing participants’ mental models of themselves, other humans, and AI agents and introduce a framework for analyzing these mental models. Our findings indicate that (1) people tend to over-estimate the performance of AI agents relative to their own performance; (2) people expect the different abilities of AI agents to be highly correlated, even after observing evidence otherwise; and (3) these OMMs fail to develop completely, particularly in capturing agents’ different strengths and weaknesses. We anticipate that our modeling framework, and these findings, will be useful in both understanding and improving interaction in hybrid human-AI teams.

## ACKNOWLEDGMENTS

This research was supported by NSF under awards 1900644 and 1927245, by the Irvine Initiative in AI, Law, and Society, and by the Hasso Plattner Institute (HPI) Research Center in Machine Learning and Data Science at the University of California, Irvine.

## REFERENCES

- [1] Terry A Ackerman. 1994. Using multidimensional item response theory to understand what items and tests are measuring. *Applied Measurement in Education* 7, 4 (1994), 255–278.
- [2] Kamran Alipour, Arijit Ray, Xiao Lin, Michael Cogswell, Jürgen P. Schulze, Yi Yao, and Giedrius T. Burachas. 2021. Improving users’ mental model with attention-directed counterfactual edits. (2021). arXiv:2110.06863
- [3] Janet Wilde Astington and Jennifer M Jenkins. 1995. Theory of mind development and social understanding. *Cognition & Emotion* 9, 2-3 (1995), 151–165.
- [4] Joshua Attenberg, Panos Ipeirotis, and Foster Provost. 2015. Beat the machine: Challenging humans to find a predictive model’s “unknown unknowns”. *J. Data and Information Quality* 6, 1 (2015), 1–17.
- [5] Gagan Bansal, Besmira Nushi, Ece Kamar, Eric Horvitz, and Daniel S Weld. 2021. Is the most accurate AI the best teammate? Optimizing AI for teamwork. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 11405–11414.
- [6] Gagan Bansal, Besmira Nushi, Ece Kamar, Walter S Lasecki, Daniel S Weld, and Eric Horvitz. 2019. Beyond accuracy: The role of mental models in human-AI team performance. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing (HCOMP 2019)*, Vol. 7. 2–11.
- [7] Gagan Bansal, Besmira Nushi, Ece Kamar, Daniel S Weld, Walter S Lasecki, and Eric Horvitz. 2019. Updates in human-AI teams: Understanding and addressing the performance/compatibility tradeoff. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI 2019)*, Vol. 33. 2429–2437.

- [8] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. 2021. Does the whole exceed its parts? The effect of AI explanations on complementary team performance. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–16.
- [9] Winston H F Barnes. 1944. The nature of explanation. *Nature* 153, 3890 (1944), 605–605.
- [10] Ado Abdu Bichi and Rohaya Talib. 2018. Item response theory: An introduction to latent trait models to test and item development. *International Journal of Evaluation and Research in Education* 7, 2 (2018), 142–151.
- [11] Sebastian Bordt and Ulrike Von Luxburg. 2022. A bandit model for human-machine decision making with private information and opacity. In *Proceedings of the 25th International Conference on AI and Statistics (AI-Stats 2022)*. 7300–7319.
- [12] Nathan Bos, Kimberly Glasgow, John Gersh, Isaiah Harbison, and Celeste Lyn Paul. 2019. Mental models of AI-based systems: User predictions and explanations of image classification results. In *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, Vol. 63. 184–188.
- [13] Jordan Boyd-Graber and Benjamin Börschinger. 2020. What question answering can learn from trivia nerds. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 7422–7435.
- [14] Moritz C. Buehler and Thomas H. Weisswange. 2020. Theory of mind based communication for human agent cooperation. In *2020 IEEE International Conference on Human-Machine Systems (ICHMS)*. 1–6.
- [15] Carrie J. Cai, Samantha Winter, David Steiner, Lauren Wilcox, and Michael Terry. 2019. "Hello AI": Uncovering the onboarding needs of medical practitioners for human-AI collaborative decision-making. *Proceedings of the ACM on Human-Computer Interaction* 3, Article 104 (2019), 24 pages.
- [16] Noah Castelo, Maarten W Bos, and Donald R Lehmann. 2019. Task-dependent algorithm aversion. *Journal of Marketing Research* 56, 5 (2019), 809–825.
- [17] Chelsea Chandler, Peter W Foltz, and Brita Elvevåg. 2022. Improving the applicability of AI for psychiatric applications through human-in-the-loop methodologies. *Schizophrenia Bulletin* 48, 5 (2022), 949–957.
- [18] Valerie Chen, Q. Vera Liao, Jennifer Wortman Vaughan, and Gagan Bansal. 2023. Understanding the Role of Human Intuition on Reliance in Human-AI Decision-Making with Explanations. *arXiv preprint arXiv:2301.07255* (2023).
- [19] Hao-Fei Cheng, Logan Stapleton, Anna Kawakami, Venkatesh Sivaraman, Yanghui Cheng, Diana Qing, Adam Perer, Kenneth Holstein, Zhiwei Steven Wu, and Haiyi Zhu. 2022. How child welfare workers reduce racial disparities in algorithmic decisions. In *CHI '22: Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–22.
- [20] Maria De-Arteaga, Riccardo Fogliato, and Alexandra Chouldechova. 2020. A case for humans-in-the-loop: Decisions in the presence of erroneous algorithmic scores. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–12.
- [21] Greg d'Eon, Jason d'Eon, James R Wright, and Kevin Leyton-Brown. 2022. The Spotlight: A general method for discovering systematic errors in deep learning models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 1962–1981.
- [22] Kate Donahue, Alexandra Chouldechova, and Krishnamurthy Venkatesh. 2022. Human-algorithm collaboration: Achieving complementarity and avoiding unfairness. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 1639–1656.
- [23] Jeff Druce, James Niehaus, Vanessa Moody, David D. Jensen, and Michael L. Littman. 2021. Brittle AI, causal confusion, and bad mental models: challenges and successes in the XAI program. *CoRR abs/2106.05506* (2021). [arXiv:2106.05506](https://arxiv.org/abs/2106.05506)
- [24] John Dunlosky and Janet Metcalfe. 2008. *Metacognition*. Sage Publications.
- [25] David Dunning. 2011. The Dunning–Kruger effect: On being ignorant of one's own ignorance. In *Advances in Experimental Social Psychology*. Vol. 44. 247–296.
- [26] Jean-Paul Fox. 2010. *Bayesian Item Response Modeling: Theory and Applications*. Springer, New York.
- [27] Chris Frith and Uta Frith. 2005. Theory of mind. *Current Biology* 15, 17 (2005), R644–R645.
- [28] Katy Ilonka Gero, Zahra Ashktorab, Casey Dugan, Qian Pan, James Johnson, Werner Geyer, Maria Ruiz, Sarah Miller, David R Millen, Murray Campbell, et al. 2020. Mental models of AI agents in a cooperative game setting. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–12.
- [29] Johannes Hartig and Jana Höhler. 2009. Multidimensional IRT models for the assessment of competencies. *Studies in Educational Evaluation* 35, 2 (2009), 57–63.
- [30] Patrick Hemmer, Max Schemmer, Michael Vössing, and Niklas Kühl. 2021. Human-AI complementarity in hybrid intelligence systems: A structured literature review. In *Proceedings of the Twenty-fifth Pacific Asia Conference on Information Systems*. 1–14.
- [31] Kenneth Holstein and Vincent Aleven. 2022. Designing for human-AI complementarity in K-12 education. *AI Magazine* 43, 2 (2022), 239–248.
- [32] G. Humphreys, Lloyd. 1979. The construct of general intelligence. *Intelligence* 3, 2 (1979), 105–120.
- [33] Ece Kamar. 2016. Directions in hybrid intelligence: Complementing AI systems with human intelligence. In *Proceedings of the International Joint Conference on AI (IJCAI 2016)*. 4070–4073.
- [34] Ece Kamar, Severin Hacker, and Eric Horvitz. 2012. Combining human and machine intelligence in large-scale crowdsourcing. In *Proceedings of the AAMAS Conference*, Vol. 12. 467–474.
- [35] Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Taffjord, Peter Clark, and Hannaneh Hajishirzi. 2020. UnifiedQA: Crossing format boundaries with a single QA system. *arXiv preprint arXiv:2005.00700* (2020).
- [36] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*. 22199–22213.
- [37] Todd Kulesza, Simone Stumpf, Margaret Burnett, and Irwin Kwan. 2012. Tell me more? The effects of mental model soundness on personalizing an intelligent agent. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1–10.
- [38] Aakriti Kumar, Padhraic Smyth, and Mark Steyvers. 2023. Differentiating mental models of self and others: a hierarchical framework for knowledge assessment. *PsyArXiv* (2023).
- [39] David La Barbera, Kevin Roitero, and Stefano Mizzaro. 2022. A hybrid human-in-the-loop framework for fact checking. In *Proceedings of the Sixth Workshop on Natural Language for Artificial Intelligence (NL4AI 2022)*.
- [40] Vivian Lai, Chacha Chen, Q. Vera Liao, Alison Smith-Renner, and Chenhao Tan. 2021. Towards a science of human-AI decision making: A survey of empirical studies. *CoRR abs/2112.11471* (2021). <https://arxiv.org/abs/2112.11471>
- [41] Sau Lai Lee, Ivy Yee man Lau, S. Kiesler, and Chi-Yue Chiu. 2005. Human mental models of humanoid robots. In *Proceedings of the 2005 IEEE International Conference on Robotics and Automation*. 2767–2772.
- [42] Garston Liang, Jennifer F Sloane, Christopher Donkin, and Ben R Newell. 2022. Adapting to the algorithm: How accuracy comparisons promote the use of a decision aid. *Cognitive Research: Principles and Implications* 7, 1 (2022), 14.
- [43] Jennifer A Livingston. 2003. Metacognition: An Overview.
- [44] Jennifer Marie Logg. 2017. Theory of machine: When do people rely on algorithms? *Harvard Business School working paper series* 17-086 (2017).
- [45] Jennifer M Logg. 2022. The psychology of big data: Developing a "theory of machine" to examine perceptions of algorithms. In *The Psychology of Technology: Social Science Research in the Age of Big Data*, Sandra Matz (Ed.). American Psychological Association, 349–378.
- [46] Yong Luo and Khaleel Al-Harbi. 2017. Performances of LOO and WAIC as IRT model selection methods. *Psychological Test and Assessment Modeling* 59, 2 (2017), 183.
- [47] John Mathieu, Tonia Heffner, Gerald Goodwin, Eduardo Salas, and Janis Cannon-Bowers. 2000. The influence of shared mental models on team process and performance. *Journal of Applied Psychology* 85 (04 2000), 273–283.
- [48] Michael Merry, Pat Riddle, and Jim Warren. 2021. A mental models approach for defining explainable artificial intelligence. *BMC Medical Informatics and Decision Making* 21, 1 (2021), 1–12.
- [49] Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence* 267 (2019), 1–38.
- [50] Don A Moore and Daylian M Cain. 2007. Overconfidence and underconfidence: When and why people underestimate (and overestimate) the competition. *Organizational Behavior and Human Decision Processes* 103, 2 (2007), 197–213.
- [51] Mahsan Nourani, Chiradeep Roy, Jeremy E Block, Donald R Honeycutt, Tahrima Rahman, Eric Ragan, and Vibhav Gogate. 2021. Anchoring bias affects mental model formation and user reliance in explainable AI systems. In *Proceedings of the 26th International Conference on Intelligent User Interfaces*. 340–350.
- [52] Rohan Paleja, Muyleng Ghuy, Nadun Ranawaka Arachchige, Reed Jensen, and Matthew Gombolay. 2021. The utility of explainable AI in ad hoc human-machine teaming. In *Advances in Neural Information Processing Systems*, Vol. 34. 610–623.
- [53] Charvi Rastogi, Liu Leqi, Kenneth Holstein, and Hoda Heidari. 2022. A unifying framework for combining complementary strengths of humans and ML toward better predictive decision-making. *arXiv preprint arXiv:2204.10806* (2022).
- [54] Mark D. Reckase. 1997. The past and future of multidimensional item response theory. *Applied Psychological Measurement* 21, 1 (1997), 25–36.
- [55] Beau G. Schelble, Christopher Flathmann, Nathan J. McNeese, Guo Freeman, and Rohit Mallick. 2022. Let's think together! Assessing shared mental models, performance, and trust in human-agent teams. *Proceedings of the ACM on Human-Computer Interaction* 6, Article 13 (2022), 29 pages.
- [56] Matthias Scheutz, Scott A DeLoach, and Julie A Adams. 2017. A framework for developing and using shared mental models in human-agent teams. *Journal of Cognitive Engineering and Decision Making* 11, 3 (2017), 203–224.
- [57] Yanyan Sheng and Christopher K Wikle. 2007. Comparing multiunidimensional and unidimensional item response theory models. *Educ. Psychol. Meas.* 67, 6 (2007), 899–919.
- [58] Francesc Sidera, Georgina Perpiñà, Jèssica Serrano, and Carles Rostan. 2018. Why is theory of mind important for referential communication? *Current Psychology* 37 (2018), 82–97.
- [59] Mary M Smyth, Alan F Collins, Peter E Morris, and Philip Levy. 1994. *Cognition in Action* (2nd ed.). Lawrence Erlbaum Associates.

- [60] C. Spearman. 1904. "General intelligence," objectively determined and measured. *The American Journal of Psychology* 15, 2 (1904), 201–292.
- [61] Mark Steyvers and Aakriti Kumar. 2022. Three challenges for AI-assisted decision-making. *PsyArXiv* (2022). <https://doi.org/10.31234/osf.io/gctv6>
- [62] Mark Steyvers, Heliodoro Tejeda, Gavin Kerrigan, and Padhraic Smyth. 2022. Bayesian modeling of human–AI complementarity. *Proceedings of the National Academy of Sciences* 119, 11 (2022), e2111547119.
- [63] Michael L. Thomas. 2019. Advances in applications of item response theory to clinical assessment. *Psychological Assessment* 31, 12 (2019), 1442–1455.
- [64] Philipp Tschandl, Christoph Rinner, Zoe Apalla, Giuseppe Argenziano, Noel Codella, Allan Halpern, Monika Janda, Aimilios Lallas, Caterina Longo, Josep Malvehy, John Paoli, Susana Puig, Cliff Rosendahl, H Peter Soyer, Iris Zalaudek, and Harald Kittler. 2020. Human-computer collaboration for skin cancer recognition. *Nature Medicine* 26, 8 (2020), 1229–1234.
- [65] W.J. van der Linden and R.K. Hambleton. 2013. *Handbook of Modern Item Response Theory*. Springer, New York.
- [66] Aki Vehtari, Andrew Gelman, and Jonah Gabry. 2017. Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing* 27 (2017), 1413–1432.
- [67] Dayong Wang, Aditya Khosla, Rishab Gargeya, Humayun Irshad, and Andrew H Beck. 2016. Deep learning for identifying metastatic breast cancer. *arXiv preprint arXiv:1606.05718* (2016).
- [68] Qiaosi Wang, Koustuv Saha, Eric Gregori, David Joyner, and Ashok Goel. 2021. Towards Mutual Theory of Mind in Human-AI Interaction: How Language Reflects What Students Perceive About a Virtual Teaching Assistant. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 384, 14 pages. <https://doi.org/10.1145/3411764.3445645>
- [69] Samuel Westby and Christoph Riedl. 2022. Collective intelligence in human-AI teams: A Bayesian theory of mind approach. *ArXiv abs/2208.11660* (2022).
- [70] David Westerman, Autumn P. Edwards, Chad Edwards, Zhenyang Luo, and Patric R. Spence. 2020. I-It, I-Thou, I-Robot: The Perceived Humanness of AI in Human-Machine Communication. *Communication Studies* 71, 3 (2020), 393–408. <https://doi.org/10.1080/10510974.2020.1749683> arXiv:<https://doi.org/10.1080/10510974.2020.1749683>
- [71] Bryan Wilder, Eric Horvitz, and Ece Kamar. 2020. Learning to complement humans. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI-20)*. 1526–1533.

## A ADDITIONAL RESULTS AND FIGURES

### A.1 From Section 6

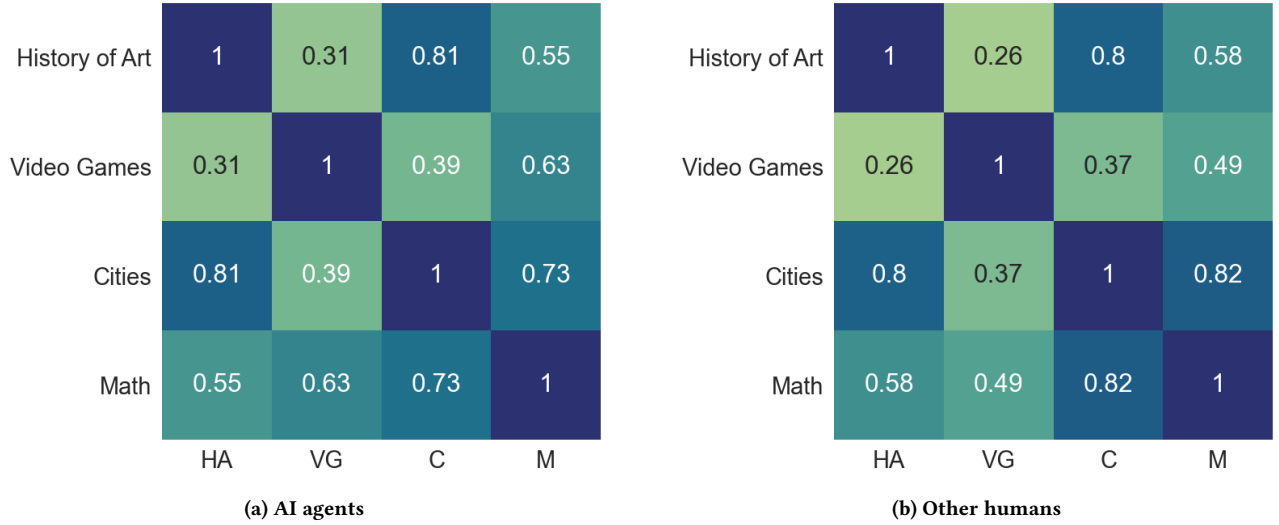
Here we include other metrics (WAIC and LOO scores) for the one-dimensional and multidimensional IRT models. Scores are computed on the log-score scale, i.e., a higher score is better. Note there are no standard errors for the baseline model, as no parameters are learned.

**Table 4: WAIC scores and standard errors of one-dimensional and multidimensional models for other humans and AI agents, across both feedback conditions.**

	Humans	AI
Baseline	-4227.04	-4144.96
One-dimensional	-3404.6 $\pm$ 28.0	-3348.5 $\pm$ 28.7
Multidimensional	-3024.9 $\pm$ 33.7	-2855.4 $\pm$ 40.3

**Table 5: LOO scores and standard errors of one-dimensional and multidimensional models for other humans and AI agents, across both feedback conditions.**

	Humans	AI
Baseline	-4227.04	-4144.96
One-dimensional	-3405.1 $\pm$ 28.0	-3349.1 $\pm$ 28.7
Multidimensional	-3031.5 $\pm$ 33.9	-2859.8 $\pm$ 40.4



**Figure 9: True ability correlations for per-topic performances, split by agent type. These are the latent correlations between abilities, computed from the true other agent performance data.**

### A.2 From Section 7

Here we include additional scores, averaged over each participant, for each of the three hierarchical models.

## B IRT MODEL DETAILS AND PRIORS

All models were fit using Stan. The underlying model (B.1) was used to model other-assessment data  $x_{i,j}^o$  in Section 6 as well as true performance  $x_{i,j}$  in the top level of the hierarchy in Section 7. The second level of the hierarchy was modeled by the self-assessment model (B.2). The three other-assessment models are detailed in B.3.

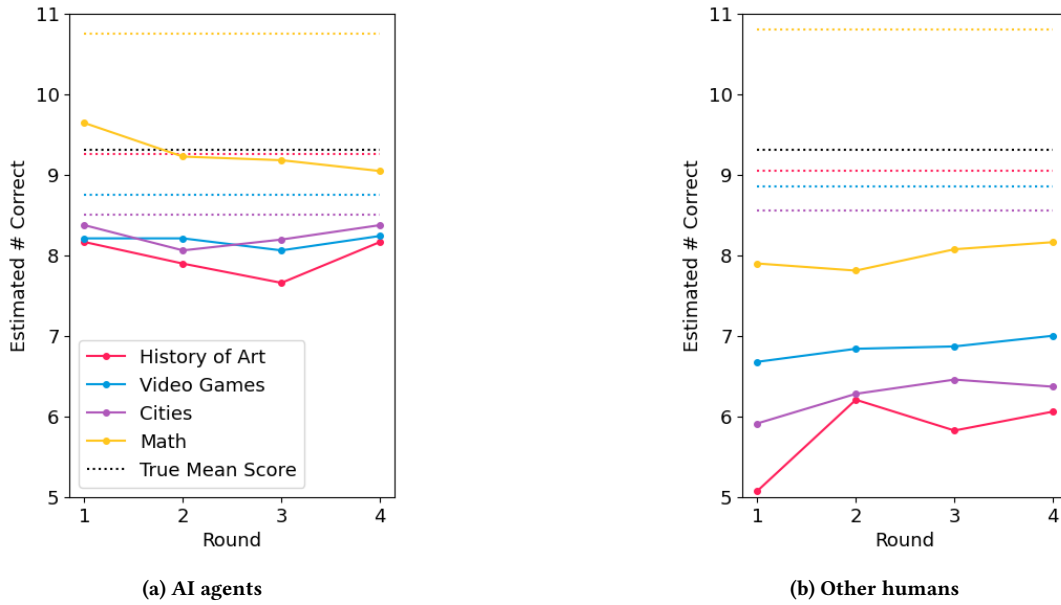


Figure 10: Average other-assessment for high accuracy agents (AI agents in (a), other humans in (b)) in the feedback condition. In both figures, the solid lines plot the average other-assessed performance across all four rounds. The average true performances of the other agent are shown as dotted horizontal lines, with the overall average in black and per-topic averages in their respective colors.

Table 6: Held-out next-round log-likelihoods (higher is better) for self differentiation models in the no feedback condition

	Humans	AI
Baseline	-2.56	-2.56
Undifferentiated	-2.68	-5.80
Differentiated by Ability	<b>-1.89</b>	<b>-1.72</b>
Fully Differentiated	-2.07	-1.77

Table 7: Average LOO scores for self differentiation models, across both feedback conditions

	Humans	AI
Baseline	-41.04	-41.04
Undifferentiated	-42.58	-67.89
Differentiated by Ability	-29.49	-30.29
Fully Differentiated	-30.58	-28.82

Table 8: Average WAIC scores for self differentiation models, across both feedback conditions

	Humans	AI
Baseline	-41.04	-41.04
Undifferentiated	-42.58	-67.89
Differentiated by Ability	-29.38	-30.17
Fully Differentiated	-29.94	-28.25

The underlying models were run with 800 warm-up iterations, 1500 samples, and three chains. Each of the self-assessment and other-assessment models (one for each participant) were run with 600 warm-up iterations, 1000 samples, and 2 chains. These hyperparameters were chosen based on chain convergence plots.



To convert latent scores  $\theta_{i,j}$  to discrete scores, we used:

$$p_{i,j} = \frac{1}{1 + \exp(-\theta_{i,j})}$$

$$x_{i,j} \sim \text{OrderedProbit}(p_{i,j}, v, \sigma)$$

where  $v$  is an array of cutoff points for conversion to discrete scores. We used 13 equally-spaced bins between 0 and 1 (converting into a score between 0 and 12, the number of questions in each problem set).

## B.1 Underlying Model

### Multidimensional

$$x_{i,j} = f(\lambda_j \cdot \mathbf{a}_i, d_j, \sigma)$$

$$\sigma \sim \text{Cauchy}(0, 2)$$

$$d_j \sim \text{N}(\mu_d, \sigma_d)$$

$$\mu_d \sim \text{N}(0, 2)$$

$$\sigma_d \sim \text{Cauchy}(0, 5)$$

$$\mathbf{a}_i \sim \text{MVN}(\mathbf{0}, \Sigma_L)$$

$$\Sigma_L = L_{\text{std}} \cdot L_{\Omega}$$

$$L_{\Omega} \sim \text{lkj\_corr\_cholesky}(1)$$

$$L_{\text{std}} \sim \text{N}(0, 2.5)$$

### One-dimensional

$$Y_{i,j} = f(a_i, d_j, \sigma)$$

$$\sigma \sim \text{Cauchy}(0, 2)$$

$$d_j \sim \text{N}(\mu_d, \sigma_d)$$

$$\mu_d \sim \text{N}(0, 2)$$

$$\sigma_d \sim \text{Cauchy}(0, 5)$$

$$a_i \sim \text{N}(0, 1)$$

## B.2 Self-Assessment Model

### Multidimensional

$$x_{i,j}^s = f(\lambda_j \cdot \mathbf{a}_i^s, d_j^s, \sigma^s)$$

$$\sigma^s \sim \text{Cauchy}(0, 2)$$

$$d_j^s \sim \text{N}(\gamma \cdot d_j + \Lambda, \sigma_{d,i})$$

$$\gamma \sim \text{N}(0, 1)$$

$$\Lambda \sim \text{N}(0, 1)$$

$$\sigma_{d,i} \sim \text{Cauchy}(0, 2)$$

$$a_{i,k}^s \sim \text{N}(a_{i,k}, \sigma_{a,i})$$

$$\sigma_{a,i} \sim \text{Cauchy}(0, 2)$$

### One-dimensional

$$x_{i,j}^s = f(a_i^s, d_j^s, \sigma^s)$$

$$\sigma^s \sim \text{Cauchy}(0, 2)$$

$$d_j^s \sim \text{N}(\gamma \cdot d_j + \Lambda, \sigma_{d,i})$$

$$\gamma \sim \text{N}(0, 1)$$

$$\Lambda \sim \text{N}(0, 1)$$

$$\sigma_{d,i} \sim \text{Cauchy}(0, 2)$$

$$a_i^s \sim \text{N}(a_i, \sigma_{a,i})$$

$$\sigma_{a,i} \sim \text{Cauchy}(0, 2)$$

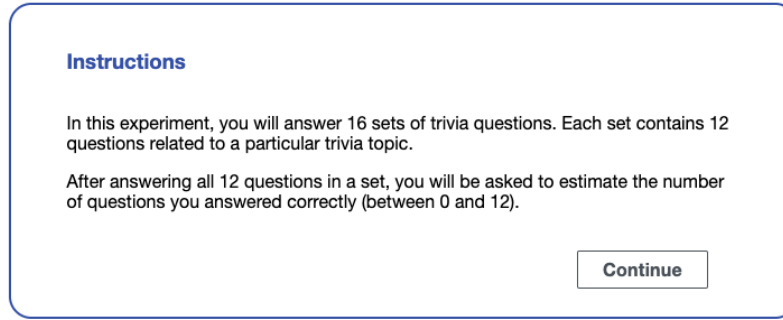


Figure 11: First page of instructions (for all participants).

### B.3 Other-Assessment Models

#### B.3.1 Undifferentiated.

##### Multidimensional

$$x_{i,j}^o = f(\lambda_j \cdot a_i^s, d_j^s, \sigma^s)$$

Input data:  $a_i^s, d_j^s, \sigma^s$

##### One-dimensional

$$x_{i,j}^o = f(a_i^s, d_j^s, \sigma^s)$$

Input data:  $a_i^s, d_j^s, \sigma^s$

#### B.3.2 Differentiated by Ability.

##### Multidimensional

$$x_{i,j}^o = f(\lambda_j \cdot a_i^o, d_j^s, \sigma^s)$$

$$a_{i,k}^o = a_{i,k}^s + \delta_{i,k}$$

$$\delta_{i,k} \sim N(0, 1)$$

Input data:  $a_i^s, d_j^s, \sigma^s$

##### One-dimensional

$$x_{i,j}^o = f(a_i^o, d_j^s, \sigma^s)$$

$$a_i^o = a_i^s + \delta_i$$

$$\delta_i \sim N(\mu_{\delta_i}, \sigma_{\delta_i})$$

$$\mu_{\delta_i} \sim N(0, 1)$$

$$\sigma_{\delta_i} \sim \text{Cauchy}(0, 2)$$

Input data:  $a_i^s, d_j^s, \sigma^s$

#### B.3.3 Fully Differentiated.

##### Multidimensional

$$x_{i,j}^o = f(\lambda_j \cdot a_i^o, d_j^o, \sigma^s)$$

$$d_j^o \sim N(\mu_d^o, \sigma_d^o)$$

$$\mu_d^o \sim N(0, 2)$$

$$\sigma_d^o \sim \text{Cauchy}(0, 5)$$

$$a_{i,k}^o \sim N(0, 1)$$

Input data:  $\sigma^s$

##### One-dimensional

$$x_{i,j}^o = f(a_i^o, d_j^o, \sigma^s)$$

$$d_j^o \sim N(\mu_d^o, \sigma_d^o)$$

$$\mu_d^o \sim N(0, 2)$$

$$\sigma_d^o \sim \text{Cauchy}(0, 5)$$

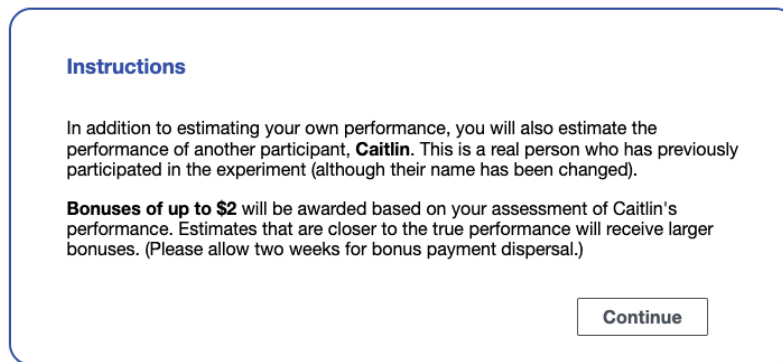
$$a_i^o \sim N(0, 1)$$

Input data:  $\sigma^s$

## C EXPERIMENT SETUP

The instructions provided to participants at the beginning of the experiment are shown in Figures 11 through 14. The second page of instructions depends on whether the participant was assigned another human or an AI agent to assess.

To match AI agent performance to that of human performance, we first chose the five highest-accuracy and five lowest-accuracy participants in our pilot study. We then ran six models: three versions of UnifiedQA (base, large, and 3B, see documentation) and three versions of Zero-shot-CoT (one with GPT3-XL version 1 and method “zero\_shot,” two with method “zero\_shot\_cot” and GPT3-XL versions 1 and 3, respectively; see documentation). Testing this wide variety of models enabled us to match accuracy relatively closely to the humans. The models used for each topic are shown in Table 9.



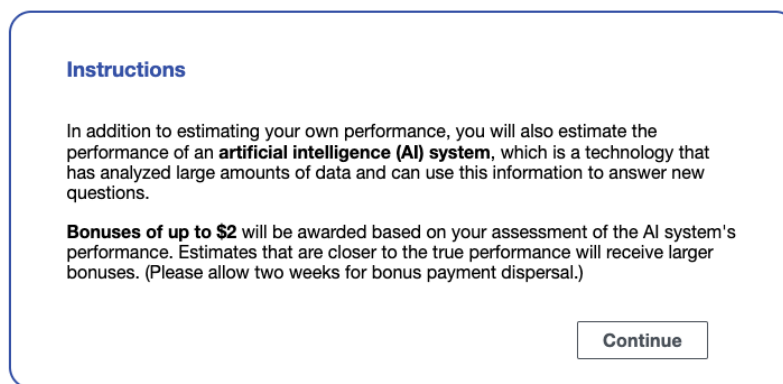
**Instructions**

In addition to estimating your own performance, you will also estimate the performance of another participant, **Caitlin**. This is a real person who has previously participated in the experiment (although their name has been changed).

**Bonuses of up to \$2** will be awarded based on your assessment of Caitlin's performance. Estimates that are closer to the true performance will receive larger bonuses. (Please allow two weeks for bonus payment dispersal.)

Continue

Figure 12: Second page of instructions shown to participants in the “other human” condition.



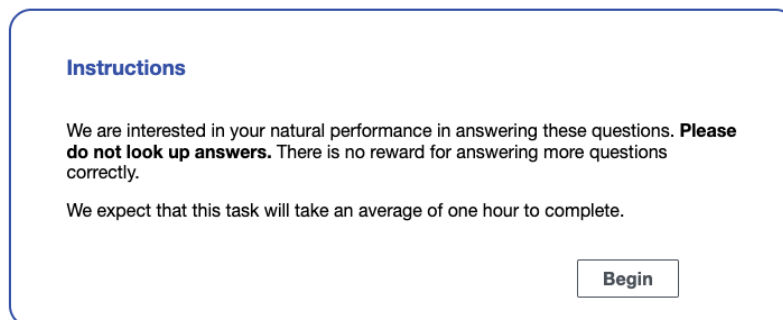
**Instructions**

In addition to estimating your own performance, you will also estimate the performance of an **artificial intelligence (AI) system**, which is a technology that has analyzed large amounts of data and can use this information to answer new questions.

**Bonuses of up to \$2** will be awarded based on your assessment of the AI system's performance. Estimates that are closer to the true performance will receive larger bonuses. (Please allow two weeks for bonus payment dispersal.)

Continue

Figure 13: Second page of instructions shown to participants in the “AI agent” condition.



**Instructions**

We are interested in your natural performance in answering these questions. **Please do not look up answers.** There is no reward for answering more questions correctly.

We expect that this task will take an average of one hour to complete.

Begin

Figure 14: Third page of instructions (for all participants).

How many of the 12 questions in this set do you think **you** answered correctly?

0 1 2 3 4 5 6 7 8 9 10 11 12

How many of the 12 questions in this set do you think **Sonia** answered correctly?

0 1 2 3 4 5 6 7 8 9 10 11 12

Next

Your estimated performance was 6 and you actually answered 8 correctly.

0 1 2 3 4 5 6 7 8 9 10 11 12

You estimated Sonia's performance to be 8 and they actually answered 5 correctly.

0 1 2 3 4 5 6 7 8 9 10 11 12

Next

**Figure 15: Example performance estimation questions. In this example, the participant estimates their own performance and the performance of another person, Sonia (left). They receive feedback about their own, and Sonia's, actual performance (right).**

**Table 9: Models used for AI Agent**

	Topic	Human	Model	UQA (base)	UQA (large)	UQA (3B)	ZS (001)	ZSC (001)	ZSC (003)
High accuracy	Art	75.4	77.1				✓		
	Video Games	73.8	72.9					✓	
	Cities	71.3	70.8			✓			
	Math	90.0	89.6						✓
Low accuracy	Art	30.0	29.2	✓					
	Video Games	51.3	52.1		✓				
	Cities	33.8	31.3	✓					
	Math	60.0	62.5				✓		