# Differentiating Mental Models of Self and Others: A Hierarchical Framework for Knowledge Assessment

Aakriti Kumar[1], Padhraic Smyth[2], and Mark Steyvers[1]
[1] Department of Cognitive Sciences, University of California, Irvine
[2] Department of Computer Science, University of California, Irvine

Developing an accurate model of another agent's knowledge is central to communication and cooperation between agents. In this article, we propose a hierarchical framework of knowledge assessment that explains how people construct mental models of their own knowledge and the knowledge of others. Our framework posits that people integrate information about their own and others' knowledge via Bayesian inference. To evaluate this claim, we conduct an experiment in which participants repeatedly assess their own performance (a metacognitive task) and the performance of another person (a type of theory of mind task) on the same image classification tasks. We contrast the hierarchical framework with simpler alternatives that assume different degrees of differentiation between mental models of self and others. Our model accurately captures participants' assessment of their own performance and the performance of others in the task: Initially, people rely on their own self-assessment process to reason about the other person's performance, leading to similar self- and other-performance predictions. As more information about the other person's ability becomes available, the mental model for the other person becomes increasingly distinct from the mental model of self. Simulation studies also confirm that our framework explains a wide range of findings about human knowledge assessment of themselves and others.

*Keywords:* metacognition, theory of mind, mind reading, other-assessment, Bayesian modeling

Understanding and comparing the knowledge states of others, to our own knowledge, is a fundamental skill that supports social interaction in daily life. Does Akira know what I know? Would Georgina perform better than me on this task? Will this problem be as difficult for Keith as it is for me? Humans constantly make predictions about their abilities at different tasks and how well other people might fare at the same task relative to themselves. For an individual making predictions about the difficulty of a task for others, a potential starting point is to base it on their own experience with the task (Nickerson, 1999) such as remembering information (Jameson et al., 1993; Koriat & Ackerman, 2010) or solving problems (Kelley & Jacoby, 1996). One's mental model about oneself may often lead to accurate predictions about others. However, previous research has not explored how the mental model of another person can be differentiated to account for specific information learned about them. When we observe another person over time, what is the process by which an initial undifferentiated mental model of that person becomes tailored toward them?

Our research combines ideas from (a) metacognition, which includes processes used to draw inferences about one's own knowledge states, and (b) theory of mind (also known as mind reading), which includes processes used to draw inferences about other people's knowledge states. Recent computational perspectives have suggested that reasoning processes about self and others are closely intertwined (Fleming, 2021). For example, a recent model for metacognition has been motivated by considering self-evaluation as a "second-order" computation distinct from simpler first-order accounts in which the same internal state guides decisions and self-evaluation (Fleming & Daw, 2017). Such second-order computation is also required when assessing the knowledge states of other people. Similarly, computational models for mind reading have been motivated by inverse planning—the process by which other people's goals and beliefs are inferred by applying one's own mental model to the observed actions (Aboody et al., 2021; Baker et al., 2009, 2017; Berke & Jara-Ettinger, 2021; Tauber & Steyvers, 2011). Empirical studies have provided increasing support for commonalities between metacognition and theory of mind based on shared cognitive resources (Nicholson et al., 2021), overlapping brain structures (Vaccaro & Fleming, 2018), and overlapping developmental trajectories (Gopnik & Astington, 1988; Paulus et al., 2014; but see Baer et al., 2021). Taken together, there is substantial evidence for a close correspondence between reasoning about self and others.

In this article, we present a hierarchical framework for knowledge assessment that explains how people assess their own knowledge and the knowledge of others. The framework is inspired by the connection between metacognition and theory of mind and has significant implications for understanding knowledge assessment in general. We focus on the relationship between *self-assessment* (i.e., predicting one's performance on a task) and *other-assessment* (i.e.,

predicting how well another person performs on the same task). There are two types of empirical results that the hierarchical framework is designed to address. First, the model can be used to explain the relationship between self- and other-assessment in situations where there is a lack of information about the other person being judged. For example, people are asked to assess the percentage of randomly selected students who know the answer to a given question (Nickerson et al., 1987; Tullis, 2018) or their relative placement in a population (Dunning, 2011; Moore & Healy, 2008). These studies have shown that people tend to predict that they are better than others on easy tasks but worse than others on challenging tasks (Moore & Cain, 2007). In these tasks, people consider comparisons to randomly sampled other individuals from a population. In later sections, we show how our framework may be applied to these experimental settings and demonstrate its ability to explain the empirical results observed in the literature. Second, the hierarchical framework also accounts for situations where people learn to make predictions about a *specific* person as information about that person becomes available. Our framework can also explain how people assess a specific other person by observing their performance on a task over time. To test our framework's predictions, we conduct a behavioral experiment where participants classify images and assess their own performance and the performance of a specific other person on this task. This experimental setup allows us to investigate two distinct aspects of assessing others: how individuals assess another individual without any explicit information about the other's ability and how this assessment changes as information about the other's performance becomes available. We also apply our framework to explain other assessment in paradigms where no information is provided about the other person (Moore & Healy, 2008; Tullis, 2018). Throughout this article, we assume that performance is indicative of a person's knowledge or ability. However, our proposed framework could also be applied to other domains that are not related to knowledge. For example, inferring a person's strength when observing them perform specific exercises in a gym or assessing the skill of drivers by observing them in challenging parking situations.
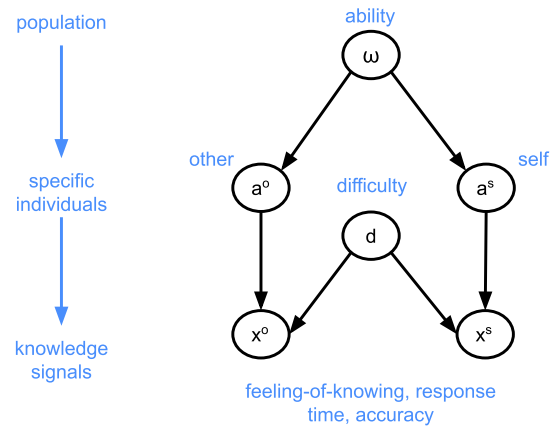
In the following sections, we provide a detailed overview of our modeling framework. We then present data from a knowledge assessment task in which people assess their own performance and the performance of one other person on an image classification task. We apply our proposed framework and simpler alternative models to this empirical data and demonstrate that the predictions of our hierarchical model closely match the trends observed in the data. We also show how our framework supports other findings in the empirical literature on knowledge assessment. Finally, we discuss the significance and implications of this framework for future research.

## A Hierarchical Framework for Knowledge Assessment

We propose a hierarchical framework for knowledge assessment that describes the computational problem that people solve when assessing themselves or another person. We posit that both self-assessment and other-assessment are inference problems that people solve through Bayesian inference. Figure 1 illustrates the different levels of the framework and the graphical model corresponding to it. The central idea underlying our framework is that reasoning about the performance of oneself or another person occurs at three different levels:

**Figure 1**

*Three Levels of the Hierarchical Model Used to Reason About One's Own as Well as Other People's Performance*



*Note.* People may have access to different kinds of knowledge signals such as feeling of knowing, response time, and accuracy when assessing their own knowledge or another person's knowledge. See the online article for the color version of this figure.

1. Population level: The top level corresponds to the population level (ω), which encodes information about the population of individuals to which the self and the other belong.

2. Individual-specific level: The middle level pertains to information about specific people (including self and others) such as the ability of self and others ($a^s$, $a^o$), the difficulty of the task perceived by self and others ($d$).

3. Knowledge-signals level: The bottom-most level concerns knowledge signals ($x$), which include observed performance outcomes for self and/or others and internal metacognitive signals that people may have access to when doing a task.

We assume that people can reason across the three levels and make inferences about self- or other-performance $a^s$, $a^o$, as well as task difficulty $d$ using the observed knowledge signals $x$. To enable reasoning across abilities of people and difficulties of items in tasks, the hierarchical framework adopts concepts from item response theory (IRT; Fox, 2010; van der Linden & Hambleton, 2013) to describe the relationship between $x$ and $a_s$, $a_o$, $d$. IRT has recently been used to model self-assessment (Jansen et al., 2020, 2021). Similar to the model by Jansen et al. (2021), we hypothesize that people make errors in their self-assessment such that their predicted performance deviates from the actual performance that would be predicted by an item-response model. Specifically, we assume that people combine a *subjective* estimate of ability with a *subjective* estimate of task difficulty in order to estimate the performance on a task.

To support inferences about ability and task difficulty, our work builds on previous research (Koriat, 1997; Moore & Healy, 2008; Nickerson, 1999; Thomas & Jacoby, 2013), which identifies a variety of signals that people use for assessment. In our framework, we assume that people may have access to two kinds of knowledge

signals ($x$) while performing a task. The first kind is based on *external signals*, such as feedback on people's assessment of self or others, information about the correct or optimal solution to a problem, or information about the other's performance. For example, in some tasks, people may receive feedback about their accuracy, which could be used as an external signal to infer their ability and predict future performance. The second kind of signal is *internal signals* that arise from reflecting on one's internal metacognitive processing. These include how long it takes people to arrive at a solution (Thomas & Jacoby, 2013; Tullis, 2018), their confidence in their response (Hart, 1965; Leibert & Nelson, 1998; Nelson & Narens, 1980), or their feeling of knowing about the problem at hand (Koriat, 2000). We use feeling of knowing to refer to the intuition that one may have about being able to solve a problem or answer a question without actually attempting to solve the problem or answer the question (e.g., when reading a general knowledge question, one may feel the question is answerable based on the familiarity with the words in the question).

Knowledge signals allow people to make estimates of individual-specific parameters such as the ability of self and others, and the perceived difficulty of the task. Depending on the available signals, our framework suggests two ways in which people may infer the ability of others:

1. In the absence of specific information about others (e.g., the inference is about a randomly sampled person from the population), people may use the knowledge signals regarding their own performance and metacognition ($x^s$) to reason about the ability of others. This corresponds to inferring $p(a^o|x^s)$.

2. If some information about the other person is available, people may also consider a combination of their own and others' knowledge signals to infer $p(a^o|x^s, x^o)$.

The first inference problem maps directly onto previous research where no information is provided about others (Moore & Healy, 2008; Nickerson, 1999; Tullis, 2018). The second inference problem has not been studied previously. In the next section, we present results from an experimental paradigm where participants track the performance of a specific other person and are provided with an increasing amount of information about the other person's performance. The framework also extends to assessing multiple other people. Note that, in many real-world contexts, people already have an estimate of their own ability on a variety of tasks: They gather information about their ability over time through varied interactions with other agents and environments. Hence, $a^s$ may be partially or fully observed in these cases. In comparison, people typically have less information about other people's abilities. Therefore, in most cases, $a^o$ is unobserved and must be inferred. As a result, people's assessment of their own abilities and knowledge will be less noisy than their assessment of others (Moore & Healy, 2008).

People must also reason about the task at hand when doing self- or other-assessment. External signals such as accuracy may enable people to better assess the difficulty ($d$) of the task at hand. Internal signals such as the time it takes people to solve a problem may provide additional information about the difficulty of the task and help predict how others would fare at the same task. For example, people may infer that questions that take them longer to answer are more difficult and may take others longer to answer as well.

Together, these internal and external signals provide information that people may use to infer task difficulty (Kelley & Jacoby, 1996).

The top level of the hierarchy formalizes the assumption that any person's ability, including one's own, is a sample from a population's ability distribution, which is denoted by ω. Note that ω may vary across tasks and population composition. Consider a chemistry teacher who is about to begin teaching a lesson on stoichiometry to a group of students who have never studied it. She has, however, observed other students of the same grade in the past and can easily make inferences about how well the new batch of students might fare on a test before and after her lesson. This is because the teacher assumes that any new student may be considered a random sample from the population of all students. She would also have a reasonable understanding of what questions the students might find difficult. On the other hand, if asked to compare her own knowledge of stoichiometry to another chemistry teacher, she would think about the population of chemistry teachers (which also includes herself) and her placement in this population. Therefore, people's assessment of the ability of others starts with assumptions about the population they are evaluating. In this article, we focus on people's assessment of others from the same population as themselves. However, it is straightforward to extend our framework to model how people assess individuals from different populations or even artificial agents. One way to do this is to add another level to the current hierarchy: Two populations may be considered samples from a superpopulation of agents.

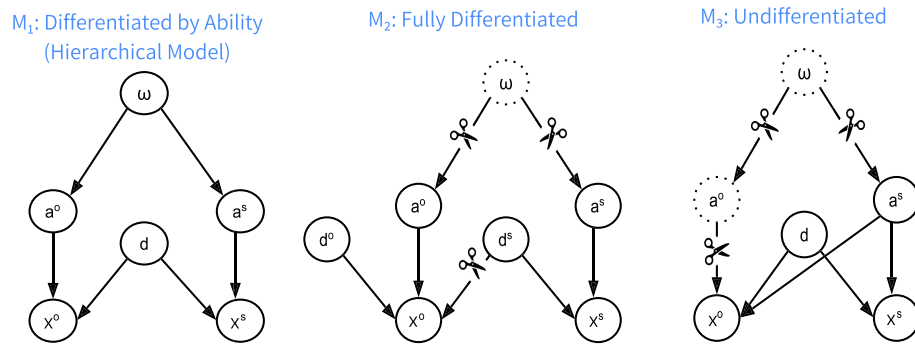## Three Instantiations of the Hierarchical Framework

Within this hierarchical approach to knowledge assessment, we explore three classes of models for connecting the subjective estimates of self and others as illustrated in Figure 2. These models correspond to different substantive assumptions about the psychological process of other assessments in terms of the assumed connections between the different layers of the hierarchy. The first instantiation, *differentiated by ability* model, is equivalent to the full hierarchical model. The second instantiation, the *fully differentiated* model, assumes that self- and other-assessment are distinct processes. The *undifferentiated* model assumes no distinction between self- and other-assessment. We will also refer to these models with the shorthand notation $M_1$, $M_2$, and $M_3$, respectively.

### Differentiated by Ability Model ($M_1$)

This model maps directly to the proposed hierarchical model of knowledge assessment. One way to formalize the reasoning process in this model is that people separately assess their own ability ($a^s$) and the ability of another person ($a^o$). However, because the hierarchical structure imposes connections between the self and other ability (e.g., with no knowledge of the other person, the best estimate of another person equals that one of one's own ability, $a^o = a^s$), it is conceptually convenient to assume that people evaluate the ability of others relative to their own abilities. Specifically, $\delta = a^o - a^s$ captures the *differential ability*, the amount by which the ability of others is different from one's own ability. Hence, we refer to this model as the *differentiated by ability* model.[1] As shown in Figure 2,

---

[1] Note that assessing differential ability $\delta$ and $a^s$ is equivalent to separately assessing $a^s$ and $a^o$.

**Figure 2**
*Graphical Models Connecting Self- and Other-Assessment*



*Note.* Schematic graphical models connecting the subjective estimates of self and other, corresponding to different substantive assumptions about the psychological process of other assessment: (a) differentiated by ability model ($M_1$), which is equivalent to the full hierarchical model; (b) fully differentiated model ($M_2$), which ignores population-level information; and (c) undifferentiated model ($M_3$), which ignores the individual-specific level of the full framework. See the online article for the color version of this figure.

this model considers inference at all three levels: population, specific individuals, and knowledge signals. As more information becomes available via external knowledge signals such as performance feedback, it is possible to learn whether the other person is better ($\delta > 0$) or worse ($\delta < 0$) relative to themselves.

Additionally, it assumes that estimates of perceived difficulty of the problem ($d$) are the same for both self and the other person. Hence, the participant uses their perceived item difficulty when estimating the other person's score on the same task. This is a key feature of the model. In contrast to the next model ($M_2$), it allows a person to draw meaningful insights from their experience with the task. When predicting the other's score for a target problem, the prediction can be informed by information gained about differential ability from previous problems and the participant's own perceived problem difficulty for the target problem. Therefore, this model predicts correlated scores between self- and other-estimated scores.

An equivalent formulation of the differentiated by ability model is a "differentiated by difficulty" model. Intuitively, the differentiated by difficulty model suggests that people assume equal ability for self and others but would experience the same task as having different difficulties. Due to the interconnected relationship between the ability and difficulty parameters, the two models would make similar predictions.

### Fully Differentiated Model ($M_2$)

This model assumes that other-assessment is not informed by any self-assessed estimates, consistent with a *fully differentiated* model of the other. As shown in Figure 2, this model assumes that inference about self and others is disjointed. As a consequence, there is no information sharing at the individual level. The fully differentiated model suggests that people draw no information from their own experience with the task when reasoning about another person. According to this model, in the absence of feedback, the participant possesses no meaningful information that can be used to inform predictions of the other person's performance. The participant starts with arbitrary priors about the other person's ability and perceived item difficulty and proceeds to learn about the other by

solely observing their scores (in the feedback condition) and ignoring any insights from their own experience. As more observations become available over time, the estimated other ability can be updated and can inform the prediction for the next set of problems. Note that, because people do not rely on their experience with the task to assess the other person, this model does not allow the person to learn any meaningful estimates of difficulty as experienced by the other person. Both ability and difficulty estimates of the other are evaluated independently of the ability and difficulty estimates of the self.

### Undifferentiated Model ($M_3$)

The last model assumes that the predicted other scores are highly constrained as the process of other-assessment uses the same information as the process used for self-assessment. As shown in Figure 2, this formulation ignores inference at the specific individual or the population levels of the proposed hierarchical framework. Therefore, this model suggests that people rely only on their assessment of themselves to make predictions about the other person. Overall, this model predicts no differentiation in ability as more information about the other person becomes available.

### Overview of Experiments and Modeling

Up to this point, we have explained the hierarchical framework and model variants primarily at a conceptual level. In the next sections, we will apply the framework to specific empirical paradigms. First, we will describe an empirical paradigm based on an image classification task where participants sequentially make predictions about the performance of themselves as well as the performance of another person. We evaluate how the self- and other-predictions differentiate over time as more information about the other person becomes available and test which of the three instantiations of the hierarchical model best accounts for the observed data. Second, we will use the hierarchical model to account for previous empirical findings about other assessment in tasks where no specific information about the other person is available

and participants reason about the other person and relative placement in the population using a combination of internal and external knowledge signals.
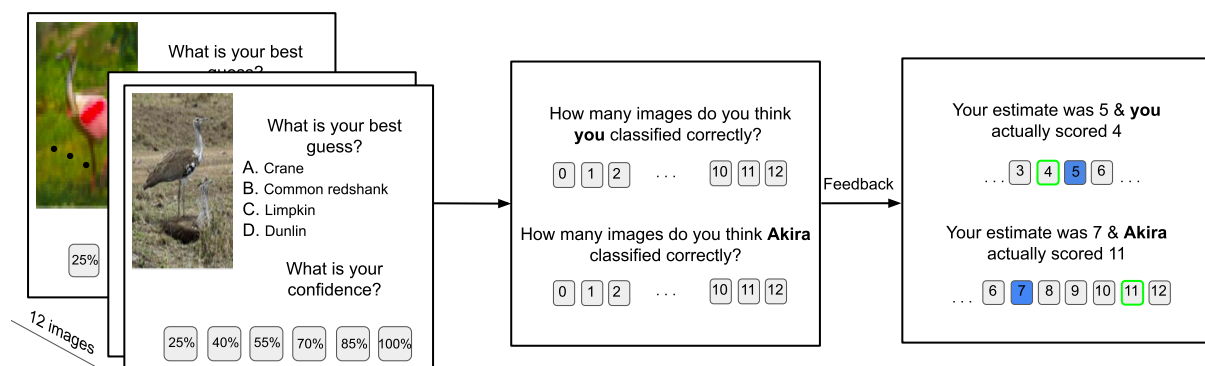
## A Sequential Knowledge Assessment Task

We develop an empirical paradigm similar to observer paradigms (Jameson et al., 1993) where there are multiple rounds of assessing one's own performance as well as the performance of another target person, allowing people to update their mental models of the target person. In this empirical paradigm, participants go through a series of problem sets (See Appendix C for details), where each problem set consists of a series of classification problems involving images of different species of animals (see Figure 3 for examples). After each problem set, participants self-assess their own performance ("How many items do you think you answered correctly?") as well as the performance of a target person who previously performed the task ("How many items do you think Akira answered correctly?"). The target person is referenced with a made-up name, but the associated data are based on an actual person who performed the experiment. In the no-feedback condition of the experiment, no information is provided about the actual performance of the target person, and the assessment is based on a priori predictions. In the feedback condition, the performance of the target person can be used by the participant to update their mental model of the other person's ability. In the example in Figure 3, when the participant is predicting how many items Akira answered correctly in the first problem set (involving birds), no feedback has been presented yet. However, after learning that Akira answered nine out of 12 items correctly while the participant themselves answered only seven items correctly, this provides an opportunity for the participants to adjust their mental model of the other person. This differentiated mental model can then be applied in the assessment phase for the second classification problem set (dogs) and further refined after receiving feedback. We apply an instantiation of the proposed framework to behavioral data collected via the sequential knowledge assessment

task, extending the work by Jansen et al. (2021) on other-assessment. We assume that other-assessment proceeds in a similar fashion as self-assessment by combining a subjective estimate for the perceived ability of the other person with estimates of the perceived difficulty for the other person. We use this framework to assess the degree of differentiation between the mental model of self (containing ability and problem difficulty estimates for self) and the mental model of others (containing ability and problem difficulty estimates for the other person). Consistent with previous research that has shown that one's own perceived difficulty in retrieving information or solving problems can be used to predict the difficulty experienced by others (Jameson et al., 1993; Kelley & Jacoby, 1996; Nickerson, 1999; Nickerson et al., 1987), we show that the subjective estimates of problem difficulty are shared between the self and other mental models. In addition, we show that the other person model differentiates from the self model based on differences in perceived ability. As information becomes available about the other person's performance, the differential ability can be updated, leading a person to upgrade or downgrade the predictions relative to their own ability.

## Notation

Before describing the computational model, we introduce some notation and define the scope of the model. In our empirical paradigm, each person $i$ is paired with a single other person. That is, each person reasons about their own performance and one other person's performance throughout the experiment. Therefore, we will omit from the notation which specific other individual person $i$ the self is reasoning about. We instead use the superscripts $s$ (self) and $o$ (other) to denote both the true scores of a person or of the assigned other person and subjective estimates of a person about their own or the other person's performance, respectively. We will use subscript $j$ to index the problem set, where $j \in \{1, \ldots, L\}$.

## Figure 3
*Illustration of the Empirical Paradigm for Self- and Other-Assessment*



*Note.* Participants go through a series of classification problem sets requiring participants to discriminate between different types of animals in a four-alternative forced-choice task. After classifying 12 images that constitute a problem set, participants proceed to the assessment phase, where they estimate the number of items they and another person answered correctly. The assessment phase is followed by feedback (if provided) on the actual number of items answered correctly. Numbers in blue and green show estimates and true scores, respectively. The scores of the other (target) person are based on selected participants who previously went through the experiment. A number of different names, including Akira, are used to reference the other person. See the online article for the color version of this figure.

For example, $x_{i,j}^s$ represents the number of items person $i$ answered correctly in problem set $j$, and $x_{i,j}^o$ represents the number of items answered correctly in problem set $j$ by the other person paired with $i$. $\hat{x}_{i,j}^s$ represents the number of items person $i$ *estimates* they answered correctly on problem set $j$. Similarly, $\hat{x}_{i,j}^o$ represents the estimated performance of the other person from the viewpoint of person $i$, that is, how many items person $i$ *believes* the other person answered correctly for problem set $j$. Both true and estimated scores are limited to the number of classification items ($M$) within each set, $x_{i,j} \in \{0, \ldots, M\}, \hat{x}_{i,j} \in \{0, \ldots, M\}$, where $M = 12$ throughout our experiments. In the empirical paradigm, the order in which the problem sets are presented varies across participants. We will use subscript $t = 1, 2, \ldots, T$ to refer to the order in which problem sets are presented and $j$ to refer to the specific type of problem set. For example, the bird problem set in Figure 3 could correspond to $t = 1$ and (say) $j = 4$. For person $i$ in this particular example and for $t = 1$, the number of estimated and true self and other answered correctly are $\hat{x}_{i,t}^s = 5, \hat{x}_{i,t}^o = 7, x_{i,t}^s = 4, x_{i,t}^o = 11$, with $M = 12$.

## Modeling Actual Performance

To formalize actual performance, we start with a model from IRT (Fox, 2010; van der Linden & Hambleton, 2013), which accounts for the observed performance differences across people and problem sets. The IRT model will also form the basis for the two other parts of the model (self- and other-assessment). To simplify the application of the IRT model across the three parts, we will use a basic Rasch model (Rasch, 1960) extended for ordered polytomous categories (i.e., the responses $x \in \{0, \ldots, M\}$). The key assumption of the Rasch modeling approach is that the number of items answered correctly, $x_{i,j}$ for person $i$ and problem $j$, is modeled by combining two latent factors, the ability $a_i$ of each person $i$ and the difficulty $d_j$ for problem set $j$:

$$\theta_{i,j} = a_i - d_j$$
$$p_{i,j} = \frac{1}{1 + \exp(-\theta_{i,j})}$$
$$x_{i,j} \sim \text{OrderedProbit}(p_{i,j}, v, \sigma). \quad (1)$$

Note that $a_i$ and $d_j$ represent the objective ability of person $i$ and the objective difficulty of problem $j$ measured using the IRT model. $\theta_{i,j}$ represents the latent score of person $i$ on problem set $j$ on a logit scale ($-\infty < \theta < \infty$), which is modeled as a sum of $a_i$, the ability of person $i$, and $d_j$, the difficulty for problem set $j$. Therefore, a higher score is expected for people with high ability or problems with low difficulty. The variable $p_{i,j}$ represents the latent score for person $i$ and problem set $j$ is converted to a value between 0 and 1. The ordered probit model[2] is a simple probabilistic process that maps the latent score $p_{i,j}$ to a discrete score, $x_{i,j} \in \{0, \ldots, M\}$. In this process, normally distributed noise with zero mean and standard deviation $\sigma$ is added to the latent score $p_{i,j}$ and the placement of the resulting value in a set of intervals (defined by the cutoff points $v$) determines the observed score. The variable $\sigma$ represents the uncertainty in mapping from latent to observed scores (see Appendix for details).

In this particular model, we have assumed that ability is one-dimensional—all variations in ability can be characterized by changes along a single overall ability scale. We could also consider multidimensional extensions of this model, analogous to multidimensional item response theory (MIRT; Reckase, 2009) that allows for differences in ability along a number of dimensions.

## Modeling Self-Assessment

For the self-assessment model, we assume that each person $i$'s estimate of their own ability $a_i^s$ and estimate of the problem difficulty for problem set $j$, $d_{i,j}^s$, are noisy and distorted versions of the true values. Both $a_i^s$ and $d_{i,j}^s$ may be interpreted as subjective estimates made by each person $i$ on problem $j$. These subjective estimates are related to the objective measures of ability ($a_i$) and difficulty ($d_j$) from Equation 1 according to:

$$a_i^s \sim N(a_i, \sigma_{a,i})$$
$$d_{i,j}^s \sim N(\gamma d_j + \lambda, \sigma_{d,i}), \quad (2)$$

where $\gamma$ and $\lambda$ parameter are scaling parameters that can capture systematic deviations of people's estimates from the true values of difficulty ($d_j$). Specifically, when $\lambda > 0$, problem difficulty will be overestimated leading to underestimates of scores. Similarly, when $\lambda < 0$, problem difficulty will be underestimated leading to overestimates of scores. The linear transformation of the problem difficulty is similar to the linear in log-odds models that have been used to model distortions in probability estimation in a variety of cognitive tasks (Turner et al., 2014; Zhang & Maloney, 2012).

An estimated score $\hat{x}_{i,j}^s$ by person $i$ for problem set $j$ is produced by combining the self-estimated ability and problem difficulty by following the same general process as in Equation 1:

$$\theta_{i,j}^s = a_i^s - d_{i,j}^s$$
$$p_{i,j}^s = \frac{1}{1 + \exp(-\theta_{i,j}^s)}$$
$$\hat{x}_{i,j}^s \sim \text{OrderedProbit}(p_{i,j}^s, v, \sigma^s). \quad (3)$$

Overall, there are two sources of noise that can produce distortions in self-estimation. The subjective ability might not reflect the true ability, and the subjective problem difficulty might systematically deviate from the actual problem difficulty.

Note that the self-assessment model in Equations 2 and 3 is similar to the IRT model in Equation 1 but that it plays a very different role in our approach conceptually. The IRT model in Equation 1 serves the purpose of a data analysis model to estimate the true abilities and true item difficulties, whereas the self-assessment model in Equations 2 and 3 formulates a cognitive model to explain the process of self-assessment. We use the ordered probit model as a link function to map a person's subjective latent probability of being correct, $p_{i,j}^s$, to a score between 0 and 12. However, as we will show in a later section of the article, we may easily modify this to accommodate cases where different knowledge signals are available (e.g., feeling of knowing or response time).

---

[2] There are alternative generative models for ordered responses including the graded response model (Greene & Hensher, 2010). We have found that the use of this alternative construction does not change the qualitative results.

## Modeling Other-Assessment

For this model, we make the assumption that the way people reason about the other person's performance is through the lens of their own self-assessment process. That is, once a person $i$ has an estimate of the ability of the other person ($a_i^o$) and an estimate of the problem difficulty for problem set $j$ as experienced by the other person ($d_{i,j}^o$), we assume that scores for the other person can be predicted by applying the same cognitive model as Equation 3:

$$\theta_{i,j}^o = a_i^o - d_{i,j}^o$$
$$p_{i,j}^o = \frac{1}{1 + \exp(-\theta_{i,j}^o)}$$
$$\hat{x}_{i,j}^o \sim \text{OrderedProbit}(p_{i,j}^o, v, \sigma^s). \tag{4}$$

Note that in this model, $a_i^o$ and $d_{i,j}^o$ are not the true ability and problem difficulty of the other. Instead, they represent $i$'s estimate of the true ability of other and the estimate of the difficulty for the other.

## Hypotheses About the Relationship Between the Self and Other Model

Now that the basic models for self- and other-assessment have been formalized, we specify how the three hypotheses, the differentiated by ability ($M_1$), fully differentiated ($M_2$), and undifferentiated model ($M_3$), translate to different computational assumptions about how the estimates of the other ability and problem difficulty are formed. The underlying computational assumptions of the three hypotheses are summarized in Table 1 in terms of the notation above. Note that these relationships describe different *beliefs* held by the person making inferences about the other person. In other words, these are psychological assumptions about how people use available information to draw inferences in their cognitive model of the other person.

## $M_1$: Differentiated by Ability Model

The differentiated by ability model ($M_1$) assumes that for each type of problem set $j$, the difficulty for another person is the same as the difficulty for one's self (i.e., $d_{i,j}^o = d_{i,j}^s$). However, it allows for the possibility that there is a difference, $\delta_i$ in ability between self and other from the viewpoint of person $i$. This differential ability is inferred as information about the performance of the other person becomes available over time.

The inference process can be stated as a sequential updating problem. After $t$ problem sets, the person $i$ has received information about the other person's performance $x_{i,1}^o, \ldots, x_{i,t}^o$ (e.g., if after $t = 3$ rounds of problem sets, the other person scored 11, 7, and 8 correct out of 12, we have $x_{i,1}^o = 11$, $x_{i,2}^o = 7$, and $x_{i,3}^o = 8$). On the basis of this information, a prediction for the performance on the next problem set, $x_{i,t+1}^o$, can be made by first making an inference about the differential ability $\delta_i$ from the viewpoint of person $i$:

### Table 1

*Model-Based Hypotheses About the Relationship Between Self- and Other-Mental Model Parameters*

| Model | Hypothesized dependencies | |
| --- | --- | --- |
| | $a_i^o$ and $a_i^s$ | $d_{i,j}^o$ and $d_{i,j}^s$ |
| $M_1$: Differentiated by ability | $a_i^o = a_i^s + \delta_i$ | $d_{i,j}^o = d_{i,j}^s$ |
| $M_2$: Fully differentiated | unrelated | unrelated |
| $M_3$: Undifferentiated | $a_i^o = a_i^s$ | $d_{i,j}^o = d_{i,j}^s$ |

*Note.* Each hypothesis is associated with a different cognitive model for other-assessment. $M_1$ = differentiated by ability model; $M_2$ = fully differentiated model; $M_3$ = undifferentiated model.

$$p(\delta_i | x_{i,1}^o, \ldots, x_{i,t}^o) \propto p(x_{i,1}^o, \ldots, x_{i,t}^o | \delta_i, d_{i,1}^s, \ldots, d_{i,t}^s) p(\delta_i)$$
$$= \left( \prod_{\tau=1}^{t} p(x_{i,\tau}^o | \delta_i, d_{i,\tau}^s) \right) p(\delta_i). \tag{5}$$

Note that the second line follows from the first because of conditional independence. The term in the product can be evaluated by Equation 4 by using the model assumption $a_i^o = a_i^s + \delta_i$. In the next step, on the basis of the posterior estimates of $a_i^o$, the score of the other person for the next problem set presented at time $t + 1$, $p(x_{i,t+1}^o | a_i^o, d_{i,t+1}^o)$, can be predicted by applying Equation 4. Here, $d_{i,t+1}^o$ is the same difficulty as inferred by the self using the self-assessment model ($d_{i,t+1}^s$). The term $p(\delta_i)$ reflects the person $i$'s prior about the differential ability. We assume that this prior is centered around zero, such that at the start of learning, the mental model of self and others are undifferentiated.

## $M_2$: Fully Differentiated Model

The most unconstrained of the three hypotheses is the fully differentiated model ($M_2$). In this model, the estimates in the mental self model are unrelated to the estimates in the mental other model (i.e., $a_i^o$ is unrelated to $a_i^s$ and $d_{i,j}^o$ is unrelated to $d_{i,j}^s$). This model posits that people use no insights from their experience with the task when assessing the other person.

A prediction for the performance on the next problem set $t + 1$, $\hat{x}_{t+1}^o$, can be made by making an inference about the ability of the other person ($a_i^o$) and difficulty for the other person ($d_{i,1}^o, \ldots, d_{i,t}^o$):
(See below)

The terms $p(a_i^o)$ and $p(d_i^o)$ reflect a person's priors about the other person and we have assumed independence between these priors. Note that the second line follows from the first because of conditional independence. The score of the other person for the next problem set, $p(x_{i,t+1}^o | a_i^o, d_{i,t+1}^o)$, can be predicted by applying Equation 4 to the posterior estimates of $a_i^o$ and drawing a sample from the posterior of $d_i^o$.

$$p(a_i^o, d_{i,1}^o, \ldots, d_{i,t}^o | x_{i,1}^o, \ldots, x_{i,t}^o) \propto p(x_{i,1}^o, \ldots, x_{i,t}^o | a_i^o, d_{i,1}^o, \ldots, d_{i,t}^o) p(d_{i,1}^o, \ldots, d_{i,t}^o) p(a_i^o)$$
$$= \left( \prod_{\tau=1}^{t} p(x_{i,\tau}^o | a_i^o, d_{i,\tau}^o) p(d_{i,\tau}^o) \right) p(a_i^o). \tag{6}$$

Note that the flexibility of this other-assessment model allows for the possibility that a problem set has differing levels of difficulty across people. When the same type of problem set occurs over time, this model will allow a person to potentially make accurate predictions for the other person's performance. However, in an environment where problem sets do not repeat (as in our empirical paradigm), this model does not generalize well as the information acquired for each type of problem set is not utilized in the future.

### $M_3$: Undifferentiated Model

The most constrained of the three models is the undifferentiated model ($M_3$). In this model, the mental models of self and others are the same and remain undifferentiated as new information becomes available about the performance of the other individual. Therefore, the process for producing predictions for the problem set presented at time $t$ for self ($\hat{x}_{i,t}^s$) and other ($\hat{x}_{i,t}^o$) in Equations 3 and 4 is based on the same parameters. Note that in this model, the predicted self and other scores can still deviate from each other because of the noise process of producing discrete scores in Equations 3 and 4.

## Experiments

We conduct two image classification experiments to investigate self- and other-assessment and develop and test the computational models. In Experiment 1, we collect behavioral data from 68 participants on the basic experimental paradigm that only includes self-assessment. Experiment 2 follows the same experimental paradigm but also includes other-assessment of participants from Experiment 1. There were 128 individuals in total serving as "self" in Experiment 2. Specifically, the best- and worst-performing 16 participants from Experiment 1 served as the "other" individuals that participants in Experiment 2 are learning about.

## Method

### Participants

Participants were recruited through Amazon Mechanical Turk. Sixty-eight and 128 participants were recruited for Experiment 1 and Experiment 2, respectively. To be eligible for the studies, participants were required to meet the following criteria: (a) have greater than or equal to 80% human intelligence task (HIT) approval rate for all requesters' HITs; (b) be located in the United States; and (c) be 18-years-old or older. All participants provided informed consent before taking part in our study and were compensated $6 for their participation. The median time to complete the experiment was 33 min.

### Images

There were 192 unique images in total used in the experiments, divided equally into four categories (birds, dogs, primates, and reptiles). Each category was associated with $T = 4 \times 4 = 16$ problem sets in total, with each problem set containing $M = 12$ individual classification problems. In each classification instance, the goal is to classify images according to four different labels corresponding to a specific category. For example, for one of the bird problem sets, the labels are *crane*, *common redshank*, *limpkin*, *dunlin*, and for one of

the dog problem sets, the labels are *Afghan hound*, *Ibiza hound*, *Norwegian elkhound*, *redbone coonhound* (see Appendix A for a list of the 16 classification problem sets). The images and labels for the classification problems are based on the ImageNet Large Scale Visual Recognition Challenge 2012 database (Russakovsky et al., 2015). ImageNet is an image data set where the labels for each image are hierarchically organized according to the WordNet hierarchy (Miller, 1995). We selected 16 classification problem sets equally divided among the four categories. For each classification problem set, we randomly selected 12 images (three images per label) from the validation set of ImageNet. Each image was center cropped and scaled to $256 \times 256$ pixels.

### Procedure

In both Experiments 1 and 2, participants went through 16 problemmsets where each problem set included 12 classification problems of a particular category as well as a prediction task where participants assessed their own performance (Experiments 1 and 2) and also assessed another person's performance (Experiment 2 only). For each problem set, a participant first classified 12 individual images (Figure 3). For each image, the participant selected a label from four response alternatives (e.g., *little blue heron*, *oystercatcher*, *dowitcher*, and *great egret*). The response alternatives remained the same during each problem set. The participant also selected a discrete confidence level from six alternatives (25%, 40%, 55%, 70%, 85%, and 100% confidence). The 25% and 100% confidence levels had additional text labels "Guessing" and "Absolutely Certain," respectively. No feedback was provided during this classification phase. The confidence ratings and individual classifications were not used for the purpose of this research.

At the end of each problem set, the 12 images from the preceding classification task were presented simultaneously on the screen. In both Experiments 1 and 2, participants were instructed to predict the number of images they classified correctly by selecting a response option between 0 and 12 (self-assessment). In Experiment 2, they were also asked to predict the performance of another person by selecting a number between 0 and 12 (other-assessment). This person was referred to by a name, sampled randomly from a set of seven male and seven female names (e.g., "Vince," "Glenda"). The participant was told that this was not the real name of the other person but that the other person was an actual person who participated previously in the experiment (the same name was used throughout the experiment).

In Experiment 1, after the predictions were made for each problem set $t$, participants were provided feedback and were told the actual number of correct responses (e.g., "You classified 8 out of 12 images correctly"). Participants were given an option to see which individual images they classified incorrectly. The correct label was not shown. After this feedback, participants proceeded to the next problem set $t + 1$. In Experiment 2, in the feedback condition, feedback was provided about the number of correct self- as well as other-responses (e.g., "Vince scored 6 out of 12 images correctly"). In the no-feedback condition, this feedback about self- or other-performance was omitted.

Overall, each participant provided 192 image classifications with corresponding confidence levels and provided 16 predictions about their performance across 16 different types of classification problem sets.

## Design

The 16 best-performing and 16 worst-performing participants from Experiment 1 served as the other person to learn about in Experiment 2. We will refer to these two groups of other people as top and bottom, respectively. In the feedback condition, a participant in Experiment 2 received feedback about the particular other person assigned to the participant. In the no-feedback condition, no such information was provided. The assignment of the 16 top and 16 bottom participants from Experiment 1 to the 128 participants in Experiment 2 was counterbalanced across the two feedback conditions—each target participant from Experiment 1 was assigned to exactly four participants in Experiment 2, two in the feedback and two in the no-feedback conditions. This study was not preregistered. All data sets analyzed in this work can be accessed from https://osf.io/68347/.

## Metrics for Assessment Performance

For both self-assessment and other-assessment, we report results based on three different metrics to provide a more comprehensive picture of assessment performance (Dunning & Helzer, 2014). Note that because our assessment task of estimating the number of items scored correctly does not relate to a binary detection task, various standard metacognition measures such as metacognitive sensitivity and efficiency (Fleming & Lau, 2014) cannot be applied.

The first metric is the coefficient of predictive ability (CPA; Gneiting & Walz, 2021), a rank-based measure that generalizes the area under the curve (AUC) to ordinal and continuous variables (for details, see Appendix D). In our context, the CPA evaluates how well people can discriminate in their assessment between different true scores. More specifically, the CPA is a weighted probability that under random sampling of problem sets, a problem set with a higher true score is self-assessed with a higher score than a problem set with a lower true score.[3] The weights in CPA are based on the distance between the ranks of the true scores. Therefore, a person who is able to assign different scores to closely ranked true scores will achieve a higher CPA. The CPA measure is theoretically appropriate for a number of reasons: The CPA is equivalent to AUC when applied to binary outcomes and is equivalent to Kendall's tau rank-order correlation when there are no ties in the true scores. It is also closely related to Goodman Kruskal's $\gamma$ coefficient that has been used to assess metacognitive sensitivity (Nelson, 1984). Because of the rank-based nature, CPA is insensitive to bias. Any changes to the estimated scores that preserve ranking will result in the same CPA. The CPA attains values between 0 and 1. A value of 1 is attained when there is a perfect correspondence between estimated and true scores. A value of 1/2 is attained when the estimated scores are independent of the true scores.

Second, we report a bias measure to measure the systematic deviations between the true and estimated score, defined as Bias = $(1/N) \sum_{i=1}^{N} (\hat{x}_i - \bar{x})$, where $\hat{x}$ is the estimated score through self- or other-assessment and $\bar{x}$ is the mean of true scores across problem sets. If the assessment scores are consistently overestimating or underestimating the true performance, the bias score will be positive and negative, respectively.

Third, to facilitate comparison to previously reported results on assessment (e.g., Zell & Krizan, 2014), we also report the Pearson correlation coefficient ($\rho$) between the true and estimated scores.

## Model Inference

We used Markov chain Monte Carlo (MCMC) sampling to infer model parameters for the cognitive models presented in Figure B1 and obtain samples from the posterior distribution. We chose the Stan computing environment for posterior inference (Stan Development Team, 2020). Model inference proceeds in a sequential fashion. We begin with actual performance assessment, followed by self-assessment, and finally other-assessment. We start by estimating the parameters ($a$, $d$, $\sigma$) that account for the actual performance of the participants using the true scores $x^s$. These parameters were estimated using a standard one-parameter IRT model described in A Hierarchical Framework for Knowledge Assessment section on modeling actual performance. In the next stage of our inference, we treat the posterior means of $a$, $d$, $\sigma$ as observed data to infer the parameters of our self-assessment model ($a^s$, $d^s$, $\sigma^{a,i}$, $\sigma^{d,i}$, $\sigma^s$, $\lambda$, $\gamma$) using participant's estimates of their true scores ($\hat{x}^s$). Inference on the self-assessment model gives us the estimated perceived ability of self ($a^s$) and perceived difficulty of items ($d^s$) for every individual. We ignore learning over time when estimating these self-assessment parameters as we did not observe any such learning in our empirical data. Finally, the posterior means of the parameters from the self-assessment model serve as the starting point for the other-assessment models.
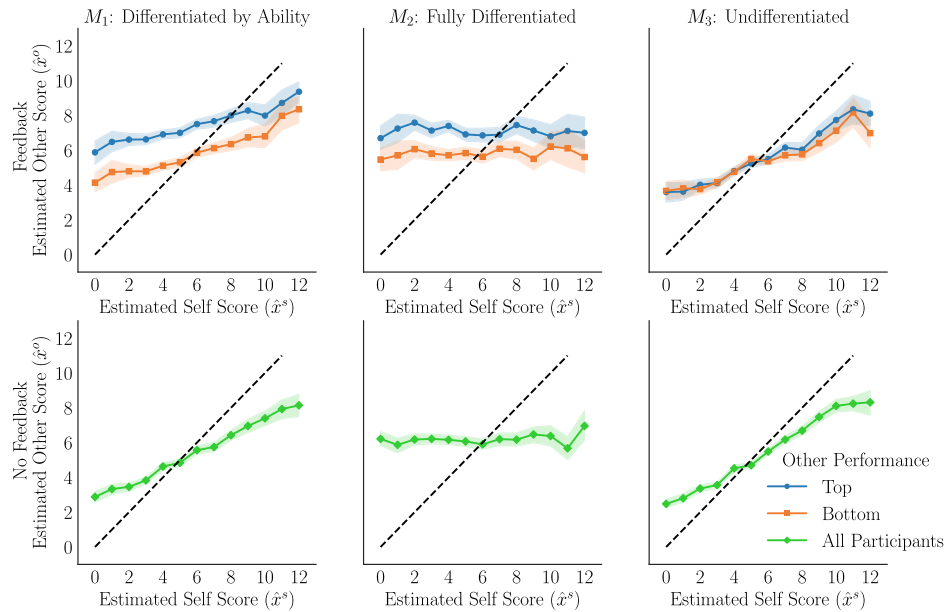
We use the three variants of the other-assessment model to simulate participants' estimates of the other person's scores. To do inference, we condition on $a^s$, $d^s$, $\sigma^s$, and $x^o$. Figure B1 shows the graphical models corresponding to each model variant. At the first time step, depending on the variant of the other-assessment model, we either use priors for $a^o$ and $d^o$ ($M_3$) or values of $a^s$ and $d^s$ ($M_1$, $M_3$) to predict the participant's first estimate of the other person's performance (here, the participant has not received any information about the other person). At each subsequent time step, participants may learn about the other person in the feedback condition. Simulating from the undifferentiated model ($M_3$) requires no learning: We simply use self estimates ($a^s$ and $d^s$) to predict the participant's estimated scores of the other person on each time step. To simulate the participant's estimates using the fully differentiated model ($M_2$), we use the mean posterior estimates of $a^o$ and $d^o$ from the previous time step to predict the estimated scores of the other person. For the differentiated-by-ability model ($M_1$), we use the mean posterior estimates of $a^o$ from the previous time step and $d^s$ for the current item to predict the participant's estimated score of the other person $\hat{x}^o$ (See Appendix B for details).

Our experimental and modeling setup allows us to simulate a participant's estimate of any other person's score, that is, we can use a participant's inferred self-ability and item difficulties from the self-assessment model to predict their estimates of any randomly picked other person's scores. For Figures 4 and 5, we increased the number of simulated other-assessments fourfold in order to more clearly visualize the differences in model predictions from the three different linkage hypotheses. In these simulations, for every participant, we simulate their other assessment separately for four randomly assigned participants as their "other persons." We then use the other-assessment procedure described above to make predictions about the participant's estimates of the new others' scores.

---

[3] Ties between the self-assessed scores are resolved at random.

**Figure 4**

*Model Predictions for the Relationship Between Estimated Other Score and Estimated Self-Performance*



*Note.* The results are separated by the feedback condition and performance levels of the other person. Note that in the no-feedback condition, participants cannot differentiate between top and bottom performers. Dashed line indicates exact equivalence between estimated self and other scores. The colored areas show 95% confidence intervals. See the online article for the color version of this figure.

## Empirical Results

### Classification Performance

Participants substantially differed in overall performance. From the worst to the best-performing participant, the mean proportion correct varied between 33% and 81% across Experiments 1 and 2. Classification performance improved slightly within each problem set. Across the first, middle, and last four classification items in a problem set, average performance was 53%, 55%, and 57%, respectively. This improvement is likely due to participant strategies of adjusting their classifications after seeing a larger range of images. Across problem sets, no apparent learning took place

Implementing the IRT model requires careful attention to the selection of priors on both ability and difficulty to avoid potential identifiability issues. For the actual performance model, we used normal priors of ability and difficulty IRT parameters: $a_i \sim \mathcal{N}(0,1)$, $d_j \sim \mathcal{N}(\mu_d, \sigma_d)$, where $\mu_d \sim \mathcal{N}(0,1)$, $\sigma_d \sim$ Cauchy $(0, 5)$. Additionally, for the self-assessment model, we used normal priors for $\lambda \sim \mathcal{N}(0,1)$, $\gamma \sim \mathcal{N}(0,1)$ and Cauchy priors for standard deviation parameters $\sigma_{a,i}$, $\sigma_{d,i} \sim$ Cauchy $(0, 5)$. Finally, for the differentiated-by-ability model, we use a normal prior on $\delta_i \sim \mathcal{N}(\mu_\delta, \sigma_\delta)$ where $\mu_\delta \sim \mathcal{N}(0,1)$ and $\sigma_\delta \sim$ Cauchy $(0, 5)$. Throughout the inference process, we ran the sampler with two chains with a burnin of 1,000 iterations before taking 1,000 samples per chain. The chains mixed appropriately based on Rhat values (close to 1). Stan code for self- and other-assessment models can be accessed from https://osf.io/68347/.

(keep in mind that each problem set involved new classification problems with a unique set of labels). The average accuracy grouped by four consecutive problem sets was 56%, 56%, 53%, and 55%.
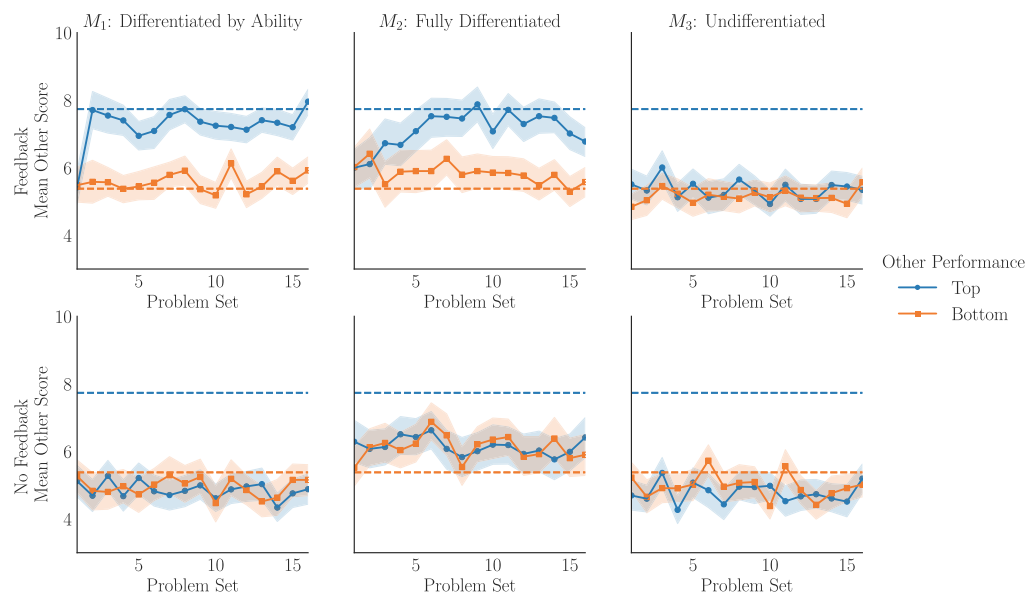
## Assessment Performance

While many metrics have been introduced to evaluate metacognition, they are typically applied to binary decision tasks (Fleming & Lau, 2014). Given that the self- and other estimated and true scores are based on discrete counts with more than two outcomes, we adopt a relatively new measure, the CPA (Gneiting & Walz, 2021), to assess metacognitive sensitivity, the ability to discriminate between different true scores.

Table 2 shows the self- and other-assessment performance based on CPA as well as bias (see the Method section for details), and Pearson correlation coefficient ($\rho$) between true and estimated scores.[4] According to the CPA as well as the Pearson correlation, participants' self- and other-assessment is well above chance level (note that chance level for CPA is 0.5). For self-assessment, the Pearson correlation coefficients are in the 0.5–0.7 range, which is well above the 0.2–0.3 range reported for many other self-assessment tasks (Zell & Krizan, 2014).

---

[4] Although not used in this research, we also collected confidence scores for individual questions in each problem set. There is a close correspondence between the mean of the estimated probabilities across items and the estimated score (.78 [$SD$ = .14] for Experiment 1 and .82 [$SD$ = .13] for Experiment 2), suggesting that participants' estimates are based on aggregates of individual confidence scores.

**Figure 5**
*Model Predictions for the Mean Estimated Score of the Other Person Over Problem Sets*



*Note.* The results are separated by the feedback condition and performance levels of the other person. Dashed lines show the mean true score across the top and bottom performing other people. The colored areas show 95% confidence intervals. See the online article for the color version of this figure.

Figure 6 shows the self-estimated score as a function of the true score for a particular problem set. The data for this analysis are combined across Experiments 1 and 2. The results show a small range of true scores associated with a pattern of overestimation. For a larger range of true scores, there was a pattern of underestimation. Generally, this pattern of systematic deviations is consistent with previous findings in self-assessment (Jansen et al., 2021; Kruger & Dunning, 1999) and is consistent with the general pattern of over- and underestimation in subjective assessment tasks (Zhang & Maloney, 2012). However, it is important to note that there were few problem sets where participants produced low true scores that are associated with the overestimation pattern (see the marginal distribution at the top of the figure). Overall, there was a tendency to underestimate performance, as revealed by the negative bias values in Table 2. Across Experiments 1 and 2, there were 169 participants with more under- than overestimates in the self-assessment and only 19 participants with more over- than underestimates.

Other-assessment is a more challenging task than self-assessment leading to somewhat lower performance. However, participants' accuracy in assessing other participants (i.e., the participants in Experiment 1) is not far off from the ability of those participants to predict their own performance (i.e., see self-assessment results from Experiment 1, top/bottom performers). Across participants, feedback improves other-assessment on all performance metrics including bias.[5]

Figure 7 demonstrates that individual participants are tracking the performance of other people in the feedback condition. In the feedback condition, when participants make predictions about the other person for the first problem set, no feedback has been provided yet and the results show that predictions are the same across the top and bottom other performers. However, the estimated mean scores

diverge within a few problem sets depending on the type of other person they are learning about. In the no-feedback condition, participants' estimated scores cannot (by definition) reflect differences between other people. Instead, without feedback, estimates have to be based on prior knowledge only. Generally, these prior predictions underestimate true performance (i.e., negative bias).

Finally, the other assessment shows patterns of over- and underestimation that are similar to self-assessment. Figure 6b shows that for particular problem sets that lead to low (high) true scores, participants tend to over (under) estimate performance. This pattern is similar across feedback conditions.

## Relationship Between Self- and Other-Assessment

Figure 8 shows that there is a close correspondence between self- and other-assessment. In the no-feedback condition, there is a strong tendency to link the estimate of the other score to the estimate of the self score, suggesting that when people believe a problem is challenging for themselves, they believe it is likely to be challenging for other people as well. In the feedback condition, the results show the same pattern but the predictions are differentiated by the type of other person they are learning about with higher predicted scores for a top performer. Therefore, in the feedback condition, the results suggest that two factors affect the other-

---

[5] At the individual participant level, discrimination (CPA) and correlation (C) are higher in the absence of feedback, which suggests that feedback lowers the ability to discriminate between different levels of performance. However, it should be noted that each participant in the feedback condition tracks the performance of either a top or bottom performing other person. Therefore, for those participants, there is a restricted range of scores to discriminate, which reduces CPA and C.

**Table 2**

*Self- and Other-Assessment Performance Across Experiments and Conditions*

| Type/condition | Across participants | | | Per participant | | | |
|---|---|---|---|---|---|---|---|
| | CPA | Bias | $\rho$ | *M* CPA | *M* bias | *M* $\rho$ | *N* |
| Self-assessment | | | | | | | |
| Exp. 1, feedback (all) | 0.75 | −1.41 | 0.52 | 0.79 (0.011) | −1.41 (0.19) | 0.62 (0.019) | 68 |
| Exp. 1, feedback (TB) | 0.75 | −1.24 | 0.53 | 0.80 (0.015) | −1.24 (0.30) | 0.62 (0.029) | 32 |
| Exp. 2, feedback | 0.82 | −1.41 | 0.65 | 0.82 (0.011) | −1.41 (0.14) | 0.64 (0.022) | 64 |
| Exp. 2, no feedback | 0.78 | −1.54 | 0.57 | 0.80 (0.009) | −1.54 (0.21) | 0.64 (0.018) | 64 |
| Other-assessment | | | | | | | |
| Exp. 2, feedback | 0.70 | −0.08 | 0.40 | 0.63 (0.013) | −0.08 (0.14) | 0.28 (0.027) | 64 |
| Exp. 2, no feedback | 0.63 | −0.60 | 0.27 | 0.69 (0.016) | −0.60 (0.26) | 0.41 (0.032) | 64 |

*Note.* For the analysis per participant, the statistics are calculated at the individual participant level and then averaged; numbers between parentheses are standard errors. *N* is the number of participants. For the analysis across participants, we ignore individual differences and report a single outcome across participants and problem sets. TB = subset of participants who were part of the top and bottom performers; CPA = coefficient of predictive ability; Exp. = experiment.

assessment, the estimated overall performance of the other person and the perceived problem difficulty.

## Discussion of Empirical Results

Our empirical results are consistent with the hypothesis that participants are developing and updating a mental model that allows them to make inferences about the overall level of performance of the other person. Figure 7 shows that participants' estimates of top and bottom other performers diverge within a couple of feedback rounds. This suggests that people employ an efficient mental representation of the other that enables them to quickly distinguish their own performance from the other person's.

Our results are consistent with previous studies of predicting general knowledge in self and others (Jameson et al., 1993). Target participants in Experiment 1 were more accurate in assessing
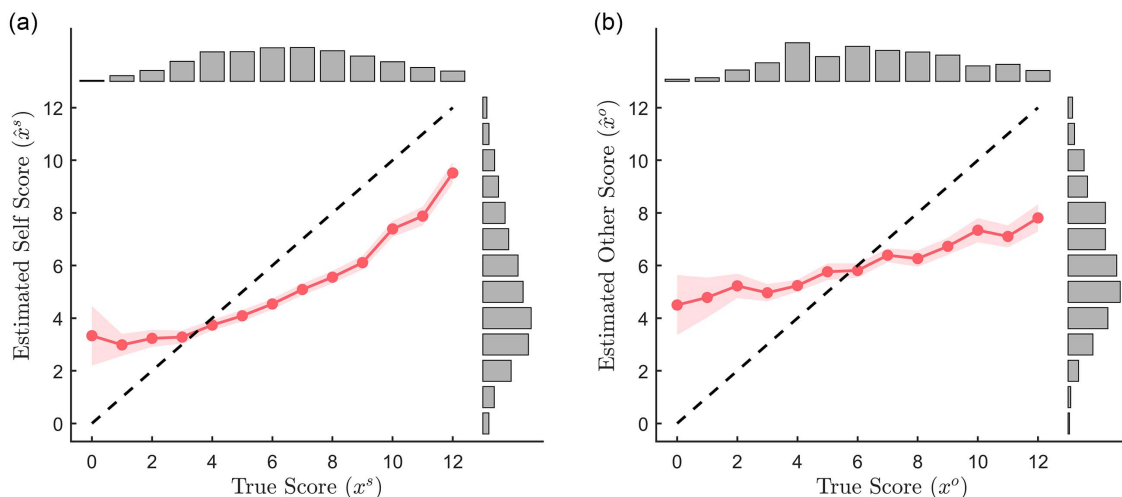
themselves than the observers in Experiment 2 who assessed the targets and received feedback. In turn, the observers who received feedback were more accurate than the observers who did not receive feedback. However, without feedback, performance is still well above chance. Figure 8 hints that observers without feedback use their own perceived ability and their self-assessed problem difficulty as predictors, assuming that what is difficult for them is also difficult for another person. This guessing strategy is effective in situations where the perceived problem difficulty for self correlates with the actual problem difficulty faced by other people (Fussell & Krauss, 1991; Jameson et al., 1993; Nickerson et al., 1987).

## Model-Based Results

Our primary modeling objective is to understand the mechanisms at play when humans make inferences about the ability and performance of other individuals. To do so, we simulate the three qualitatively
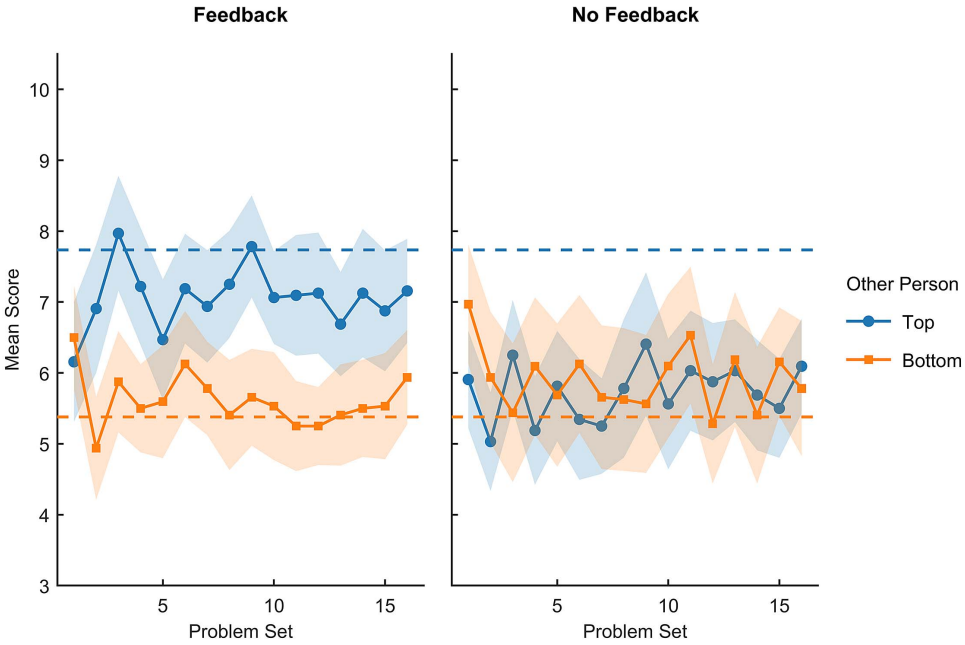
**Figure 6**

*Mean Estimated Self Score (a) and Other Score (b), Each as a Function of Actual Performance for a Particular Problem Set*



*Note.* For the self-scores, the data are combined across Experiments 1 and 2. Histograms show the marginal distribution of scores. The colored areas shows 95% confidence intervals. See the online article for the color version of this figure.

**Figure 7**

*Mean Estimated Score of the Other Person Across Feedback Conditions and Performance Levels of the Other Person*
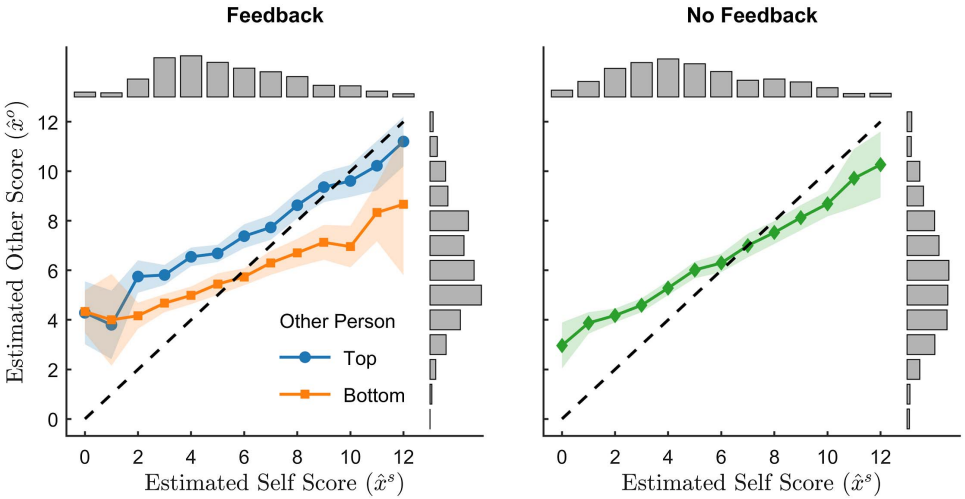


*Note.* Dashed lines show the mean true score across the top and bottom performing other people. Note that the no-feedback condition (right panel) shows the a priori predictions of participants. The colored areas show 95% confidence intervals. See the online article for the color version of this figure.

different models described above and relate them to the key empirical findings in our experiments. We use two methods to evaluate model adequacy. First, we perform a qualitative model evaluation by assessing the models' ability to replicate the qualitative patterns we observed in the empirical data. We do this through posterior predictive

simulation. For all three hypotheses, we use the existing behavioral data from the set of participants and problem sets to estimate posterior distributions of the parameters. We then simulate the behavior of new participants and new problem sets by sampling from the posterior predictive distribution (i.e., these are predictions for a replication of

**Figure 8**

*Estimated Score for the Other Person ($\hat{x}^o$) Conditional on the Estimated Self Score ($\hat{x}^s$)*



*Note.* The results for the feedback condition are separated by the overall performance of the other person. Histograms show the marginal distribution of scores. The colored areas show 95% confidence intervals. See the online article for the color version of this figure.

the experiment with a new set of participants and new problem sets). We use this simulated data to compare the qualitative predictions of our models to our empirical findings on (a) the relationship between self- and other-assessment, and (b) people's ability to differentiate between good and bad performances of other participants when given feedback. Our second method for model evaluation is through out-of-sample predictive checks using cross-validation. In this approach, we use the posterior distributions for the actual set of participants and problem sets in the experiments and compare the model predictions for held-out problem sets against the observed data.

### Relationship Between Self- and Other-Assessment

Previous investigations of neural activity during self- and other-assessment (Frith & Frith, 1999; Jenkins et al., 2008; Mitchell et al., 2005) have revealed a close correspondence between people's metacognition and their theory of mind. Our empirical results also indicate that self-assessment is closely tied to other-assessment. Figure 4 shows the relationship between self- and other-assessment as predicted by the three models. These results are based on a combination of experimental data and simulated data. We simulate participants' assessment of others' performance for four randomly assigned participants as their "other persons."

Compared to the observed empirical data in Figure 8, we see that the differentiated by ability model ($M_1$) most closely captures the trend observed in the empirical data in both the feedback and no-feedback conditions. When feedback is provided, it predicts a strong association between the self and other estimates while allowing for learning of differential ability of the other. This is consistent with what we see in our empirical data where people's estimates of their own performance are closely tied to their performance of the other. People draw on their experience with the task to make inferences about the other person's experience and assume that their subjective difficulty on any item must be commensurate to the difficulty experienced by the other person. Throughout the experiment, their estimates of the other person's performance are anchored by their own scores.

In contrast, without any informative priors about ability or difficulty, the fully differentiated model ($M_2$) fails to predict any association between self- and other-assessment. Alternatively, the undifferentiated model ($M_3$) relies too heavily on priors and predicts that people's estimates of others' performances are closely tied to their assessment of their own performance. Note that in the case of no feedback, $M_1$ is similar to $M_3$. With no information to learn from, people are forced to rely heavily on their own metacognitive

assessments of their ability and difficulty of each item as a prior for the other person. Hence, both models predict similar trends between self and other scores in the no-feedback condition.

### Differentiating Between Good and Bad Performers

In Figure 7, we observed that participants are able to distinguish between good and bad performances of other participants in the feedback condition. On the first trial, people use their prior beliefs about the other person's ability and difficulty to estimate others' scores. Subsequently, in the presence of feedback, people adjust their beliefs about the other participant's ability to make their estimates. The corresponding model predictions are shown in Figure 5. The results show that the differentiated by ability model ($M_1$) accurately emulates this behavioral pattern. The simulated participants' estimates of the good and bad performances diverge after they receive a single data point as feedback. On the other hand, while $M_2$ does better than $M_3$ at capturing the dependence of other-assessment on self-assessment (Figure 4), it does not capture people's ability to learn and differentiate between good and bad performances by the other. This is an important feature of the feedback condition in our experiment—people quickly learn the differential ability of the other person. Both $M_2$ and $M_3$ fail to capture this critical empirical feature.

### Quantitative Assessment of Model Performance

Table 3 shows how well each of the three models is able to capture the other-assessments in the empirical data. The sequential nature of our models allows us to make out-of-sample predictions for other-assessment at each time step. For example, when making a prediction at time $t + 1$, the model only receives information about the other person's true performance up to time $t$.

The table shows the mean squared error ($MSE$) and Pearson correlation ($\rho$) between the predicted estimates of other-performance as evaluated by the models and the actual estimates of other-performance made by participants in the experiment. These values indicate how closely model estimates resemble the true data. We only compare the models of their performance on the feedback condition. Overall, we see that the differentiated by ability model ($M_1$) outperforms the two other models ($M_2$ and $M_3$). This model provides the best quantitative fit to the data when the correspondence is assessed for each individual participant as well as across participants. Other statistics such as CPA follow the same trends as shown in Table 3 (see Appendix for details). We focused on $MSE$ because it is a standard way to evaluate the predictive performance of models.

**Table 3**
*Other-Assessment Across Models $M_1$, $M_2$, and $M_3$*

| Model | Across participants | | Per participant | | |
| --- | --- | --- | --- | --- | --- |
| | $MSE$ | $\rho$ | $M$ $MSE$ | $M$ $\rho$ | $N$ |
| $M_1$: Differentiated by ability | **8.92** | **0.39** | **8.92 [5.515, 12.324]** | **0.359 [0.241, 0.478]** | 64 |
| $M_2$: Fully differentiated | 15.95 | 0.15 | 15.95 [12.726, 19.172] | 0.076 [−0.067, 0.219] | 64 |
| $M_3$: Undifferentiated | 10.60 | 0.26 | 10.60 [7.254, 13.945] | 0.276 [0.137, 0.414] | 64 |

*Note.* For analysis per participant, the statistics are calculated at the individual participant level and then averaged; numbers between parentheses are 95% confidence intervals. $N$ is the number of participants. For the analysis across participants, we ignore individual differences and report a single outcome across participants and problem sets. $MSE$ = mean squared error; $M_1$ = differentiated by ability model; $M_2$ = fully differentiated model; $M_3$ = undifferentiated model. The differentiated by ability model ($M_1$) shown in bold outperforms the other models.

## Discussion of Model-Based Results

We contrasted three models and assessed the ability of the models to capture the qualitative patterns as well as match the human predictions in a quantitative way. The best-performing model was the differentiated by ability ($M_1$) model. It is a model with relatively few parameters that makes an assumption that there is a simple link between the mental model of self and other. Model $M_1$ learns only one differential ability parameter linking self- to other-assessment. Note that this is one of many ways to formulate how self- and other-assessment are tied together. Our claim is that for simpler tasks and with small amounts of data, this link between self- and other-assessment remains low dimensional. How quickly these models grow in complexity needs to be explored in future work.

Predictions from the differentiated by ability model ($M_1$) replicate the qualitative pattern we see in our empirical results while also being quantitatively closest to the observed data as shown in Table 3. The other two models ($M_2$ and $M_3$) fail to simultaneously capture the relationship between estimated self and other scores (Figure 4) and the divergence of estimated scores for top and bottom performers (Figure 5). In contrast, in the absence of feedback, people only have their own encounters with the task to rely on. This reliance is best captured by models $M_1$ and $M_3$. In $M_3$, the estimated ability and problem difficulty are assumed to be the same for the other person, leading a person to predict similar performance in self- and other-assessments.

## Explaining Previous Empirical Findings on Knowledge Assessment

Up to this point, we have shown how the hierarchical knowledge assessment model can explain a variety of findings from an empirical paradigm that we specifically designed to test how people differentiate between their own and others' performance. However, the hierarchical model can also be applied to other empirical paradigms. In this section, we demonstrate the model's ability to explain how people's assessment of other's performance changes as different knowledge signals are made available to them (Tullis, 2018) and how people place themselves relative to others (Moore & Healy, 2008). For each of the experiments, we qualitatively compare model predictions from the hierarchical model to the observed data. The details of the simulations are presented in Appendices E and F.

## Metacognitive Cue Utilization for Knowledge Assessment

The availability of certain performance-related signals influences people's assessment of their performance on a task (Jost et al., 1998; Nelson et al., 1998; Tullis, 2018). In addition to assessing one's own knowledge, Nickerson proposes that the same signals may also guide one's assessment of others. For example, when asked to assess another person's performance on a task without doing the task themselves, a person may rely on a vague feeling of knowing about the task. However, if the person does the task themselves before assessing another person, they have access to additional information about their performance through signals such as the time it takes for them to perform the task. This information may enable the person to make a more informed assessment of another person's performance on the same task. Tullis (2018) proposed a theory of knowledge estimation as cue utilization that builds upon these previous

accounts on self- and other-knowledge assessment (Koriat, 1997; Nickerson, 1999; Thomas & Jacoby, 2013). In this theory, the degree of overlap between self-assessment and other-assessment depends on the cues available to oneself. These cues may depend on an individual's interactions with the task, information about the specific other person being assessed, or general information about the population.

Through a series of experiments, Tullis (2018) demonstrated that the bases and accuracy of the assessment of others depends on the conditions under which the assessment is elicited. In Experiment 1 by Tullis (2018), participants judged the percentage of other participants who would know the answer to a series of trivia questions. There were two experimental conditions. In the *answer before* condition, on each trial, participants first answered the trivia question and then subsequently estimated the proportion of other participants who would know the answer. In the *answer after* condition, participants first estimated for each trivia question the proportion of other participants who would know the answer and then answered the trivia questions. Experiment 2 included four experimental conditions. As in Experiment 1, participants were either required to answer trivia questions before estimating other participants' performance or they were asked to estimate the other participants' performance without needing to answer the question themselves. In addition, feedback was manipulated: participants either did or did not receive corrective feedback about the correct answer after answering each question. Table 4 describes the four conditions in Experiment 2 and the corresponding metacognitive signals available to the participants.

The left panels of Figures 9 and 10 summarize the key empirical findings. Results are reported as γ correlations between (a) predictions of other's knowledge and the time needed for the person to answer the question themselves and (b) predictions of other's knowledge and the accuracy of the participant themselves. Figure 9A shows that participants' predictions of others' knowledge were more strongly tied to their own performance when they were required to answer trivia questions themselves before estimating others' knowledge on the same questions. This is consistent with our hypothesis that people draw information through the process of answering questions when assessing others. The results also show that participants' assessment of others' improved when they were provided feedback about the accuracy of their answers (left panel of Figure 9B). This additional cue helped participants better assess the
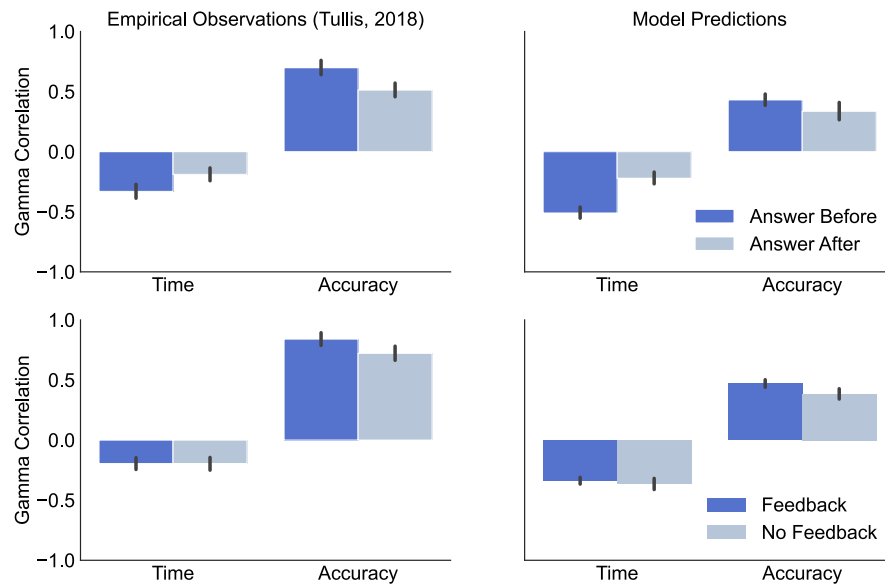
**Table 4**

*Assumptions About the Types of Knowledge Signals Available to People for the Different Conditions in Experiments 1 and 2 in Tullis (2018)*

| Condition | Types of knowledge signals |
|---|---|
| Exp. 1 | |
|   Answer after | FK |
|   Answer before | FK, RT |
| Exp. 2 | |
|   Answer not required, feedback not given | FK |
|   Answer not required, feedback given | FK, ACC |
|   Answer required, feedback not given | FK, RT |
|   Answer required, feedback given | FK, RT, ACC |

*Note.* FK = feeling of knowing; RT = response time; ACC = accuracy; Exp. = experiment.

**Figure 9**

*Observed and Model-Predicted Correlations Between a Person's Prediction of Others' Knowledge and the Time Needed for the Person to Answer the Question Themselves and Their Accuracy*
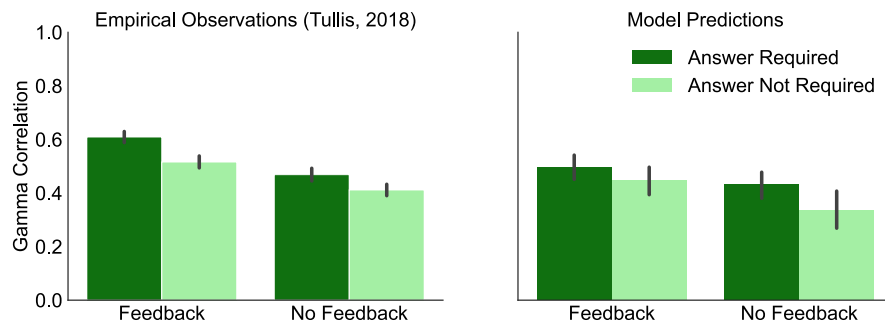


*Note.* The observed data are from Tullis (2018). The top row shows the results from the answer before and answer after conditions in Experiment 1. The bottom row shows results for the feedback and no-feedback conditions in Experiment 2. See the online article for the color version of this figure.

difficulty of each question and hence make better assessments of others' performance. Moreover, negative γ correlations between participant's predictions for others and the time they took to answer the questions suggest that participants expected others to perform worse on questions that took them longer to answer. This supports our assumption that participants use response time as a signal to assess the difficulty of problems and therefore to inform their assessment of others. However, there was no significant difference in this effect between the feedback and no-feedback conditions.

To apply the hierarchical knowledge assessment framework to the other-assessment task presented in Tullis (2018), we will assume that the experimental conditions determine which metacognitive cues or knowledge signals are available to a person when assessing themselves and others. We will use $x_{i,j}^{FK}$, $x_{i,j}^{RT}$, and $x_{i,j}^{ACC}$ to denote the three types of knowledge signals potentially available to participant $i$ for problem j: *feeling of knowing* (FK), *response time* (RT), and *performance feedback* (accuracy [ACC]), respectively. We assume that these knowledge signals are produced according to:

**Figure 10**

*Observed and Model-Predicted Correlations Between a Person's Prediction About Others' Knowledge and the (Sign Reversed) Difficulty of the Questions*



*Note.* The observed data are from Experiment 2 from Tullis (2018) across the feedback and no-feedback conditions. Note that the difficulty of a question for the empirical observations was based on the empirical proportion of participants that answered the question correctly. For the model predictions, difficulty of the questions is the inferred latent difficulty. See the online article for the color version of this figure.

$$x_{i,j}^{FK} \sim f(p_{i,j}^{s}, \eta), \quad x_{i,j}^{RT} \sim g(p_{i,j}^{s}, v), \quad x_{i,j}^{ACC} \sim h(p_{i,j}^{s}), \qquad (7)$$

where functions $f$, $g$, and $h$ link the knowledge signals to a person $i$'s estimate about their probability of being correct on problem $j$ ($p_{i,j}^{s}$) and $\eta$, $v$ encode the noise in the mapping to the observed knowledge. The mappings encode simple monotonic relationships between the probability correct and the knowledge signals. For example, feeling of knowing ($x_{i,j}^{FK}$) is modeled as linearly related to $p_{i,j}^{s}$—the more likely a person is correct, the stronger their feeling of knowing. In contrast, we expect people's response times $x_{i,j}^{RT}$ to be inversely related to $p_{i,j}^{s}$—the longer it takes people to solve a problem the harder they think it is. Note that in this experimental setup, participants only have access to their estimates of their response time. They do not observe the response time of other participants.

In Experiment 1 in Tullis (2018), in the answer after condition, participants judge other participants' performance before answering the question themselves, and hence participants only have a FK signal available to make knowledge assessments, that is, $x_{i,j}^{s} = \{x_{i,j}^{FK}\}$. In contrast, in the answer before condition, participants are required to answer the questions before evaluating others. Therefore, they have access to their response time in addition to the FK signal, that is, $x_{i,j}^{s} = \{x_{i,j}^{FK}, x_{i,j}^{RT}\}$. Table 4 details the assumptions about the types of knowledge signals available to people across different conditions and experiments.

In the experimental task, participants have to estimate the percentage of other participants who know the answer to a series of trivia questions. This can be thought of as assessing the performance of an average person instead of a specific individual. Since participants do not have access to any knowledge signals ($x^{o}$) pertaining to the other person, they can only make estimates about an average other person. In the absence of $x^{o}$, our modeling setup assumes that $a^{o}$ is a random draw from the population and hence represents the ability of an average person. Therefore, we frame the inference problem for the participant to estimate $a^{o}$ and problem difficulty $d$ on the basis of the observed knowledge signals $x^{s}$. Since we do not have access to the raw experimental data from the article, we first simulate experimental data for Experiments 1 and 2 using simple assumptions about individual differences in ability, variability of question difficulty, as well as basic assumptions about the functional forms used in Equation 7. Next, we apply the differentiated by ability model to simulate the inference process on the basis of the simulated experimental data (see Appendix E for details). The qualitative results shown here do not depend critically on the choice of simulation parameters.

Our model's predictions closely track the qualitative trends observed in the experimental data for Experiments 1 and 2, as demonstrated in Figure 9. In Figure 9A, the model predictions are consistent with the empirical observation that participants in the answer before condition showed a significantly stronger negative correlation between the time they took to answer a question and their accuracy of other assessment than participants in the answer after condition (i.e., participants estimated lower scores for others on questions that took them longer to answer). Additionally, the model predicts a positive correlation between participants' accuracy and their predictions of others' knowledge (i.e., participants tend to estimate higher scores for others on questions they themselves answered correctly). Similarly, for Experiment 2 (Figure 9B), the model predicts that participants estimate lower scores for others

on questions that took them longer to answer. This effect is stronger in the feedback condition than in the no-feedback condition. Additionally, the model captures the finding that participants tend to estimate higher scores for others on questions they themselves answered correctly. Figure 10 shows that the model predicts, consistent with the empirical observations, that participants' estimates of others improved when they were required to answer the question themselves and then were provided feedback. Overall, these results show that our model is able to accurately capture knowledge assessment across different experimental conditions.

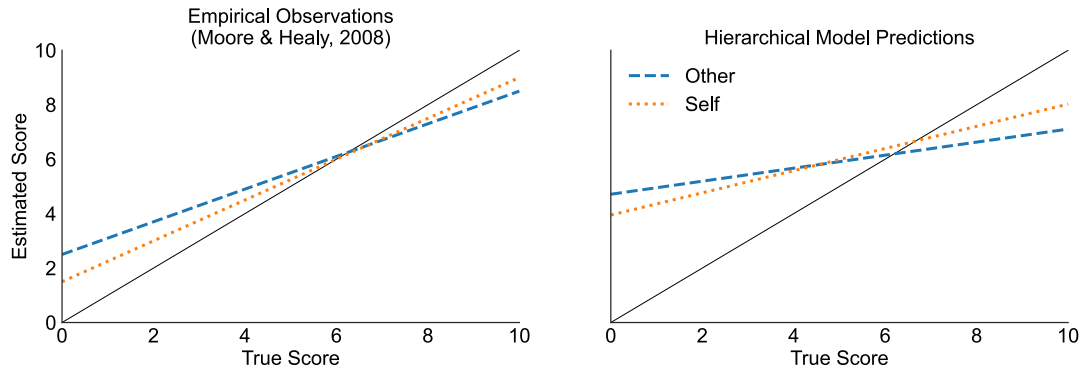## Overestimation and Overplacement

People's assessment of their own performance and the performance of others is known to be biased in several ways (Dunning, 2011; Larrick et al., 2007; Moore, 2007; Moore & Healy, 2008; Tullis, 2018). In particular, people tend to believe that they are less likely than average to exhibit extraordinary abilities and more likely than average to exhibit ordinary abilities (Moore, 2007). These beliefs about ability also depend on task difficulty.

Moore and Healy (2008) showed that on difficult tasks, people tend to overestimate their performance but incorrectly believe that they are worse than others. Whereas, on easy tasks, people tend to underestimate their performance but incorrectly believe they are better than others (Dunning, 2011; Moore & Healy, 2008). These findings can be attributed to two forms of overconfidence that people often display: *overestimation* and *overplacement*. For example, in the experimental paradigm from Moore and Healy (2008), participants answered trivia questions and predicted their own score and the score of a randomly selected participant at three different stages of the experiment. First, participants made predictions about themselves and the other participant before they had any specific information about the quiz they were about to take. Second, they answered quiz questions and then estimated their own scores and the other participant's score again. This is termed their *interim* estimate. Finally, participants were shown the correct answers to the quiz and asked to make final estimates about their performance and the other participant's performance.

The hierarchical knowledge assessment model is consistent with the theory presented by Moore and Healy (2008). The authors present a theory of overconfidence, which assumes that people have imperfect information about their own performances and even worse information about the performances of others. As a result, people's estimates of themselves are regressive, but their estimates of others are even more regressive. The left panel of Figure 11 exemplifies the theory's prediction of participants' regressive estimates about the performance of self and others. The right panel of Figure 11 demonstrates that our model predictions are consistent with the predictions of their theory of overconfidence and the empirical data presented in Moore and Healy (2008)—people's estimates of others' performance are more regressive than their estimates of their own performance. This qualitative trend is observed for a broad range of parameter values in our simulations. The main difference between the two theories is that the hierarchical model was designed to apply to a broader variety of empirical manipulations and tasks. The hierarchical framework provides explicit ways to model manipulations of problem difficulty, feedback, ordering of answering relative to other assessment, as well as situations that lead to knowledge signals specific to other people.

**Figure 11**

*Relationship Between the Estimated Performance of Self and Other and True Performance of Self and Other as Predicted by the Theory of Overconfidence (Moore & Healy, 2008) and as Predicted by the Hierarchical Model*



*Note.* See the online article for the color version of this figure.

The empirical observation columns in Table 5 show the degree of participants' overplacement and overestimation in the interim phase of the experiment. Higher positive values correspond to higher levels of overestimation and overplacement, and negative values correspond to underestimation and underplacement. The degree of overestimation was evaluated by the difference between the estimate of their performance and the person's true performance (i.e., $\hat{x}^{s,ACC} - x^{s,ACC}$). The degree of overplacement was evaluated by a difference of two differences: first, the difference between the estimated performance of self and other and second, the difference between the actual performance of self and other, that is, $(\hat{x}^{s,ACC} - \hat{x}^{o,ACC}) - (x^{s,ACC} - x^{o,ACC})$. This can be understood as the difference between a person's estimate of how much better they are when compared to another person $(\hat{x}^{s,ACC} - \hat{x}^{o,ACC})$ and the true difference between the two people $(x^{s,ACC} - x^{o,ACC})$. The empirical results show that participants tend to overestimate their performance on hard problems and underestimate their performance on easier problems. Furthermore, participants overplace their performance on easy problems and underplace their performance on difficult problems.

We simulated the hierarchical knowledge assessment model for the interim stage of the experiment using the same setup and simulation parameters as used for the simulations of the Tullis (2018) experiments (see Appendix F for details). At the interim stage of the experiment, we assume that participants have access to feeling of knowing and response time signals, similar to the answer before

condition in Experiment 1 of Tullis (2018), that is, $x^s_{i,j} = \{x^{FK}_{i,j}, x^{RT}_{i,j}\}$. We use the model to simulate the knowledge signals available to participants in the experiment. We also simulate a distribution of problem difficulty and refer to the highest 33% difficulty values as hard, the lowest 33% as easy, and the rest as medium. Next, we simulate the task faced by the participant: the problem of inferring $x^{s,ACC}$ and $x^{o,ACC}$ (i.e., producing estimates $\hat{x}^{s,ACC}, \hat{x}^{o,ACC}$) given the available knowledge signals $x^s$. Finally, to analyze the model predictions, we assess the degree of overestimation and overplacement using the same evaluation approach used to analyze the empirical data. The model prediction in Table 5 demonstrates our model's ability to capture the relationship between task difficulty and people's tendency to overplace or overestimate their performance. In line with the empirical observations, our model predicts that people underplace but overestimate their performance on difficult problems, and people overplace and underestimate their performance on easy problems.

At first glance, it may seem that the quantitative predictions of our model in Table 5 significantly diverge from the empirical observations. However, it is important to recognize that the empirical effects of overplacement and underplacement reported in Moore and Healy (2008) are relatively small. Our primary objective was not to achieve exact quantitative matches but rather to demonstrate that the hierarchical model makes qualitatively accurate predictions for self- and other-assessment phenomena reported in the literature. Additionally, it is worth noting that we use the same

**Table 5**

*Overestimation and Overplacement*

| Difficulty | Overestimation | | Overplacement | |
|---|---|---|---|---|
| | Empirical observations | Model predictions | Empirical observations | Model predictions |
| Easy | −.22 (.93) | −1.09 (2.27) | .48 (2.59) | .06 (2.53) |
| Medium | .01 (1.27) | 2.2 (3.04) | .04 (3.91) | −.01 (3.36) |
| Hard | .79 (1.50) | 4.33 (2.66) | −1.36 (2.39) | −.41 (2.78) |

*Note.* Empirical observations from Moore and Healy (2008) and model predictions for overestimation and overplacement when making self- and other-knowledge assessment at the interim phase for three different question difficulties (standard deviations in parentheses).

set of assumptions and parameter values to simulate data for all the experiments from Tullis (2018) and Moore and Healy (2008).

## Discussion

Knowing what other agents know is central to communication and cooperation between agents. Much of the current computational work on the theory of mind has focused on inferring beliefs and goals of other people by observing intentional behavior in spatial environments (Baker et al., 2009, 2017). However, developing an accurate model of another agent not only requires an understanding of their goals and beliefs, which can explain their movements in a physical environment, but also their knowledge states, which can explain their performance on knowledge tasks. In our theoretical framework, we focus on understanding how people assess the knowledge states of other people in the absence of any physical or verbal cues—they only receive quantitative feedback about their assessment of the other person's performance. The key idea of our work is that people combine their own experience on a task with information received about the other person's performance to make assessments of the other's knowledge states.

Previous research to understand how humans infer knowledge states of other humans was limited to empirical studies (Jameson et al., 1993; Nelson, 1984) and descriptive theories (Nickerson, 1999). However, there is increasing interest in developing models of reasoning about other people's knowledge states (Aboody et al., 2021; Berke & Jara-Ettinger, 2021). Aboody et al. (2021) presented a computational account of how people infer knowledge of another person based on the expectation that the other person maximizes epistemic utility when making choices. In this research, we take a complementary view of knowledge assessment of others. Our framework formalizes how humans construct mental models of other humans' knowledge solely based on the observed quantitative performance of the other person. We developed and tested three computational models on the basis of a simple empirical paradigm where the participant is asked to make inferences about the other person. As the experiment progresses, limited information about the other person is made available to the participant. For example, after receiving feedback about their first prediction, there is only one data point about the other person that is available to the participant. Still, despite the small amount of information, participants are able to update their mental model of the other person and improve their predictions over subsequent prediction rounds. We suggest that there are two main components that drive people's estimation of the other person's performance. The first is people's tendency to generalize their experience with the task to the other person's behavior. This explains the close association between people's self and other estimates—people use their estimates of task difficulty to adjust their beliefs about the other person's performance. The second component is their capacity to distinguish between their own ability and the other person's ability. This is made apparent by people's quickly diverging estimates of top and bottom other performers in our experiment.

### Sparse Data Encourages Linking Mental Models of Self and Other

From a computational perspective, people are often faced with situations where not many observations are available about another individual, making it difficult to learn detailed and complex mental models of that individual. Instead, a simpler mental model with few parameters to estimate might be effective (at least in the initial interaction with the individual). In this research, we contrasted three computational models for the inference of knowledge states. The models varied in the degree to which the mental models of self and other are differentiated. In the simplest mental model of other ($M_3$; undifferentiated), no parameters need to be updated as the mental model for the other person is the same as the mental model for self. In the most complex mental model of other ($M_2$; fully differentiated), not only the ability of the other person needs to be estimated but also the experienced difficulty for each type of problem. This model allows for the possibility that what is easy for one's self could be challenging for the other and vice versa. We found evidence for a computational model with an intermediate level of complexity ($M_1$; differentiated by ability) that involves just a single parameter: the relative ability of the other individual. This simple mental model allows one to quickly extrapolate how likely it is that an individual can successfully perform a task with very few observations.

Our results support our claim that in the presence of feedback, people learn about the other person's ability relative to their own while also drawing information from their own experience from the task. The differentiated by ability model that best accounts for the observed data makes the assumption that the way people reason about the other person's performance is through the lens of their own self-assessment process. This assumption is consistent with a second-order model of metacognition, which suggests that humans self-reflect and think about others using similar mental processes (Fleming & Daw, 2017). We posit that the same machinery that enables people to estimate their performance also enables them to judge another person's performance. However, we do not address the issue of the number of systems involved in metacognition and mind reading. Our results simply point out that self-knowledge can be informative and is used by people to make predictions about other people's knowledge.

### Proposals for Future Investigations

We now discuss in greater detail how the self- and other-assessment can be extended to handle other interesting situations involving social cues, multiple agents, multidimensional ability, and AI agents assessed by humans and humans assessing AI agents.

### *Utilizing Social Cues to Assess Others*

During social interactions, people have the opportunity to perceive and interpret numerous signals such as facial expressions, natural language, and voice intonations of the person they interact with. These cues can serve as Supplemental Information when evaluating the other person's knowledge.

An experiment was conducted by Jameson et al. (1993) where "observer" participants witnessed "target" participants answering trivia questions and then made predictions about the targets' performance on those questions. In contrast, "judge" participants made predictions without observing the targets answering the questions. The findings revealed that observers displayed greater accuracy in predicting the performance of the individuals compared to judges. This divergence in prediction accuracy can be attributed to specific cues exhibited by the observed person, which provide additional insights beyond mere performance statistics and task

experience. These cues may include the time taken by the target to respond to a question, their facial expressions, the confidence conveyed through their voice, and possibly other factors. Similarly, Brennan and Williams (1995) assigned participants the task of listening to responses given by others to general knowledge questions and evaluated their perception of the "feeling of another's knowing." The results demonstrated that people's assessments of others' knowledge were influenced by changes in intonation, response delays, and the use of filler phrases—confirming that individuals pay attention to the metacognitive information conveyed by speakers regarding their states of knowledge.

Most recent computational approaches to capture knowledge assessment, including our ongoing research, have focused on situations in which individuals possess limited information regarding others. However, there is a need to explore how an expanded set of behavioral and social cues can be quantified and utilized as Supplemental Information for predicting people's accuracy of other-assessment.

### Assessing Multiple Other Agents

More often than not, people work with multiple other agents to accomplish tasks. An important extension of the current work is to see how easily peoples' mental models scale to groups of others or how well can people make inferences about the knowledge states of multiple other teammates when working in a group. For example, when playing a trivia quiz with a group of people, players continuously appraise other players' expertise on a variety of domains. This mechanism of group appraisal and coordination was formalized by Wegner (1987) as a transactive memory system (TMS). TMS is a property of a group that consists of knowledge stored in each person's memory and metamemory that encodes different teammates' domains of expertise. Mei et al. (2017) mathematically formalized TMS as an appraisal network and described asymptotic properties of the team. However, how people learn such an appraisal network in practice is not well investigated. Here, we focused on assessing only one other person and the model that best described the empirical data was a low-dimensional model. It is likely that humans learn a sparse representation of ability to differentiate between multiple teammates. Such parsimony would be essential to manage cognitive overload and resource constraints.

### Multidimensional Ability

In daily life, people often interact with domain experts. For example, we expect a birder to have a wider knowledge of birds than a layperson. However, information about the birder's knowledge of birds does not necessarily position us better to assess their knowledge in related domains, such as classifying dog breeds, or unrelated domains, such as identifying Renaissance painters. An important simplification in the self- and other-assessment models is that they encode ability as a one-dimensional parameter. We focused on a simple mental model where differentiation was based on a single-dimensional ability. However, we do not rule out the possibility that people are developing increasingly complex multidimensional mental models of others, as more information is observed.

We know that humans are capable of planning based on beliefs, goals, and resource constraints (Baker et al., 2009; Gopnik & Meltzoff, 1997; Lieder & Griffiths, 2020) and can use inverse

planning to infer beliefs and goals from observed behavior of other agents (Shum et al., 2019; Tauber & Steyvers, 2011). While traditional accounts of the theory of mind provide important qualitative insights into how humans make these complex inferences about other minds (Gopnik & Meltzoff, 1997), recent work provides computational frameworks to capture human judgments across a range of social interactions (Baker, 2012; Baker et al., 2017; Shum et al., 2019). However, quantitative variation in the human ability to reason about knowledge of other agents is not well studied. A straightforward extension of the self- and other-assessment models would be to account for differences in ability across different categories presented to the participant. MIRT is often used to analyze performance on tasks where multiple abilities are at play (Ackerman et al., 2003; Hartig & Höhler, 2009). MIRT is a generalization of unidimensional IRT models where the probability of success is modeled as a function of multiple ability dimensions. Such models can also be applied to instances where mixtures of abilities are required for individual test items.

### Humans Assessing AI

Humans are increasingly interfacing with artificial agents (AI) to make joint decisions in a variety of real-world applications (Kleinberg et al., 2018; Ott et al., 2011; Patel et al., 2019; Rajpurkar et al., 2020; Wright et al., 2017). A common pitfall of such collaborative human–AI decision making is the ineffective treatment of advice from an AI agent by the human. To correctly assess and use an AI agent's advice, the human must infer the AI agent's expertise and knowledge about the task at hand to build a good mental model of the AI's ability. Our work presents a first step to understanding a human's assessment of other human's ability from a computational perspective. Future work should investigate how humans update their assessment of ability when the other agent is an AI agent.

An important assumption of the current model is that humans can generalize their subjective assessment of the difficulty of the task to the relative difficulty experienced by another human. In essence, people assume that what is difficult for them is difficult for another human. However, this assumption might not hold true when humans interact with AI agents. Extensions of the current framework may be used to investigate how humans assess ability of an AI agent that has complementary abilities to the human (finds different tasks difficult or easy when compared to the human)—can people simultaneously learn a nuanced model of ability and build a high-dimensional representation of another agent's experience in the task?

### Conclusions

How a mind understands another mind is a fundamental question in psychology. While there is prior research on how people make theory of mind judgments about intentions and goals of other agents, there is relatively little investigation of how people assess knowledge of other agents. In this work, we develop a theoretical framework that describes the underlying computation that people employ when assessing the knowledge of other agents. Our empirical results and model predictions demonstrate that people's evaluation of the other person's performance (a theory of mind computation) is linked to their evaluation of their own performance

(a metacognitive computation). The models presented in the article provide a starting point for a more comprehensive exploration of how humans assess other agents.

# References

Aboody, R., Davis, I., Dunham, Y., & Jara-Ettinger, J. (2021). *I can tell you know a lot, although I'm not sure what: Modeling broad epistemic inference from minimal action* [Conference session]. Proceedings of the 43rd Annual Conference of the Cognitive Science Society.

Ackerman, T. A., Gierl, M. J., & Walker, C. M. (2003). Using multidimensional item response theory to evaluate educational and psychological tests. *Educational Measurement: Issues and Practice*, 22(3), 37–51. https://doi.org/10.1111/j.1745-3992.2003.tb00136.x

Baer, C., Malik, P., & Odic, D. (2021). Are children's judgments of another's accuracy linked to their metacognitive confidence judgments? *Metacognition and Learning*, 16, 485–516. https://doi.org/10.1007/s11409-021-09263-x

Baker, C. L. (2012). *Bayesian theory of mind: Modeling human reasoning about beliefs, desires, goals, and social relations* [PhD thesis]. Massachusetts Institute of Technology.

Baker, C. L., Jara-Ettinger, J., Saxe, R., & Tenenbaum, J. B. (2017). Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour*, 1(4), Article 0064. https://doi.org/10.1038/s41562-017-0064

Baker, C. L., Saxe, R., & Tenenbaum, J. B. (2009). Action understanding as inverse planning. *Cognition*, 113(3), 329–349. https://doi.org/10.1016/j.cognition.2009.07.005

Berke, M., & Jara-Ettinger, J. (2021). *Thinking about thinking through inverse reasoning* [Conference session]. Proceedings of the Annual Meeting of the Cognitive Science Society.

Brennan, S. E., & Williams, M. (1995). The feeling of another's knowing: Prosody and filled pauses as cues to listeners about the metacognitive states of speakers. *Journal of Memory and Language*, 34(3), 383–398. https://doi.org/10.1006/jmla.1995.1017

Dunning, D. (2011). The dunning–kruger effect: On being ignorant of one's own ignorance. In J. M. Olson & M. P. Zanna (Eds.), *Advances in experimental social psychology* (Vol. 44, pp. 247–296). Elsevier. https://doi.org/10.1016/B978-0-12-385522-0.00005-6

Dunning, D., & Helzer, E. G. (2014). Beyond the correlation coefficient in studies of self-assessment accuracy: Commentary on zell & krizan (2014). *Perspectives on Psychological Science*, 9(2), 126–130. https://doi.org/10.1177/1745691614521244

Fleming, S. M. (2021). *Know thyself: The science of self-awareness*. Basic Books.

Fleming, S. M., & Daw, N. D. (2017). Self-evaluation of decision-making: A general bayesian framework for metacognitive computation. *Psychological Review*, 124(1), 91–114. https://doi.org/10.1037/rev0000045

Fleming, S. M., & Lau, H. C. (2014). How to measure metacognition. *Frontiers in Human Neuroscience*, 8, Article 443. https://doi.org/10.3389/fnhum.2014.00443

Fox, J.-P. (2010). *Bayesian item response modeling: Theory and applications*. Springer Science & Business Media.

Frith, C. D., & Frith, U. (1999). Interacting minds–a biological basis. *Science*, 286(5445), 1692–1695. https://doi.org/10.1126/science.286.5445.1692

Fussell, S. R., & Krauss, R. M. (1991). Accuracy and bias in estimates of others' knowledge. *European Journal of Social Psychology*, 21(5), 445–454. https://doi.org/10.1002/ejsp.2420210507

Gneiting, T., & Walz, E.-M. (2021). Receiver operating characteristic (ROC) movies, universal roc (UROC) curves, and coefficient of predictive ability (CPA). *Machine Learning*, 111, 2769–2797. https://doi.org/10.1007/s10994-021-06114-3

Gopnik, A., & Astington, J. W. (1988). Children's understanding of representational change and its relation to the understanding of false belief and the appearance-reality distinction. *Child Development*, 59(1), 26–37. https://doi.org/10.1111/j.1467-8624.1988.tb03192.x

Gopnik, A., & Meltzoff, A. N. (1997). *Words, thoughts, and theories*. Mit Press.

Greene, W. H., & Hensher, D. A. (2010). *Modeling ordered choices: A primer*. Cambridge University Press.

Hart, J. T. (1965). Memory and the feeling-of-knowing experience. *Journal of Educational Psychology*, 56(4), 208–216. https://doi.org/10.1037/h0022263

Hartig, J., & Höhler, J. (2009). Multidimensional IRT models for the assessment of competencies. *Studies in Educational Evaluation*, 35(2–3), 57–63. https://doi.org/10.1016/j.stueduc.2009.10.002

Jameson, A., Nelson, T. O., Leonesio, R. J., & Narens, L. (1993). The feeling of another person's knowing. *Journal of Memory and Language*, 32(3), 320–335. https://doi.org/10.1006/jmla.1993.1017

Jansen, R. A., Rafferty, A. N., & Griffiths, T. (2020). *A rational model of sequential self-assessment* [Paper presentation]. 42nd Annual Meeting of the Cognitive Science Society: Developing a Mind: Learning in Humans, Animals, and Machines, CogSci 2020, Virtual, Online.

Jansen, R. A., Rafferty, A. N., & Griffiths, T. L. (2021). A rational model of the dunning–kruger effect supports insensitivity to evidence in low performers. *Nature Human Behaviour*, 5(6), 756–763. https://doi.org/10.1038/s41562-021-01057-0

Jenkins, A. C., Macrae, C. N., & Mitchell, J. P. (2008). Repetition suppression of ventromedial prefrontal activity during judgments of self and others. *Proceedings of the National Academy of Sciences*, 105(11), 4507–4512. https://doi.org/10.1073/pnas.0708785105

Jost, J. T., Kruglanski, A. W., & Nelson, T. O. (1998). Social metacognition: An expansionist review. *Personality and Social Psychology Review*, 2(2), 137–154. https://doi.org/10.1207/s15327957pspr0202_6

Kelley, C. M., & Jacoby, L. L. (1996). Adult egocentrism: Subjective experience versus analytic bases for judgment. *Journal of Memory and Language*, 35(2), 157–175. https://doi.org/10.1006/jmla.1996.0009

Kleinberg, J., Lakkaraju, H., Leskovec, J., Ludwig, J., & Mullainathan, S. (2018). Human decisions and machine predictions. *The Quarterly Journal of Economics*, 133(1), 237–293. https://doi.org/10.1093/qje/qjx032

Koriat, A. (1997). Monitoring one's own knowledge during study: A cue-utilization approach to judgments of learning. *Journal of Experimental Psychology: General*, 126(4), 349–370. https://doi.org/10.1037/0096-3445.126.4.349

Koriat, A. (2000). The feeling of knowing: Some metatheoretical implications for consciousness and control. *Consciousness and Cognition*, 9(2), 149–171. https://doi.org/10.1006/ccog.2000.0433

Koriat, A., & Ackerman, R. (2010). Metacognition and mindreading: Judgments of learning for self and other during self-paced study. *Consciousness and Cognition*, 19(1), 251–264. https://doi.org/10.1016/j.concog.2009.12.010

Kruger, J., & Dunning, D. (1999). Unskilled and unaware of it: How difficulties in recognizing one's own incompetence lead to inflated self-assessments. *Journal of Personality and Social Psychology*, 77(6), 1121–1134. https://doi.org/10.1037/0022-3514.77.6.1121

Larrick, R. P., Burson, K. A., & Soll, J. B. (2007). Social comparison and confidence: When thinking you're better than average predicts overconfidence (and when it does not). *Organizational Behavior and Human Decision Processes*, 102(1), 76–94. https://doi.org/10.1016/j.obhdp.2006.10.002

Leibert, T. W., & Nelson, D. L. (1998). The roles of cue and target familiarity in making feeling of knowing judgments. *The American Journal of Psychology*, 111(1), 63–75. https://doi.org/10.2307/1423537

Lieder, F., & Griffiths, T. L. (2020). Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources. *Behavioral and Brain Sciences*, 43, Article E1. https://doi.org/10.1017/S0140525X1900061X

Mei, W., Friedkin, N. E., Lewis, K., & Bullo, F. (2017). Dynamic models of appraisal networks explaining collective learning. *IEEE Transactions on Automatic Control*, 63(9), 2898–2912. https://doi.org/10.1109/TAC.2017.2775963

Miller, G. A. (1995). Wordnet: A lexical database for English. *Communications of the ACM*, *38*(11), 39–41. https://doi.org/10.1145/219717.219748

Mitchell, J. P., Banaji, M. R., & Macrae, C. N. (2005). The link between social cognition and self-referential thought in the medial prefrontal cortex. *Journal of Cognitive Neuroscience*, *17*(8), 1306–1315. https://doi.org/10.1162/0898929055002418

Moore, D. A. (2007). Not so above average after all: When people believe they are worse than average and its implications for theories of bias in social comparison. *Organizational Behavior and Human Decision Processes*, *102*(1), 42–58. https://doi.org/10.1016/j.obhdp.2006.09.005

Moore, D. A., & Cain, D. M. (2007). Overconfidence and underconfidence: When and why people underestimate (and overestimate) the competition. *Organizational Behavior and Human Decision Processes*, *103*(2), 197–213. https://doi.org/10.1016/j.obhdp.2006.09.002

Moore, D. A., & Healy, P. J. (2008). The trouble with overconfidence. *Psychological Review*, *115*(2), 502–517. https://doi.org/10.1037/0033-295X.115.2.502

Nelson, T. O. (1984). A comparison of current measures of the accuracy of feeling-of-knowing predictions. *Psychological Bulletin*, *95*(1), 109–133. https://doi.org/10.1037/0033-2909.95.1.109

Nelson, T. O., Kruglanski, A. W., & Jost, J. T. (1998). Knowing thyself and others: Progress in metacognitive social psychology. In V. Y. Yzerbyt, G. Lories, & B. Dardenne (Eds.), *Metacognition: Cognitive and social dimensions* (pp. 69–89). Sage Publications. https://doi.org/10.4135/9781446279212.n5

Nelson, T. O., & Narens, L. (1980). Norms of 300 general-information questions: Accuracy of recall, latency of recall, and feeling-of-knowing ratings. *Journal of Verbal Learning and Verbal Behavior*, *19*(3), 338–368. https://doi.org/10.1016/S0022-5371(80)90266-2

Nicholson, T., Williams, D. M., Lind, S. E., Grainger, C., & Carruthers, P. (2021). Linking metacognition and mindreading: Evidence from autism and dual-task investigations. *Journal of Experimental Psychology: General*, *150*(2), 206–220. https://doi.org/10.1037/xge0000878

Nickerson, R. S. (1999). How we know–and sometimes misjudge–what others know: Imputing one's own knowledge to others. *Psychological Bulletin*, *125*(6), 737–759. https://doi.org/10.1037/0033-2909.125.6.737

Nickerson, R. S., Baddeley, A., & Freeman, B. (1987). Are people's estimates of what other people know influenced by what they themselves know? *Acta Psychologica*, *64*(3), 245–259. https://doi.org/10.1016/0001-6918(87)90010-2

Ott, M., Choi, Y., Cardie, C., & Hancock, J. T. (2011). *Finding deceptive opinion spam by any stretch of the imagination*. PsyArXiv. https://doi.org/10.48550/arXiv.1107.4557

Patel, B. N., Rosenberg, L., Willcox, G., Baltaxe, D., Lyons, M., Irvin, J., Rajpurkar, P., Amrhein, T., Gupta, R., Halabi, S., Langlotz, C., Lo, E., Mammarappallil, J., Mariano, A. J., Riley, G., Seekins, J., Shen, L., Zucker, E., & Lungren, M. P. (2019). Human–machine partnership with artificial intelligence for chest radiograph diagnosis. *NPJ Digital Medicine*, *2*(1), Article 111. https://doi.org/10.1038/s41746-019-0189-7

Paulus, M., Tsalas, N., Proust, J., & Sodian, B. (2014). Metacognitive monitoring of oneself and others: Developmental changes during childhood and adolescence. *Journal of Experimental Child Psychology*, *122*, 153–165. https://doi.org/10.1016/j.jecp.2013.12.011

Rajpurkar, P., O'Connell, C., Schechter, A., Asnani, N., Li, J., Kiani, A., Ball, R. L., Mendelson, M., Maartens, G., van Hoving, D. J., Griesel, R.,

Ng, A. Y., Boyles, T. H., & Lungren, M. P. (2020). CheXaid: Deep learning assistance for physician diagnosis of tuberculosis using chest x-rays in patients with HIV. *NPJ Digital Medicine*, *3*(1), Article 115. https://doi.org/10.1038/s41746-020-00322-2

Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Institute of Educational Research.

Reckase, M. D. (2009). *Multidimensional item response theory* (pp. 79–112). Springer.

Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., & Fei-Fei, L. (2015). ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, *115*(3), 211–252. https://doi.org/10.1007/s11263-015-0816-y

Shum, M., Kleiman-Weiner, M., Littman, M. L., & Tenenbaum, J. B. (2019). *Theory of minds: Understanding behavior in groups through inverse planning* [Conference session]. Proceedings of the AAAI Conference on Artificial Intelligence.

Stan Development Team. (2020). *RStan: The R interface to Stan* (R package Version 2.21.2). http://mc-stan.org

Tauber, S., & Steyvers, M. (2011). *Using inverse planning and theory of mind for social goal inference* [Conference session]. Proceedings of the Annual Meeting of the Cognitive Science Society.

Thomas, R. C., & Jacoby, L. L. (2013). Diminishing adult egocentrism when estimating what others know. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *39*(2), 473–486. https://doi.org/10.1037/a0028883

Tullis, J. G. (2018). Predicting others' knowledge: Knowledge estimation as cue utilization. *Memory & Cognition*, *46*(8), 1360–1375. https://doi.org/10.3758/s13421-018-0842-4

Turner, B. M., Steyvers, M., Merkle, E. C., Budescu, D. V., & Wallsten, T. S. (2014). Forecast aggregation via recalibration. *Machine Learning*, *95*(3), 261–289. https://doi.org/10.1007/s10994-013-5401-4

Vaccaro, A. G., & Fleming, S. M. (2018). Thinking about thinking: A coordinate-based meta-analysis of neuroimaging studies of metacognitive judgements. *Brain and Neuroscience Advances*, *2*. https://doi.org/10.1177/2398212818810591

van der Linden, W. J., & Hambleton, R. K. (2013). *Handbook of modern item response theory*. Springer Science & Business Media.

Wegner, D. M. (1987). Transactive memory: A contemporary analysis of the group mind. In B. Mullen & G. R. Goethals (Eds.), *Theories of group behavior* (pp. 185–208). Springer.

Wright, D. E., Lintott, C. J., Smartt, S. J., Smith, K. W., Fortson, L., Trouille, L., Allen, C. R., Beck, M., Bouslog, M. C., Boyer, A., Chambers, K. C., Flewelling, H., Granger, W., Magnier, E. A., McMaster, A., Miller, G. R. M., O'Donnell, J. E., Simmons, B., Spiers, H., … Young, D. R. (2017). A transient search using combined human and machine classifications. *Monthly Notices of the Royal Astronomical Society*, *472*(2), 1315–1323. https://doi.org/10.1093/mnras/stx1812

Zell, E., & Krizan, Z. (2014). Do people have insight into their abilities? A metasynthesis. *Perspectives on Psychological Science*, *9*(2), 111–125. https://doi.org/10.1177/1745691613518075

Zhang, H., & Maloney, L. T. (2012). Ubiquitous log odds: A common representation of probability and frequency distortion in perception, action, and cognition. *Frontiers in Neuroscience*, *6*, Article 1. https://doi.org/10.3389/fnins.2012.00001

(*Appendices follow*)

# Appendix A

## The Ordered Probit Model

The ordered probit model, $x \sim$ OrderedProbit $(p, v, \sigma)$, is a generative model that maps a (latent) value $p$ to one of $M + 1$ discrete scores $x \in \{0, \dots, M\}$. In this process, noise is added to the latent value resulting in a new latent value, $p' = p + \varepsilon$, where $\varepsilon \sim N(0, \sigma)$ and the resulting discrete score is determined by the interval where $p'$ lies:

$$x = \begin{cases} 0 & \text{if } p' \leq v_1 \\ 1 & \text{if } v_1 < p' \leq v_2 \\ 2 & \text{if } v_2 < p' \leq v_3 \\ M & \text{if } p' > v_M \end{cases}. \quad (A1)$$

The ordered vector $v = [v_1, \dots, v_M]$ represents the transition points between different discrete scores. With this construction, the probability of producing a score $x = k$ conditional on the latent value $p$ is:
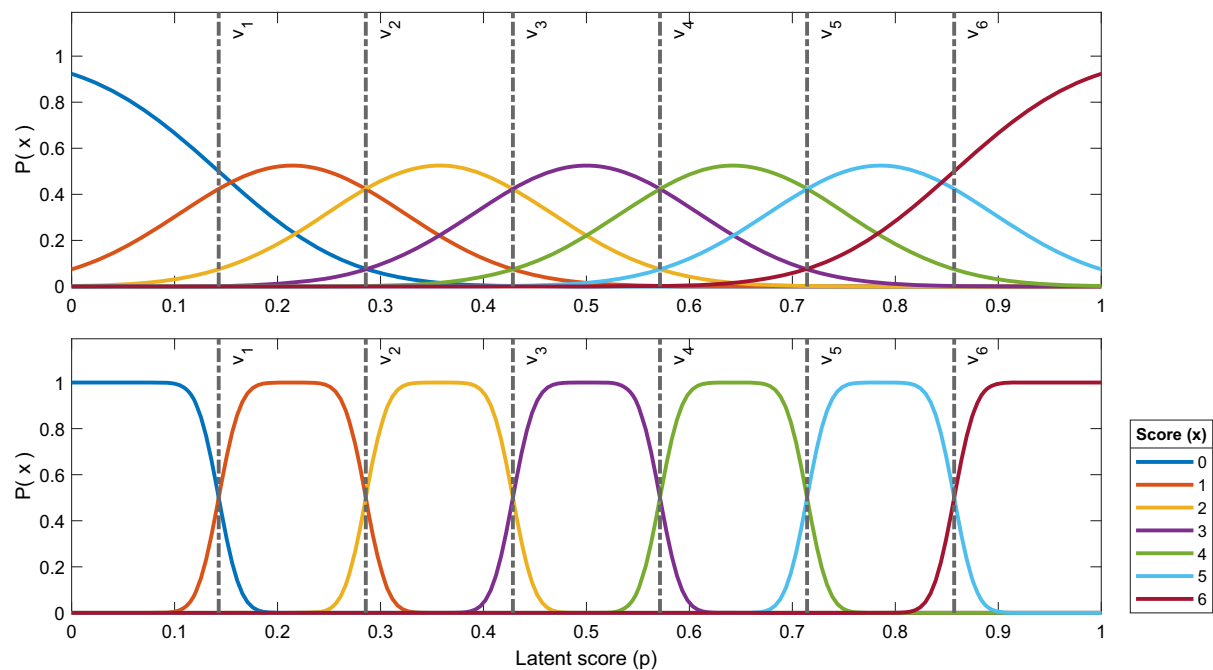
$$P(x = k|p, \sigma) = \Phi((v_{k+1} - p)/\sigma) - \Phi((v_k - p)/\sigma), \quad (A2)$$

where $\Phi$ is the cumulative standard normal distribution and $v_0 = -\infty$.

To simplify the model, we divide the 0–1 range into $M + 1$ equal intervals (i.e., $v = [1/(M + 1), 2/(M + 1), \dots, M/(M + 1)]$). With this construction, when $M = 12$ (as in our experiment), a latent value $p' = 1/12$ will result in a score $x = 1$, $p' = 2/12$ will result in a score $x = 2$, and so forth. Figure A1 shows an example of how the latent scores are mapped to scores when $M = 6$. Note that the higher value of the parameter $\sigma$ (top panel) results in a noisier mapping of latent probabilities to discrete scores.

**Figure A1**
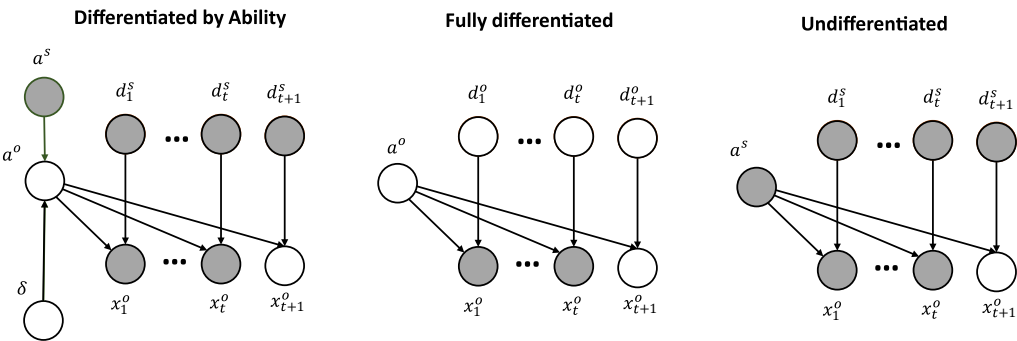*Illustration of the Ordered Probit Model When M = 6*



*Note.* Top and bottom panels are produced with $\sigma = 1/10$ and $\sigma = 1/60$, respectively. See the online article for the color version of this figure.

(*Appendices continue*)

## Appendix B

## Graphical Models

shows the graphical models for the prediction problem corresponding to the three assumptions about the relationship between self- and other-assessment. These graphical models illustrate the relationships between the observed and unobserved variables. Note that what is observable or unobserved is all from the perspective of the person reasoning about the other person.

**Figure B1**

*Graphical Models Corresponding to Three Different Other-Assessment Models for Predicting the Performance of Another Person*



*Note.* Shaded nodes show information that is known from the perspective of the person reasoning about the other person. Unshaded nodes show latent variables that need to be inferred. The key variable to infer is $x_{t+1}^o$, the performance of the target person on problem $t + 1$. See the online article for the color version of this figure.

## Appendix C

## Classification Problems

shows a list of the 16 types of classification problems used in the experiments along with the four response options for each classification problem.

**Table C1**

*List of the Classification Problems by Basic Category*

| No. | Category | Response options |
| --- | --- | --- |
| 1 | Bird | Crane (bird), Common redshank, Limpkin, Dunlin |
| 2 | Bird | Little blue heron, Oystercatcher, Dowitcher, Great egret |
| 3 | Bird | Bustard, Spoonbill, Hornbill, Bittern |
| 4 | Bird | Hummingbird, Bald eagle, Vulture, Kite |
| 5 | Dog | Shetland Sheepdog, Old English Sheepdog, Rottweiler, Komondor |
| 6 | Dog | Lhasa Apso, Airedale Terrier, West Highland White Terrier, Kerry Blue Terrier |
| 7 | Dog | Norwich Terrier, Irish Terrier, Scottish Terrier, Norfolk Terrier |
| 8 | Dog | Afghan Hound, Ibizan Hound, Norwegian Elkhound, Redbone Coonhound |
| 9 | Primate | Macaque, Titi, White-headed capuchin, Guenon |
| 10 | Primate | Langur, Black-and-white colobus, Marmoset, Common squirrel monkey |
| 11 | Primate | Gorilla, Chimpanzee, Gibbon, Baboon |
| 12 | Primate | Ring-tailed lemur, Geoffroy's spider monkey, Howler monkey, Siamang |
| 13 | Reptile | Green iguana, Desert grassland whiptail lizard, European green lizard, Carolina anole |
| 14 | Reptile | Ring-necked snake, Eastern hog-nosed snake, Vine snake, Worm snake |
| 15 | Reptile | Smooth green snake, Night snake, Kingsnake, Saharan horned viper |
| 16 | Reptile | Indian cobra, Sea snake, Water snake, Garter snake |

(*Appendices continue*)

# Appendix D

## Coefficient of Predictive Ability

Coefficient of Predictive Ability (CPA) is a rank-based measure that generalizes the area under the curve (AUC) to ordinal and continuous variables. For binary outcomes, CPA equals AUC, and for continuous outcomes, CPA relates linearly to Spearman's coefficient. We direct the readers to Gneiting and Walz (2021) for a detailed discussion on CPA.

Consider data of the form:

$$(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R} \times \mathbb{R}, \tag{D1}$$

where $x_i$ and $y_i$ are real numbers, for $i = 1, \ldots, n$. Let $z_1 < \ldots < z_m$ denote the $m \leq n$ unique values of $y_1, \ldots, y_n$ and define $n_c = \sum_{i=1}^{n} 1\{y_i = z_c\}$ such that $n_1 + \ldots + n_m = n$. We can reorder and write (Equation D1) as

$$(x_{11}, z_1), \ldots, (x_{1n_1}, z_1), \ldots, (x_{m1}, z_m), \ldots, (x_{mn_m}, z_m) \in \mathbb{R} \times \mathbb{R}, \tag{D2}$$

where $x_{i1}, x_{i2}, \ldots, x_{in_i}$ represent the $n_i$ different values of $x$ corresponding to $y = z_i$. This allows us to compute the CPA as the following:

$$\text{CPA} = \frac{\sum_{i=1}^{m-1} \sum_{j=i+1}^{m} \sum_{k=1}^{n_i} \sum_{l=1}^{n_j} (j-i) s(x_{ik}, x_{jl})}{\sum_{i=1}^{m-1} \sum_{j=i+1}^{m} (j-i) n_i n_j}, \tag{D3}$$

where $s$ is:

$$s(x, x') = 1\{x < x'\} + \frac{1}{2} 1\{x = x'\}. \tag{D4}$$

# Appendix E

## Simulation Details for Tullis (2018)

Tullis (2018) explored how people use a variety of metacognitive cues to infer the proportion of other people who know the answer to general knowledge questions. This section provides details on the simulation studies we conducted to apply our proposed hierarchical model to the data from Experiments 1 and 2. Since we do not have access to the raw experimental data from the article, we simulate experimental data for Experiments 1 and 2 and then apply our model to simulate the inference process of others' performance.

To simulate data at the participant level, we randomly generated ability levels, $a_i \sim N(0, 1)$, for 128 simulated participants who are performing the assessment, as well as 128 other participants to serve as a set of other participants. At the question level, we randomly generated the difficulty levels for 40 questions, $d_j \sim N(\mu_d, \sigma_d)$, where $\mu_d = 1$ and $\sigma_d$ are simulation parameters that determine overall mean performance and variability in question difficulty. For the self-assessed abilities, we use the same process as in Equation 2, to model the self-assessed abilities, $a_i^s \sim N(a_i, \sigma_a)$, where parameter $\sigma_a$ determines the noise in self-assessment. We use the IRT model in Equation 1 to calculate $p_{i,j}$, the true probability of correctly answering a question for every person $i$ on every question $j$.

The true probability of being correct ($p$) is used to generate different knowledge signals, including feeling of knowing ($x^{FK}$), response time ($x^{RT}$), and accuracy ($x^{ACC}$). We assume feeling of knowing is a random draw from a normal centered around $p_{i,j}$ and with an individual-specific variance $\delta_i$:

$$x_{i,j}^{FK} = p_{i,j} + N(\mu_{FK}, \delta_i), \delta_i \sim \text{Uniform}(0, \eta). \tag{E1}$$

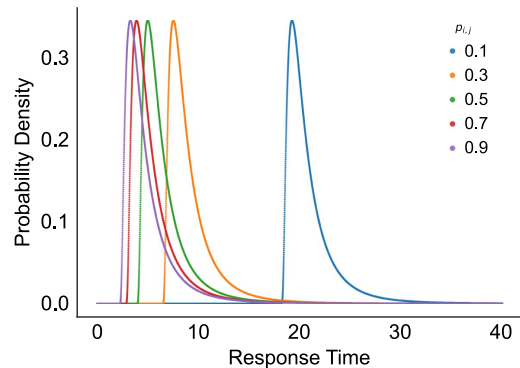Lower values of $\delta_i$ correspond to less noise in a participant's feeling of knowing and simulation parameter $\eta$ determines the degree of noise. To simulate response times, we assume an inverse relationship between RT and $p_{i,j}$:

$$x_{i,j}^{RT} \sim \text{LogNormal}\left(\frac{K}{p_{i,j} + \varepsilon_{i,j} + .01}, v\right)$$

$$\varepsilon_{i,j} \sim N(\omega, \zeta_i), \zeta_i \sim \text{Uniform}(a, b), \tag{E2}$$

where $\varepsilon_{i,j}$ is individual-specific noise in response time signals and .01 is added to avoid numerical instabilities. Simulation parameter $v$ determines the noise in the relationship between RT and accuracy. Figure E1 shows the RT distribution for different values of $p_{i,j}$. Our assumption results in people having higher RT for problems they have a lower probability of answering correctly and lower RT for problems they have a higher probability of answering correctly.

**Figure E1**

*Simulated Response Time Distributions for Different Values of $p_{i,j}$ and $K = 2$, $v = 2$*



*Note.* See the online article for the color version of this figure.

*(Appendices continue)*

We model participants' correctness on each problem $j$ as a Bernoulli draw with probability $p_{i,j}$

$$x_{i,j}^{ACC} \sim \text{Bern}(p_{i,j}). \qquad (E3)$$

To simulate the different experimental conditions of Experiments 1 and 2, we follow the logic of Table 4 that determines which knowledge signals are available in each condition. Next, we apply the hierarchical model of knowledge assessment on the simulated data. Based on the observed knowledge signals $x$ and the long-term self-estimate of ability $a^s$, the goal for the participant is to infer $x^{o,ACC}$ (which in this setup represents the performance of a randomly sampled person from the population). We used MCMC sampling to infer model parameters for the cognitive model presented in Figure 1A with different metacognitive signals $x$ and obtain samples from the posterior distribution of $a^o$. We used the Stan computing environment for posterior inference (Stan Development Team, 2020).

For simulating the experimental data, we use model parameters $K = 2$, $\mu_d = 1$, $\sigma_d = 2$, $\sigma_a = 0.5$, $\eta = .5$, $v = 2$, $\omega = .5$, $a = .03$, and $b = .06$. To create a stronger sense of feeling of knowing when participants were asked to answer questions before evaluating the performance of others, we used a value of $\mu_{FK} = .3$. For the answer after condition, where participants assessed performance before providing their own answers, we used $\mu_{FK} = .5$. As we do not have the raw experimental data available, the goal was not to pursue quantitative model fits and instead show that the model can capture the results from Tullis (2018) at a qualitative level. We found that a wide range of parameter values produce qualitatively similar model predictions. Note that we used the same set of parameters to generate model predictions for all the experiments from Tullis (2018) and Moore and Healy (2008) in Appendix F.

## Appendix F

## Simulation Details for Moore and Healy (2008)

This section provides details on the simulation studies we conducted to apply the hierarchical model to the experiment from Moore and Healy (2008). The authors present a synthesis of different ways in which overconfidence has been defined in the literature including the overestimation of one's actual performance and the overestimation of one's performance relative to others. The experimental results show that these forms of overconfidence manifest differently depending on the difficulty of the task. Since we do not have access to the raw data, we simulate data for the experiment presented in the article, including different levels of difficulty, and apply the hierarchical model to predict how people assess their own performance and place themselves relative to others.
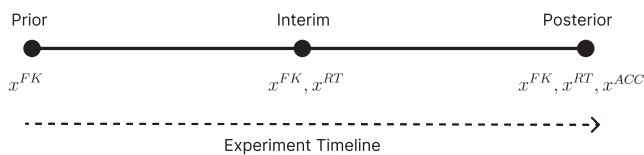
In the experiment, 82 participants answer 10 questions in 18 categories of trivia questions and predict their own score and the score of 1 randomly selected previous participant (RSPP) at three different stages of the experiment for each category. Figure F1 shows the timeline of the experiment and the hypothesized metacognitive signals available to participants when assessing their own performance and the performance of another person. First, participants made prior predictions about themselves and the RSPP before they had any specific information about the quiz they were about to take. Second, they answered 10 quiz questions from a category and then estimated their own scores and the RSPP's score again. This is termed their "interim" estimate. Next, participants are shown the correct answers to the quiz and asked to make "posterior" estimates about their performance and the RSPP's performance. Finally, they were given feedback about their own scores and the RSPP's scores.

We focus our model predictions on the interim stage of the experiment. We use the same process used for the Tullis data (Appendix E) with the same simulation parameters ($\mu_{FK} = .3$, $\mu_d = 1$, $\sigma_d = 2$, $\sigma_a = 0.5$, $\eta = .5$, $v = 2$) to generate the experimental data for 10 questions and 82 participants. Next, we apply the hierarchical model from Figure 1A, Equations E1–E2 and the same setup as used in Appendix E to obtain the participant's self and other estimates of the number of questions scored correctly out of 10 trivia questions, $\hat{x}^{o,ACC}$ and $\hat{x}^{s,ACC}$. We use a binomial link function to simulate these scores, $x^{ACC} \sim \text{Bin}(10, p_{i,j})$. On the basis of the simulated actual scores ($x^{s,ACC}$ and $x^{o,ACC}$) and the person estimated self- and other-performance ($\hat{x}^{o,ACC}$ and $\hat{x}^{s,ACC}$), we calculate two empirical measures used by Moore and Healy (2008). First, we assess the degree of *overestimation*, based on the participant's actual score subtracted from their estimated score, $\hat{x}^{s,ACC} - x^{s,ACC}$. Second, we assess the degree of *overplacement,* which measures whether a participant's assessment of themselves relative to others is in line with the actual observed difference: $(\hat{x}_i^{ACC} - \hat{x}_j^{ACC}) - (x_i^{ACC} - x_j^{ACC})$, where $\hat{x}_i^{ACC}$ is an individual's estimate of their own expected performance, $\hat{x}_i^{ACC}$ is their estimate of another person's expected performance on the same problem, and $x_i^{ACC}$ and $x_j^{ACC}$ refer to the actual scores of the individual and the other person.

**Figure F1**

*Timeline of the Experiment in Moore and Healy (2008) With the Hypothesized Metacognitive Signals Available to Participants Shown in Parentheses*



| Prior | Interim | Posterior |
|---|---|---|
| $x^{FK}$ | $x^{FK}, x^{RT}$ | $x^{FK}, x^{RT}, x^{ACC}$ |

Experiment Timeline

*Note.* FK = feeling of knowing; RT = response time; ACC = accuracy.