

Full length article

## People can infer the magnitude of other people's knowledge even when they cannot infer its contents

Rosie Aboody <sup>a</sup>,<sup>1</sup> Isaac Davis <sup>a</sup>,<sup>\*,1</sup> Yarrow Dunham <sup>a,b</sup>, Julian Jara-Ettinger <sup>a,b</sup>

<sup>a</sup> Department of Psychology, Yale University, United States of America

<sup>b</sup> Wu Tsai Institute, Yale University, United States of America

### ARTICLE INFO

Dataset link: <https://osf.io/pjy6x>

#### Keywords:

Computational modeling  
Social cognition  
Theory of mind

### ABSTRACT

Inferences about other people's knowledge and beliefs are central to social interaction. However, people's behavior is often consistent with a range of potential epistemic states, making it impossible to tell what exactly they know. Nonetheless, we are still often able to form coarse intuitions about how much someone knows, despite being unable to pinpoint the exact contents of their knowledge. Here, we sought to explore this capacity in humans, by comparing their performance to a normative model capturing this type of broad epistemic inference. We evaluated this capacity in a graded inference task where people had to make inferences about how much an agent knew based on the actions they chose (Experiment 1), and joint inferences about how much someone knew and how much they believed they could learn (Experiment 2). Critically, the agent's knowledge was always under-determined by their behavior, but the behavior nonetheless contained information about how much knowledge they possessed or believed they could gain. Our results reveal that people can make graded inferences about how much other people know from minimal behavioral data, but, interestingly, will sometimes achieve this through simpler approximations to the normative model that get the broad inferences right. Altogether, our paper reveals that people can make quantitatively precise judgments about the magnitude of an agent's knowledge from minimal behavioral evidence.

### 1. Introduction

Imagine that you and a friend are heading to dinner, but as you get into her car and she turns the key in the ignition, the engine sputters and dies. You take out your phone and offer to look up the problem on Google. Instead, your friend exits the car, enters her garage, walks past a toolbox sitting by the entrance, and picks up a second toolbox at the back of the garage. From this action, you might get a strong intuition that she must know enough about what just happened to not need to seek help; enough about the solution to identify which tools are needed; and enough about the contents of her garage to find the tools with minimal searching. Even if you cannot tell what exactly your friend knows (What is the problem? What tools are needed?), you can still get a sense of how knowledgeable she is in this domain. Inferences like these not only enable us to make sense of others' behavior, but also help us decide when to share what we know, and from whom to learn what we do not, forming a cornerstone of complex social action.

The ability to interpret other people's behavior in terms of mental states, called a Theory of Mind, has its origins in infancy and early

childhood. From infancy, we interpret other people's behavior as goal-directed (Woodward, 1998) and infer others' goals and preferences by assuming that agents act to maximize utilities—the difference between the costs they incur and the rewards they obtain (Csibra, 2003; Jara-Ettinger et al., 2016; Liu et al., 2017). Throughout our life, this expectation enables us to make a variety of judgments, such as inferring what others like (Jern et al., 2017; Lucas et al., 2014), predicting how they might behave (Jara-Ettinger et al., 2020), and determining their social affiliations (Davis et al., 2023; Jern & Kemp, 2014; Ullman et al., 2009).

Critically, inferences about others' minds are not restricted to goals and preferences: they also include judgments about what others may or may not know. Consistent with this, research in computational social cognition has found that people can make quantitative inferences about the contents of others' beliefs based on their behavior (Baker et al., 2017). This work showed that the same assumption of utility maximization, implemented in a Bayesian framework for action understanding, captures how people determine what an agent likely believes about

\* Corresponding author.

E-mail addresses: [rosie.aboody@yale.edu](mailto:rosie.aboody@yale.edu) (R. Aboody), [isaac.davis@yale.edu](mailto:isaac.davis@yale.edu) (I. Davis), [yarrow.dunham@yale.edu](mailto:yarrow.dunham@yale.edu) (Y. Dunham), [julian.jara-ettinger@yale.edu](mailto:julian.jara-ettinger@yale.edu) (J. Jara-Ettinger).

<sup>1</sup> These authors made equal contributions.

their environment given their behavior (e.g., if an agent looking for lunch walks towards the end of the block, peeks around the corner to see a Mexican food truck, and then turns around, we can infer that the agent was hoping to see a different food truck there).

While this work shows that people can make quantitative targeted belief inferences, such as determining whether an agent knew the type of food a vendor might be selling, these inferences often require access to a relatively constrained hypothesis space and key actions that reveal the agent’s beliefs. In many everyday situations, however, there may be a wide range of different belief states compatible with the behavior we observe, making it impossible to infer the specific contents of someone’s beliefs. In cases like these, our representations of other people’s epistemic states appear to consist of amorphous estimates of how much others know, without being sure exactly what it is that they know. Returning to the example in the introduction, when your friend rejects your offer to google the problem and goes straight for a specific toolbox, it is easy to infer that she knows more about the car than you do, despite not being able to tell what exactly she knows. In cases like these, where the exact contents of an agent’s beliefs are underdetermined by their behavior, can people make quantitatively precise inferences about how much an agent knows? Or are these inferences coarse and qualitative, providing no more than unreliable hints about others’ knowledge?

Research investigating people’s ability to quantify others’ knowledge—i.e., inferences about how much people know without knowing the exact epistemic content—has generally focused on children. By early in preschool, children can represent how much others know about a domain, without needing to list the full contents of their knowledge (Landrum & Mills, 2015; Lutz & Keil, 2002). Children can also infer broad domain knowledge from cues like confidence, accuracy, and generality (Koenig et al., 2015; Koenig & Harris, 2005; Kominsky et al., 2016), and appropriately infer knowledge (or lack thereof) from agents’ mistakes (Ronfard & Corriveau, 2016). However, to our knowledge, no work has explored our capacity to infer knowledge magnitude from others’ actions, or specified the computations that might underlie such inferences.

Here we propose that such inferences are part of our broader quantitative inferential system within Theory of Mind, and are therefore supported by the same utility-maximization expectation that explains more targeted belief inferences. Specifically, given the expectation that agents choose actions which (they believe) fulfill their goals as efficiently as possible, an agent’s choice of action can reveal what that agent believes to be efficient, which can, in turn, provide indirect evidence about the magnitude of their knowledge. We propose that even when this evidence is insufficient to reveal what exactly an agent knows, adults can still leverage this expectation to infer how much an agent knows with quantitative precision (as we explain further in our computational framework).

In this paper, we present a computational model of epistemic quantification, built on the same utility-maximization framework that explains more targeted epistemic inferences. We tested this account on two sets of tasks where participants inferred how much someone knows or thinks they can learn based on very minimal observations of behavior (a single binary choice). Our work shows that people can seamlessly make finely-graded quantitative estimates of how much someone knows or expects to learn, which closely track with the predictions of our computational model. Critically, we also compared participant judgments against simpler heuristic accounts, which make coarser predictions about agents’ knowledge (e.g.: that they know “a lot” or “a little”) without considering the agent’s epistemic utilities. These alternate accounts failed to capture the graded structure of participant judgments, though our results also suggest that participants may strategically revert to simpler heuristics in cases where precise knowledge estimates are perceived to be less important.

## 2. Computational framework

Our computational framework builds on a recent family of computational models of mental-state inference structured around an expectation that agents act rationally: that is, we expect agents to act in a way that maximizes the difference between the rewards they receive and the costs they incur. By formalizing this expectation as a generative model of utility maximization, and inverting this model via Bayesian inference, we can make a wide range of inferences about an agent’s mental states based on their behavior (Baker et al., 2017; Jara-Ettinger et al., 2020; Jern et al., 2017; Lucas et al., 2014). We extend this framework by proposing that adults often expect agents’ costs to be mediated by their knowledge: that is, if there are multiple ways to reach the same goal, an agent’s choice of action can reveal what they know about the available strategies, and the relative cost of each strategy. Critically, we propose that even when the agent’s choice does not reveal precisely *what* they know (if, for example, there are many knowledge states compatible with that action), adults can still leverage this expectation to make precise inferences about *how much* the agent knows, by observing the apparent costs they choose to incur.

For simplicity, we will explain our framework within the context of our Experiment 1 paradigm. In these scenarios, an agent must choose one of two different fields for an Easter egg hunt. Each field contains a different spatial and numerical configuration of eggs (see Fig. 1), and exactly one egg in each field contains a prize, while all other eggs are empty. Suppose that the agent arrived while the fields were being set up, and was able to see the contents of some of the eggs in each field (either empty or full). We define the agent’s knowledge state  $k$  as the subset of eggs in a field that the agent saw. The cost the agent incurs depends on their search trajectory, and we assume that agents navigate efficiently in space (Csibra, 2003), searching only locations where they think they may find the prize. This implies that when the agent’s knowledge state includes information about the contents of the prize egg, the agent will always move directly towards the prize. Otherwise, the agent will search the eggs in a way that minimizes the expected search time. Let  $k_1$  denote the subset of eggs in field 1 that the agent observed, and similarly for  $k_2$ .

Given this knowledge, the agent can compute the expected cost of finding the prize in each field, which we assume is equal to the expected distance traveled before finding the prize, plus a small fixed cost  $C$  of opening each egg to check its contents. If the agent’s knowledge for a field includes the egg  $e_i$  that contains the prize, then the cost of finding the prize in that field is simply the distance  $\text{dist}(e_0, e_i)$  from the entrance  $e_0$  to the target egg  $e_i$ , plus the cost  $C$  of opening the egg. Now suppose that the agent’s knowledge specifies that eggs  $k = \{e_1, \dots, e_k\}$  are empty, and that the prize must be in one of the remaining eggs  $k^c = \{e_{k+1}, \dots, e_n\}$ . Let  $\pi$  be a path that starts at the entrance and passes through each egg in  $k^c$ , and let  $\pi_i$  denote the  $i$ th stop of  $\pi$ , so that  $\pi_0$  is the entrance to the field, and  $\pi_i$  is the  $i$ th egg on the path. The cost of traversing the entire path, stopping to check each egg, is

$$\text{cost}(\pi) = \sum_{i=1}^{|k^c|} \text{dist}(\pi_i, \pi_{i+1}) + C \quad (1)$$

where  $\text{dist}(a, b)$  is the distance from point  $a$  to point  $b$ . Most of the time, however, the agent will not have to traverse the full path, as they can stop once they find the egg containing the prize. Assuming that each egg has equal probability of containing the prize, such that  $P(\text{prize in egg } i) = 1/|k^c|$  for all  $i$ , then the expected cost of finding the prize along path  $\pi$  is equal to

$$E[\text{cost}(\pi)] = \sum_{i=1}^{|k^c|} \frac{1}{|k^c|} * \text{cost}(\pi|_i) \quad (2)$$

Here,  $\pi|_i$  is the sub-path obtained by following  $\pi$  until the  $i$ th egg, then stopping. Thus, the expected cost of finding the prize along path  $\pi$  is equal to the sum of the costs of traversing each sub-path  $\pi|_i$ , weighted by the probability that the prize is in the  $i$ th egg along path

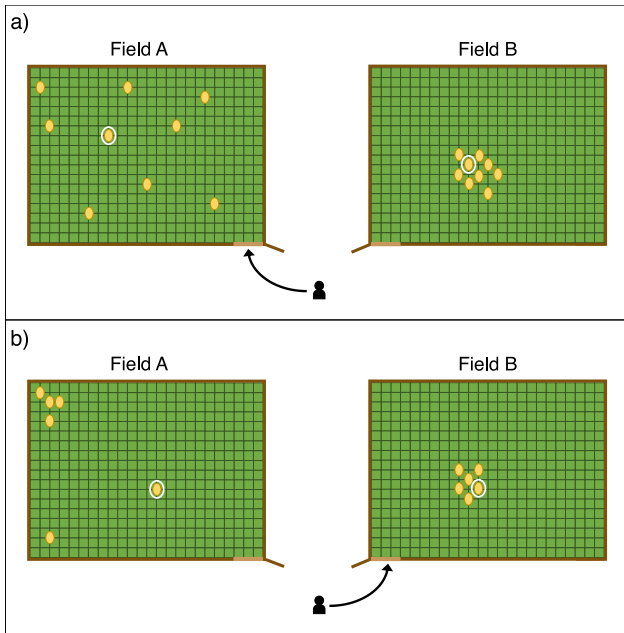


Fig. 1. Example of the experimental stimuli. The arrow indicates the agent's chosen field; eggs containing a prize are circled. Panel A depicts a strong epistemic contrast: here, you might infer that the agent knows approximately where the prize is located in their chosen field, and very little about the other field. Panel B depicts a more graded contrast: here, you might suspect that the agent knows more about the prize's location in their chosen field, but may be less certain they know a lot (because their chosen field is also much less costly to search).

$\pi$ . Given that people expect each other to act rationally and efficiently, we assume that the agent will compute the search path that minimizes the expected cost of finding the prize in field  $X$ , which we refer to as  $E[\text{cost}(X)|k]$ . If the reward of getting the prize is equal to  $R$ , then the total expected utility of an agent with knowledge state  $k$  choosing field  $X$  is equal to  $U_X = R - E[\text{cost}(X)|k]$ .

Now suppose that the agent computes the expected utility for each field,  $U_1$  and  $U_2$ . We assume that agents will generally try to maximize their expected utilities, but are not deterministic and may be prone to errors (e.g.: due to distraction or errors while computing expected costs). Thus, rather than assuming the agent will always choose the field with higher expected utility, we make a standard assumption that the agent will choose a field with probability

$$P(\text{choice} = \text{field}_i | k) \propto e^{U_i / \tau} \quad (3)$$

This is the standard softMax function, which takes a vector of real numbers (in this case, the expected utilities) and converts it into a probability vector. The “temperature” parameter  $\tau$  controls the agent's level of rationality: a very high value of  $\tau$  entails nearly uniform behavior (i.e.: choosing each option with equal probability), while a very low value entails nearly deterministic behavior (i.e.: choosing the highest utility option with probability near 1). Thus, Eq. (3) specifies the probability that an agent with knowledge states  $k_1, k_2$  (about fields 1 and 2, respectively) will choose to enter each field.

Given this generative model of the agent's behavior, a Bayesian observer can infer the agent's knowledge of each field  $k_1, k_2$  based on the field configurations and the agent's choice according to Bayes' rule:

$$P(k_1, k_2 | \text{choice}, \text{field}_1, \text{field}_2) \propto P(\text{choice} | k_1, k_2, \text{field}_1, \text{field}_2) P(k_1, k_2) \quad (4)$$

Here,  $P(k_1, k_2 | \text{choice}, \text{field}_1, \text{field}_2)$  is the posterior probability of the agent's knowledge states,  $P(\text{choice} | k_1, k_2, \text{field}_1, \text{field}_2)$  is the likelihood

of the agent's choice given these knowledge states (given by Eq. (3)), and  $P(k_1, k_2)$  is the prior probability of the agent having these knowledge states.

In our scenarios, however, the richness of the agent's possible knowledge states (all possible subsets of eggs in each field) and the coarseness of the agent's behavior (a binary choice between two fields) make the exact contents of the agent's knowledge highly underdetermined by the observed behavior. That is, there will always be a large number of possible knowledge states compatible with the agent's choice. But even when we cannot infer the precise contents of others' knowledge representations, we may still be able to infer approximately how much they know (getting a rough sense of how knowledgeable they are). Thus, given a posterior distribution over what the agent might know (Eq. (4)), we formalize the quantity of amorphous knowledge  $Q$  as the expected quantity of knowledge encoded in the probable epistemic states that the agent has, given by Eq. (5) below.

$$Q = \sum_{k \in K} |k| p(k | \text{choice}) \quad (5)$$

Here,  $K$  is the set of all possible epistemic states,  $|k|$  is a quantification of how much the agent knows in that state, and  $p(k | \text{choice})$  is the posterior probability of that knowledge state (Eq. (4)). Naturally, precisely defining the measure  $|k|$  may be highly context-sensitive. Here we focus on its application in a particular experimental context but return to the idea of how this might generalize in the discussion.

We evaluate this framework in two experimental paradigms. The first paradigm tests people's capacity to infer how much someone knows about two related environments based on which one they choose to seek a reward in. The second paradigm tests people's capacity to jointly infer how much someone knows and how much they expect to learn based on whether they seek additional information before trying to attain a reward. Additional details about the inference procedure can be found in each experiment.

### 3. Experiment 1

To test our model, we designed a task where an agent's behavior (and its costs) could reveal approximately how much they knew—but was too impoverished to reveal precisely what they knew. Specifically, participants watched an agent choose which of two fields to search for a prize hidden in an easter egg, knowing that each field had only one egg with a prize inside (and that the reward was always the same in every field).

The expected cost of locating the prize in any given field was determined by the number of eggs, their spatial distribution, and the true location of the prize. By manipulating all three variables, we test if participants infer how much others know by quantifying and comparing their expected costs—or whether participants rely on a simpler heuristic that does not require them to track or reason about others' costs when inferring epistemic states. Our procedure, stimuli, sample size, and analysis plan for our main model were preregistered (see OSF project page).

#### 3.1. Model parameters

Our main model has four parameters: the reward of obtaining the prize, the cost of checking an egg's contents upon reaching it, a prior over the agent's knowledge, and the softmax parameter ( $\tau$ ). All parameter values and model predictions were preregistered prior to data collection.

The reward function for the prize is the same across fields, and we set it as a constant  $R(a_i) = 100$ . Because the reward is constant across action plans, the difference in utilities between the two plans would be unchanged by different reward functions. We simply selected (and preregistered) a reward function large enough to ensure that no action plan could have a negative utility.

For each knowledge state sample, the cost of stopping to check an egg's contents was modeled as a continuous uniform distribution [1, 3]. This range was chosen to capture the expectation that stopping to open an egg does incur some minimal cost, but that its precise value is unknown.

To define a prior over the agent's knowledge, we assumed that the agent had a 50% chance of knowing each egg's contents. We also explicitly communicated this to participants in our task (see Procedure) to ensure that participants and the model both relied on similar epistemic priors. Finally, we selected a softmax  $\tau$  value that produced graded action predictions in proportion to each plan's expected utility ( $\tau = 3$ ): that is, we chose  $\tau$  to ensure that the model could predict both "obvious" cases where a single answer is very certain, as well as "ambiguous" cases where both options might be plausible, with one slightly more likely than the other.

We implemented our inference procedure via Monte Carlo sampling, drawing 10,000 knowledge states from each field. We then computed Eq. (5), by quantifying the amount of knowledge in an epistemic state as  $1 -$  the proportion of eggs the agent is still uncertain about (if the agent knows where the prize is, they know the rest of the eggs are empty, and thus the proportion known is 1; if the agent is unsure about half of the eggs, the proportion known is .5; and so on).

### 3.2. Alternate model

Our main model assumes that people quantify the cost of obtaining the prize in each field under different degrees of knowledge, and then reason about the knowledge states under which the agent's actions would have been utility-maximizing. However, it is possible that adults generally do not apply such complex computations when inferring others' knowledge states, and instead rely on simpler rules or heuristics. Such heuristics could get things approximately right most of the time, while requiring less effort to apply.

To address this possibility, our alternate model encoded the simple heuristic that agents tend to choose options they know more pieces of information about. Critically, this alternate model did not consider agents' knowledge states in a full mentalistic way: it did not compute the utility of each field based on the agent's knowledge state, and did not expect agents to navigate directly to an egg if they knew it contained the prize. It simply considered the proportion of eggs with known contents in each field, and expected the agent to always choose the field where this proportion was larger (or choose randomly when this proportion was equal across fields). We then generated predictions from this alternate model using the same sampling procedure as in the main model.

Our alternate model was not preregistered, but uses only one parameter: the same knowledge prior as in our main model. Because our alternate model encodes an expectation that agents will always choose fields they know more pieces of information about, we do not compute the utility of each field, and thus we do not need to specify agents' costs, rewards, or a softmax parameter.

### 3.3. Participants

40 adult participants with U.S.-based IP addresses were recruited via Amazon Mechanical Turk ( $M = 35.05$  years,  $SD = 9.23$ ). 7 additional participants were recruited but excluded from the study for failing a preregistered inclusion trial.

We based our sample size on closely related research in computational social cognition which relies on  $n = 20-30$  (e.g., (Baker et al., 2017; Jara-Ettinger et al., 2020; Ullman et al., 2009), but increased our sample to  $n = 40$  to be conservative. Our sample size was preregistered. All data were collected in January and February 2021.

### 3.4. Stimuli

Stimuli consisted of 19 test trials, plus one inclusion trial. The test trials were presented in a randomized order, and the inclusion trial was always presented last. Each trial showed an agent, and two fields. The fields each had easter eggs placed inside, and one egg in each field contained a hidden prize. This egg was circled for participants. An arrow indicated the agent's path to their chosen field, thus showing which field the agent chose to visit on each trial (see Fig. 1).

Stimuli were based on three scenarios (pairs of fields) we thought could elicit a range of model ratings. To manipulate the cost of searching each field, eggs in the first field (field A) were always wide-spread. The second field (field B) contained the same number of eggs, but these eggs were instead clustered near the middle of the field. All else being equal, a denser cluster of eggs should be less costly to search than a loosely scattered cluster, as the denser cluster requires less travel between each egg. The first scenario is shown in Fig. 1a. The second scenario was based on the first: we selected a subset of 6 eggs from each field, thus varying the number of eggs but not their position. The third scenario was in turn based on the second, but here we instead varied the position of the eggs in field A (capturing a case where most of the eggs in field A were extremely costly; see Fig. 1b).

To select the final locations of the prize in field A, we provided each scenario as input to the model, but systematically varied which egg in field A contained the prize, yielding 42 trials (21 unique scenarios  $\times$  2 choices per scenario).<sup>2</sup> We selected 24 trials (12 unique scenarios  $\times$  2 choices per scenario) that both produced a range of model responses, and were not too similar to each other. In preparation to present stimuli to participants, some trials were mirrored, and we slightly varied the position of the prize in field B amongst similar scenarios (to prevent participants from noticing similarities between trials).<sup>3</sup> We then obtained final model predictions, and excluded any trials where the model's knowledge predictions were based on less than 500 samples (that is, where the predicted choice of field was consistent with the observed choice in less than 5% of cases). This yielded 19 final trials; this criterion and our final set of stimuli were preregistered.

### 3.5. Procedure

Participants were introduced to an agent going on easter-egg hunts in a two-dimensional grid-world. Participants learned that a farmer had placed easter eggs in his fields, hiding a prize inside one egg in every field. This prize (one silver token) was always the same in every field, and the prize egg was always circled for participants.

Participants learned that because the grass in the fields was quite short, the agent could always see where the eggs were located in a field before entering it. But while the prize egg was circled for participants, the agent did not necessarily know which egg contained the prize. Participants learned that the agent had seen the farmer set up some of the eggs; it was unclear what prior over knowledge participants would bring to the task, so we specified that the agent had a 50/50 chance of knowing the contents of any given egg, and that these probabilities

<sup>2</sup> We did not expect the location of the prize in field B to strongly affect the model's predictions; to test if this was the case, we did also replicate one scenario given a different prize location in field B, yielding an additional 18 additional trials. The location of the prize in field B indeed had little effect (as all of these eggs are so close to each other), and thus we selected our final stimuli by considering primarily the location of the prize in field A.

<sup>3</sup> Despite slightly varying the prize's location in field B across similar trials in our preregistered stimuli, our model predictions were accidentally not updated accordingly prior to preregistration. Because we collected our data using the preregistered stimuli, we obtained new model predictions for any trials where the location of the prize in the stimuli did not match the coordinates originally used in the preregistered model predictions. No aspect of the model itself was modified and we used the same preregistered parameters.

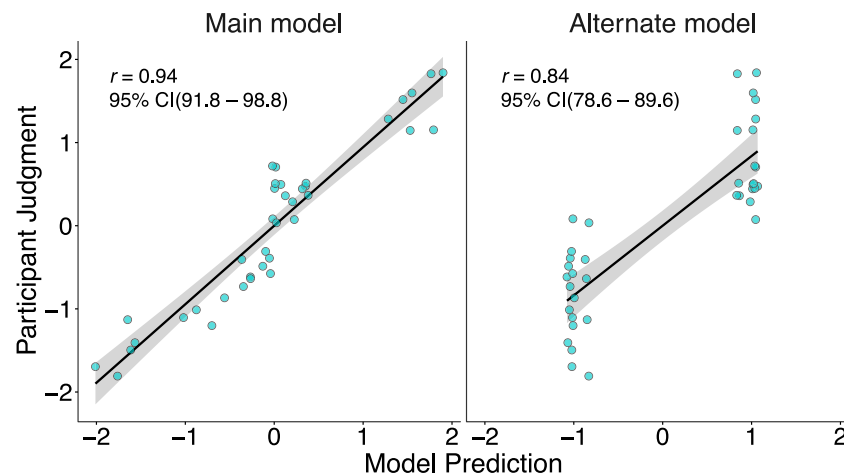


Fig. 2. Comparison of main and alternate models, with linear regressions fit to each dataset. Each point represents one knowledge rating, with model predictions on the x axis and participant judgments on the y axis. Gray bands show 95% confidence intervals in the regression.

were independent (i.e.: knowing the contents of one egg does not affect the probability of knowing the contents of any other egg). Additionally, participants were explicitly instructed that the agent did not always know the same amount about every field; the amount she knew about the location of the prize in each field could differ.

Participants learned that the agent always had to choose between two fields, and could only search the field she chose. An arrow indicated which field the agent had chosen to search (see Fig. 1). Participants were oriented to factors that might affect the agent's search decision: they were told that the agent always wanted to find the prize as quickly and easily as possible, and that the difficulty of finding the prize was determined by the number of eggs in a field, their distance from the entrance, and the amount the agent already knew about the location of the prize. Note that while this tutorial ensured participants were attentive to the main features of our task, we are interested in how participants combine these different pieces of information and reason over them to infer what others know. The tutorial did not specify how participants should weight or use any of these features in their judgments.

To access the task, participants then completed a preregistered inclusion quiz that assessed their understanding of the task instructions. Participants were given two chances to pass the inclusion quiz; those who failed on their first attempt were required to review the task introduction before trying again. Participants who failed both attempts were not given access to the task. Upon passing the inclusion quiz, participants then completed the 19 test trials (presented in a randomized order), plus one inclusion trial at the end. For each trial, participants were asked to rate, on a sliding scale from 0–100, "For the field to the LEFT: How much did [the agent] already know about the prize's location?" and an identical question for the field to the right (on a separate sliding scale): "For the field to the RIGHT: How much did [the agent] already know about the prize's location?". Critically, participants rated how much the agent knew about both fields, not just the field she had chosen. The preregistered inclusion trial always came last. It was similar to the test trials, but presented an extreme contrast where the agent's choice only made sense under the assumption that they have a high degree of knowledge about the chosen field and a low degree of knowledge about the other field (see Supplemental Materials for more details). This allowed us to make a strong prediction about the pattern of judgments an attentive participant should make: as preregistered, any participant who judged that the agent was more knowledgeable about the other, not chosen field was excluded. Finally, participants were asked what they thought the point of the task had been, and were given an opportunity to provide feedback or note any technical difficulties.

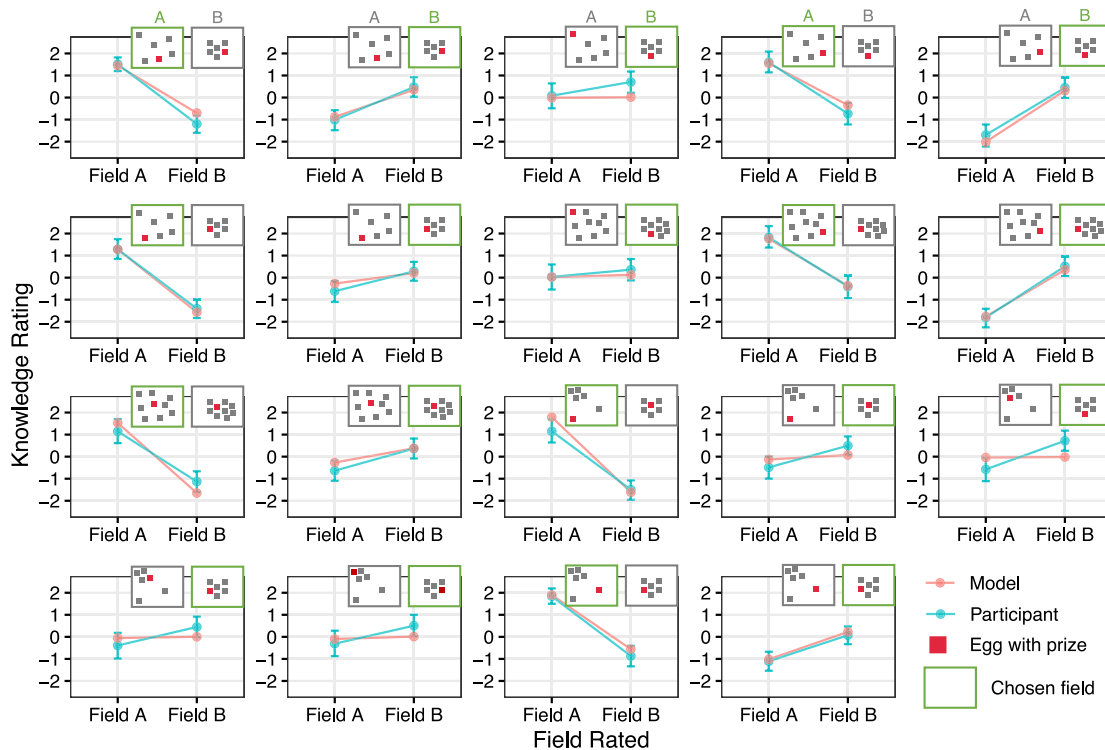
### 3.6. Results

Participants rated the agent's knowledge about both fields in 19 test trials, yielding 38 ratings. As preregistered, participant responses were averaged by question, and then z-scored; the corresponding model predictions were also z-scored.

Fig. 2 shows the overall results, revealing that our model was highly correlated with participant judgments,  $r = 0.94$  (95% CI: 91.8, 98.8). Our main model also outperformed the alternate heuristic model,  $r = 0.84$ , (95% CI: 78.5, 89.5), and a bootstrap over the correlation difference revealed that this difference was statistically significant (correlation difference, alternate model – main model =  $-0.11$ , 95% CI:  $-17.4, -4.3$ ; not preregistered). Critically, participants' judgments about the agent's knowledge were not "all-or-nothing", and reflected a graded range of inferences about the agent's degree of knowledge.<sup>4</sup> These graded inferences were captured by our main model, but not the alternate model, which produced tightly clustered, bimodal predictions (as illustrated in the right panel of Fig. 2). This supports our hypothesis that participants are considering the layout of each field, and the cost of searching each field under different epistemic states, in order to infer the agent's quantity of knowledge.

This gradedness in participant responses is further illustrated in Fig. 3, which plots the trial-by-trial correspondence between model and participant ratings, and shows how, even between trials with the same qualitative pattern of inferences (e.g.: that the agent is more knowledgeable about Field B than Field A), participants were still able to distinguish different degrees of knowledge. For example, Fig. 4a) highlights a pair of trials where participants (and both models) inferred that the agent knew more about Field B than Field A. However, this inferred difference was much larger in the trial on the right than the trial on the left, which closely tracks with the predictions of our main model, but is not captured by the alternate model. This is due to the fact that the main model is sensitive to the layout of each field, and how that layout impacts the expected search cost under different epistemic states. In the left trial, the denser configuration of Field B should make it generally easier to search even with minimal knowledge about either field. Thus, choosing Field B in this case is not a strong signal that the agent is knowledgeable about that field, as reflected by participant responses. In the right trial, however, Field A would likely be easier to search with only a little bit of knowledge about it. Since the agent chose Field B, this suggests that they were especially

<sup>4</sup> See Supplemental Materials for participant-level scatterplots illustrating the gradedness of participant responses.



**Fig. 3.** Detailed results for Experiment 1. Each panel presents one trial, with results split by the field rated (Field A or Field B, indicated on the x axis). The y axis indicates standardized knowledge ratings. Participant judgments are plotted in blue; model predictions are plotted in red. Vertical bars show 95% confidence intervals over participant judgments. The schematics show the position and number of eggs in each field, the egg with the prize, and the field the agent ultimately chose in each trial. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

un-knowledgeable about Field A, as predicted by participants and the main model. The alternate model, however, which is insensitive to the expected cost of searching each field, and simply assumes that the agent will choose whichever field they are more knowledgeable about, produced identical predictions for these two trials, essentially inferring that the agent knew “a lot” about Field B and “a little” about Field A, but without the quantitative precision of the main model. This suggests that, in line with our hypothesis, participants were in fact tracking the expected cost of finding the prize in each field under different epistemic states, and inferring quantity of knowledge by considering under what epistemic states the agent’s choice would have been utility-maximizing, rather than simply assuming the agent is more knowledgeable about whichever field they chose.

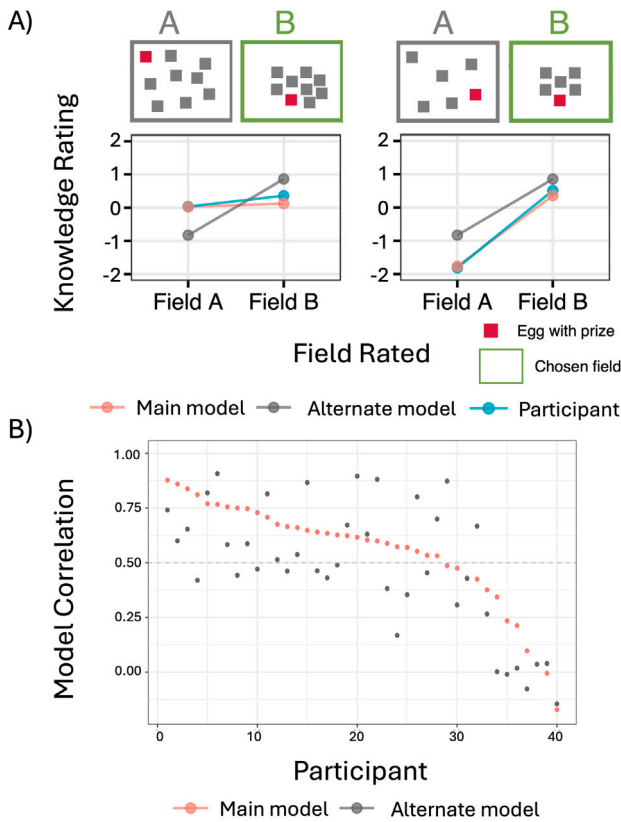
As a final check to ensure that this gradedness was genuinely reflected by participant responses, and not simply a result of averaging bimodal individual responses across participants, we computed the correlation between each model’s predictions and each participant’s responses at the individual level (not pre-registered). As illustrated in Fig. 4b, despite the additional noise typically present in un-averaged, individual-level responses, participants still showed a strong correlation with the predictions of our main model, with 70% ( $n = 28/40$ ) of participants showing a correlation of at least  $r = .5$ . Furthermore, 63% ( $n = 25/40$ ) of participants showed a stronger correlation with the main model than the alternate model, suggesting that most participants were in fact thinking about the expected search costs under different epistemic states, and not relying on a simple heuristic akin to our alternate model.

Despite the high overall fit of our main model, Fig. 3 reveals several trials where the model predictions diverged from participant responses. In particular, there were several trials where participants rated the agent as more knowledgeable about the field they chose, while the main model rated the agent as having roughly equal knowledge about both fields. In each of these trials, the agent chose a more densely packed field over a more scattered field. Because a densely packed field is

generally less costly to search than a scattered field, the denser field is often the correct choice even if the agent knows very little about both fields. Thus, choosing the denser field in these trials is not an especially strong signal that the agent knows more about that field. The fact that participants recognized this in some trials but not others may suggest that, although they did consider the expected search costs in a manner consistent with our hypothesis, they may have employed simpler heuristics to estimate those search costs. We return to this possibility in the general discussion.

#### 4. Experiment 2

Experiment 1 shows that adults are able to make precise epistemic inferences even in underdetermined scenarios—and that these inferences are well-captured by our main model. Experiment 2 both conceptually replicates and extends these findings. Specifically, in Experiment 2 we test whether our framework can capture not just adults’ inferences about how much someone knows, but also about how much they believed they could learn. To do so, we designed a task where an agent’s information-seeking choice (and its cost) could reveal approximately how much they knew and believed they could learn (but again, could not reveal these states with any precision). Specifically, participants watched agents search islands for hidden treasure (Fig. 5). Agents had the option to obtain a treasure map first, or to skip the map and go straight to the island. Importantly, the map was not always informative: sometimes it might contain a lot of information about the treasure’s location, sometimes it might contain a little, and sometimes it might contain no information at all. To elicit graded inferences, we manipulated the distance of the map (varying information’s cost), the size of the island (varying the potential difficulty of finding the treasure), and agents’ information-seeking choices (varying whether or not they pursued the map). Our procedure and sample size were preregistered.



**Fig. 4.** Panel (A) Two example trials from Experiment 1, showing main model predictions (red), participant judgments (blue), and alternate model predictions (gray). Figures above each plot show the layout of eggs in the field for that trial, with the prize egg highlighted in red, and the field chosen by the agent outlined in green. Panel (B) Individual-level correlations between participant judgments and model predictions, with participant number on the x-axis and model correlation on the y-axis. Main model shown in red and alternate model shown in gray. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

#### 4.1. Model structure and parameters

The computational model followed the same conceptual structure as Experiment 1. The key difference was that the two choices no longer referred to two possible search areas. Instead, the agent had to choose between searching the island directly, or first taking a detour to obtain a map containing potential information about the treasure's location, before returning to the island to search for the treasure. To compute the expected cost of searching the island, we divided it into  $N$  cells of equal size, and modeled the agent's knowledge as the number  $k_p \leq N$  of those cells that the agent knew might contain the treasure (e.g.:  $k_p = 1$  indicates full knowledge, as the pirate knows the treasure to be in one specific cell, while  $k_p = N$  indicates full ignorance, as the treasure could be in any of the  $N$  cells). This allowed us to compute the expected search cost of the island,  $E[\text{cost}|k]$ , given a particular knowledge state  $k$ , similarly to Experiment 1: given the cost  $C_{\text{search}}$  of searching one cell, and the number  $k$  of cells that might contain the treasure, the expected cost of finding the treasure is

$$E[\text{cost}|k] = \sum_{n=1}^k n * C_{\text{search}} * \frac{1}{k} \quad (6)$$

where  $n * C_{\text{search}}$  is the cost of searching  $n$  cells, and  $\frac{1}{k}$  is the probability of finding the treasure in the  $n$ th cell (assuming the treasure is equally likely to be in any one of the  $k$  possible cells).

To compute the utility of obtaining the map, let  $k_p$  denote the pirate's initial knowledge and  $k_m$  denote the knowledge contained in

the map. If the pirate obtains the map, their new quantity of knowledge is given by  $k' = \min(k_p, k_m)$  (recalling that lower values denote more knowledge). Thus, the overall utility of obtaining the map, then searching the island, is the expected search cost under the augmented knowledge state  $k'$ , plus the deterministic cost  $C_{\text{travel}}$  of traveling to and from the map. The expected utility of each choice is therefore  $E[\text{cost}|k_p]$  for going to the island directly, and  $E[\text{cost}|k'] + C_{\text{travel}}$  for obtaining the map first. Given the two expected utilities, we assume the agent will choose an action using a softmax decision policy with temperature parameter  $\tau$ , as with the model for Experiment 1. Given the agent's choice, we can then use Bayesian inference to jointly recover (1) the agent's original amount of knowledge about the island and (2) the amount of information they expected the map to contain, according to

$$P(k_p, k_m | \text{choice}, N, C_{\text{search}}, C_{\text{travel}}) \propto P(\text{choice} | k_p, k_m, N, C_{\text{search}}, C_{\text{travel}}) P(k_p, k_m) \quad (7)$$

where  $P(k_p, k_m)$  is a prior probability distribution over knowledge states. We defined this using a uniform prior over the probability that the map might contain each degree of knowledge, and a non-uniform prior over the probability that the pirates might have each degree of knowledge (not preregistered). This was intended to capture the possibility that adults might generally expect agents to be knowledgeable (and unlike in Experiment 1, we did not specify precisely how likely agents were to know the contents of each island square). We defined this prior using the binomial distribution with success rate  $p = .8$ .

Our main model for Experiment 2 has six parameters: the reward of obtaining the prize (set to a constant  $R(a_i) = 100$ ), the cost of sailing across one grid square, the cost of searching one island square, the softmax parameter ( $\tau$ ), and the two priors defined above. We preregistered the first four parameters prior to data collection, basing the relative cost of sailing vs. searching upon empirical estimates from a pilot sample. Our pilot sample judged that searching one island square was, on average, 2.25x more difficult than sailing across one ocean square, and thus we preregistered a sailing cost of 1, a searching cost of 2.25, and  $\tau = 4$  (based upon the range of utilities these costs produced). However, we explicitly preregistered that we would re-estimate these based on our final sample, and re-adjust our softmax parameter if needed. In our final sample, most participants judged that searching was more difficult than sailing, judging that it was on average 3.9x harder. Thus to generate our final predictions, we set the cost of searching to 3.9. Because this affected the range of possible utilities, as preregistered we adjusted our softmax parameter, setting  $\tau = 6.5$ .<sup>5</sup>

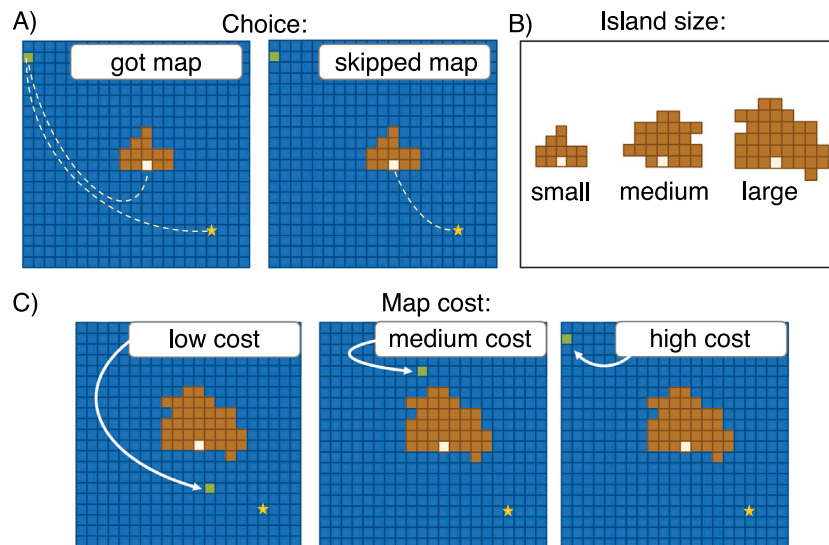
#### 4.2. Alternate model

Our preregistered alternate model is a linear regression, trained on participants' z-scored average ratings in our task. It predicts knowledge based on an interaction between agents' information-seeking choice (to retrieve the map/skip the map), and the type of knowledge (what agents know/what information they believe the map contains). The formula for this regression in R is:  $\text{lm}(\text{mean participant rating} \sim \text{choice} * \text{knowledge category})$ .

#### 4.3. Participants

40 adult participants with U.S.-based IP addresses were recruited via Amazon Mechanical Turk ( $M = 38.73$  years,  $SD = 12.23$ ). 9 additional participants were recruited but excluded from the study for failing a preregistered inclusion trial. As in Experiment 1, we based our sample size on closely related research in computational social cognition; our sample size was preregistered. All data were collected in May 2020.

<sup>5</sup> Note that Experiment 2 was conducted before Experiment 1. The preregistered procedure for Experiment 1 was simpler due to the realization that the tau parameter did not particularly matter for our predictions.



**Fig. 5.** Space of all possible experimental stimuli. We varied (A) agents' choices (to pursue information or ignore it), (B) the size of the island to be searched (small, medium, large), and (C) the cost of pursuing information (small, medium, large). This yielded 18 test trials. The first choice (panel A, left) depicts a strong epistemic contrast: here, you might infer the agents knew relatively little, and believed they stood to gain a lot of information (because they chose to incur a high cost to obtain the map, even though the island was small and thus relatively easy to search). The second choice (panel A, right) depicts a more graded contrast: while the agents clearly did not think the map was worth it, it may not be entirely clear why (did they know a lot, or did they simply believe the island would be easy to search even given ignorance?).

#### 4.4. Stimuli

Stimuli consisted of 18 test trials, plus two inclusion trials. The test trials were presented in a randomized order, and the inclusion trials were always presented last. Each trial showed a pirate ship (represented by a yellow star), a treasure map (represented by a green square), and an island (represented by brown squares); see Fig. 5. Each island had a beach (represented by a lighter brown square), which was the only point on the island pirates could land their ship. An arrow indicated agents' path, showing whether they chose to pursue added knowledge (obtaining the treasure map first), or whether they chose to search the island without obtaining the map (see Fig. 5a).

To construct our stimuli space, we varied the size of the island pirates needed to search (12, 24, or 36 grid-squares), the detour required to obtain the treasure map (adding approximately 10, 20, or 40 grid-squares to the journey), and agents' choices to obtain or skip the map. This yielded 18 test trials which systematically varied information's cost (as well as agents' information-seeking choices).

#### 4.5. Procedure

Participants were introduced to pirates searching for treasure in a two-dimensional grid-world. Participants were shown how to identify the pirate ship (marked by a star), and learned that pirates could only land on the island at the beach (this was intended to explain why the pirates sometimes took circuitous, high-cost paths to the island; e.g., see Fig. 5a). Participants learned that pirates sometimes knew a lot about the treasure's location, sometimes knew a little, and often knew something in between.

Participants learned that islands could be all different sizes, and that there was always a map somewhere in the ocean, marked by a green square. However, this map was not always helpful: sometimes it contained a lot of information about the location of the treasure, sometimes it contained only a little, and often it contained something in between. To obtain the map, pirates needed to sail to the green square first, before going to the island. An arrow indicated pirates' final choice (showing their chosen path).

Participants were oriented to factors that might affect agents' information-seeking decisions: they were told that the less pirates knew, the more work it might take to locate the treasure; the bigger

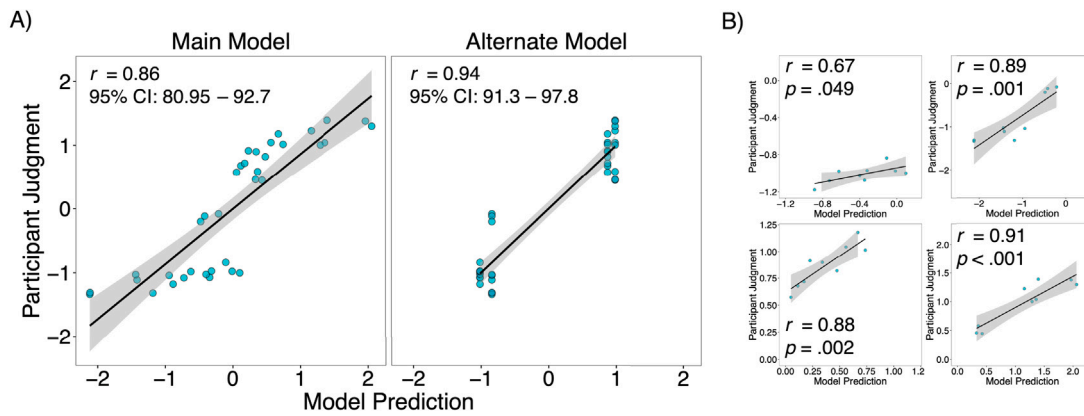
the island, the more work it might be to search for treasure; and the farther the map, the more time and effort might be required to obtain it. Participants were explicitly told that, in each case, the pirates needed to decide whether it was worthwhile to pursue the map. As before, note that while this tutorial ensured participants were attentive to the main features of our task, we are interested in how participants combine these different pieces of information and reason over them to infer what others know and believe they can learn. The tutorial did not specify how participants should weight or use any of these features in their judgments.

Before the task, participants completed three simple attention check questions that assessed their understanding of the task instructions. Participants were asked to identify how the pirate ship was marked (by a star), to recall the pirates' goal (find treasure), and finally were asked to identify both that the map was always on the green square, and that pirates could only get on an island via the beach (distinguishing these from three other incorrect statements). Participants were able to select as many answers as they chose to each question; however, attentive participants should have noticed that the first two questions could only have one correct answer. Any participants who selected more than one answer in response to these two questions was excluded (preregistered). Participants who answered any question incorrectly were corrected.

Finally, participants were again reminded that both the pirates' knowledge and the informativeness of the map might vary, and that in each case, pirates needed to decide whether it was worthwhile to pursue the map. For each trial, after observing pirates' information-seeking choices (and their expected costs), participants were asked to rate, on a sliding scale from 0–100, “How much did the pirates already know about where to find the treasure?” and “How much information did the pirates think the map had about where to find the treasure?”

Two inclusion trials always came last. These were similar to the test trials, but presented an extreme contrast where the pirate's decision only made sense if they were extremely knowledgeable (in the first trial) or mostly ignorant (in the second trial) about the treasure's location, allowing us to make a strong prediction about the pattern of judgments an attentive participant should make. Participants whose judgments differed from this pattern were excluded, as preregistered.

Participants were also asked to judge which was more difficult: to sail across one ocean square, or search one island square for treasure. After identifying which was harder, participants were asked to judge



**Fig. 6.** (A) Comparison between our model and the alternate model, with linear regressions fit to each dataset. Each point represents one knowledge rating, with model predictions on the x axis and participant judgments on the y axis. Gray bands show 95% confidence intervals in the regression. (B) Correlation between participant judgments and main model, binning participant judgments according to the alternate model's predictions. Each point represents one knowledge rating, with model predictions on the x axis and participant judgments on the y axis. Gray bands show 95% confidence intervals in the regression. This reveals meaningful variation our alternate model was not able to capture.

how much more difficult their chosen option was, in relation to the other. This choice was preregistered, with the idea that the cost our model assigned to each action (sailing vs. searching) would be scaled based upon participants' judgments. Finally, participants were asked what they thought the point of the task had been, and were given an opportunity to provide feedback or note any technical difficulties.

#### 4.6. Results

Participants rated how much the pirates knew, and how much they believed they could learn from the map, in 18 test trials. This yielded 36 final ratings. As in Experiment 1, participant responses were averaged by question, and then z-scored; the corresponding model predictions were also z-scored.<sup>6</sup>

Fig. 6 shows the overall results, revealing that our model was highly correlated with participant judgments,  $r = 0.86$  (95% CI: 81, 92.9). Similarly to Experiment 1, this correlation did not reflect only cases where both the model and participants inferred a lot of knowledge or very little knowledge. Critically, it included cases where both the model and participants were equally uncertain, in a graded manner, about how much the agent knew.<sup>7</sup> Fig. 7 plots the correspondence between model and participant ratings for each trial separately, further showing that participants' judgments were not bi-modal, and did depend on both the size of the island and the cost of obtaining the map.

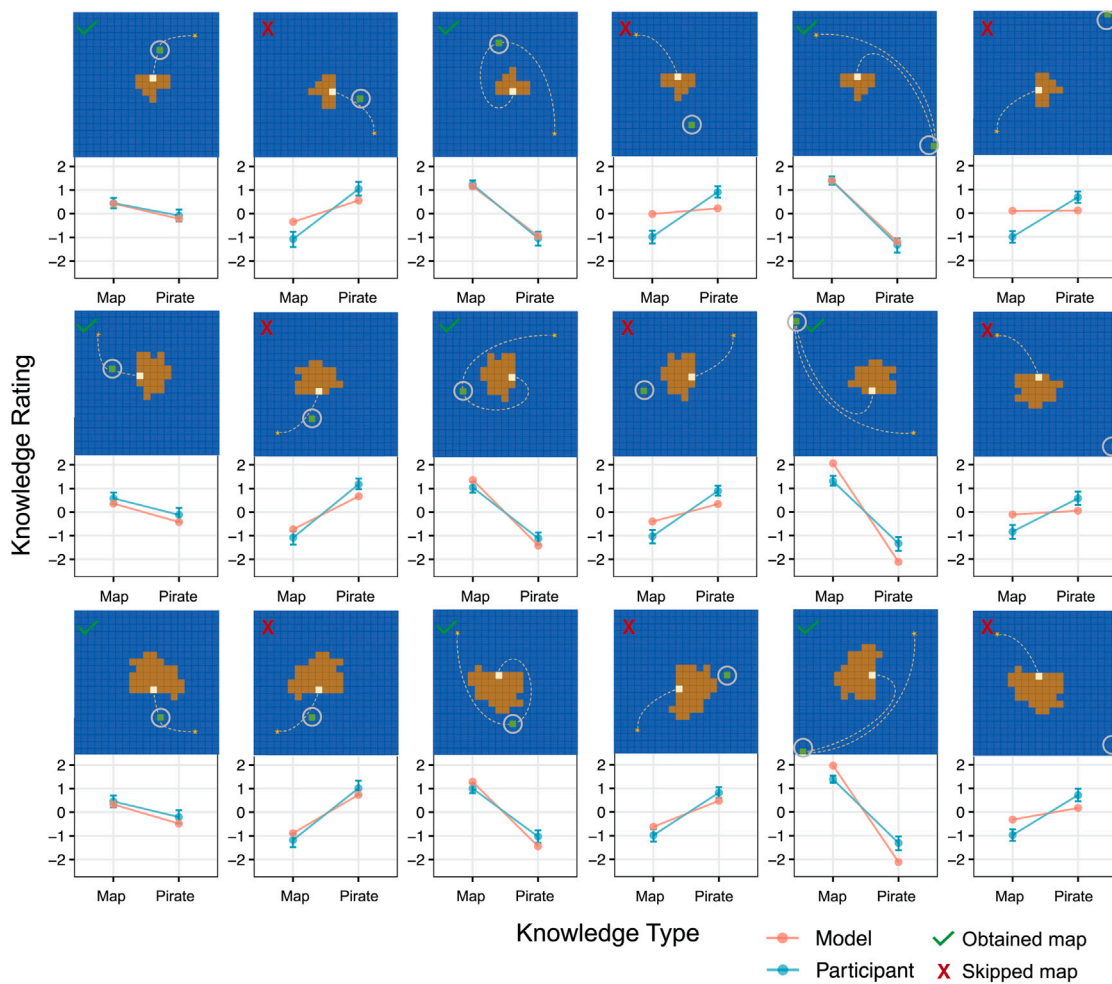
To check whether these results could be the product of a simple heuristic, we implemented an alternate model. Rather than performing full mental-state inference, our alternate model simply assumed that an agent who skipped the map both (a) knew a lot about the island and (b) thought the map contained little information (and vice versa if the agent traveled to the map). We implemented this via a linear regression that estimated how much participants considered to be “a lot” or “a little” knowledge in each case, and produced one of those four values (“a lot” and “a little” for each of “pirate knowledge” and “map information”) in every trial. Because this model was insensitive to cost, it did not consider more graded cases we expected humans might (e.g., that if the map is right on the way you might check even if you are not sure how much you will learn; whereas if the map is far away, you may choose not to obtain it even if you lack some knowledge). This alternate model showed a stronger correlation

with participant judgments,  $r = 0.94$  (95% CI: 91.4, 97.7); a bootstrap over the correlation difference revealed that the alternate model was reliably better correlated with participant judgments than the main model (correlation difference, alternate model – main model = 0.079, 95% CI: 0.9, 14.6; not preregistered). Combined with the relatively high correlation between our main model and alternate model predictions,  $r = .76$ , (95% CI: .58, .87), this suggests that following a simple heuristic will produce reasonably good answers much of the time. However, although the alternate model was better correlated with participant judgments (perhaps not unexpectedly, as it was trained on participant judgments in the first place), it produced “all or nothing” predictions that were insensitive to the size of the island and distance to the map, and therefore did not capture the graded range of intermediate values we found in participants' responses. While it is generally true that in our task, agents sought out information when they needed it and skipped it when they did not, both participants and our main model were able to make much more nuanced epistemic inferences, recognizing that some trials strongly imply a lot of knowledge, while others imply the possibility of some knowledge, etc. To better illustrate this, Fig. 8a shows a pair of trials from Experiment 2 where the alternate model makes identical predictions but participants (and our main model) produce systematically different inferences. In the trial on the left, the map is already along the pirate's route to the island, while obtaining the map in the rightmost trial requires a significant detour (and hence a much higher cost to the pirate). Participants (and our main model) were both sensitive to these features, and thus rated the map as having only slightly more information than average in the left hand trial, but significantly more information than average in the right hand trial. The alternate model, however, produced identical predictions in both of these trials.

To evaluate the alternate model's performance more rigorously, we conducted two follow-up analyses. First, following our preregistered analysis plan, we tested whether there is actually meaningful variation in participant judgments that the alternate model fails to capture (despite well-capturing the overall trajectory of participants' responses). Specifically, because the alternate model binned all predictions into four categories (“a lot” versus “a little” for each of agent knowledge and map information), we tested whether participant judgments *within* each of these categories were still well-correlated with those of our main model. If the heuristic model does in fact more accurately capture participant judgments, we would expect that in any two trials where it predicts the same “bin”, there should not be any meaningful variation in participant responses between those trials. On the other hand, if there is variation, and it correlates strongly with our main model, that would suggest structured variation in participant responses not captured by the alternate model. In other words, obtaining meaningful

<sup>6</sup> We mistakenly preregistered a slightly different z-scoring procedure—z-scoring participant ratings and then averaging by trial and prediction type. For consistency, we follow the process outlined in Experiment 1.

<sup>7</sup> This gradedness is shown more clearly in the individual-level scatterplots in the Experiment 2 section of Supplemental Materials.



**Fig. 7.** Detailed results for Experiment 2. Each panel presents one trial, with results split by the inference type (“information contained in map” on left, and “pirate’s knowledge of island” on right), indicated on the x-axis. The y axis indicates standardized knowledge ratings. Participant judgments are plotted in blue; model predictions are plotted in red. Vertical bars show 95% confidence intervals over participant judgments. Images above each panel show trial configuration, including island size, map location (green box, circled in gray), pirate’s path (dotted gray line), and whether the pirate detoured to obtain the map (green checkmark) or skipped the map (red X). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

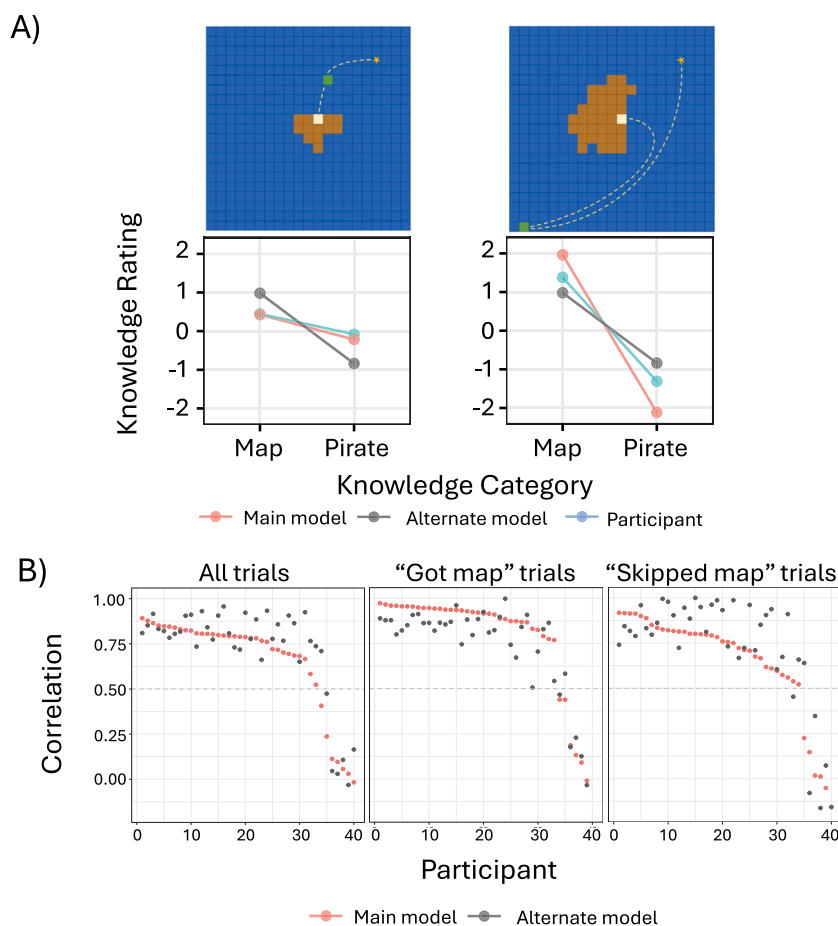
correlations within each bin suggests that there is still structure in each category that only our main model is able to capture. Consistent with this possibility, even when separating participant judgments according to the predictions of our alternate model, participants’ judgments were significantly correlated with the corresponding judgments from our main model (all  $r$ ’s between [0.67, 0.91], all  $p$ ’s < .05; see Fig. 6b). This demonstrates that our alternate model fails to capture meaningful variation in participant judgments, despite the high overall correlation between participant judgments and the predictions of our alternate model.

As an additional follow-up analysis (not pre-registered), we computed correlations between participant judgments and both model predictions for each participant individually. Across all trials, our main model had a higher correlation with only 42.5% ( $n = 17/40$ ) of participants (Fig. 8b, leftmost panel). However, this analysis also revealed a stark difference between trials: in trials where the pirate chose to obtain the map, our main model showed a better fit with 73% ( $n = 29/40$ ) of participants (Fig. 8b, center panel), compared to only 30% ( $n = 12/40$ ) of participants in trials where the pirate skips the map (Fig. 8b, rightmost panel). We also found that the main model had a significantly higher correlation with averaged participant responses than the alternate model in the “gets map” trials (Main:  $r = .98$ , 95%CI: .93, .99; Alternate:  $r = .88$ , 95%CI: .74, .96; correlation difference, alternate-main = -.10, 95%CI: -.14, -.01), while the opposite was true

in the “skips map” trials (Main:  $r = .86$ , 95%CI: .65, .95; Alternate:  $r = .99$ , 95%CI: .97, .99; correlation difference, alternate-main = .13, 95%CI: .01, .22). These results reveal that, although participants are clearly able to make graded, quantitatively precise knowledge inferences in a way that reflects our main hypothesis, they may revert to a simpler heuristic in cases where the agent chooses not seek out additional knowledge. We return to this possibility and its implications in the general discussion.

### 5. General discussion

Here we presented two experiments and a computational model designed to test people’s capacity to make amorphous epistemic inferences: quantitative estimates about how much someone knows or expects to learn, in contexts where it is not possible to infer the precise representations or contents of this knowledge. We found that people can make quantitative inferences about how much someone knows (Experiment 1), and joint inferences about how much someone knows and how much they expect to learn (Experiment 2), all from minimal observable choices. These inferences were predicted by a normative, Bayesian model that estimates amount of knowledge by considering how different epistemic states would affect an agent’s expected costs (and thus their behavior). Furthermore, these judgments reflected a fine-grained sensitivity to task features that could not be explained by



**Fig. 8.** Panel (A) Two example trials from Experiment 2, showing main model predictions (red), participant judgments (blue), and alternate model predictions (gray). Stimulus information for each trial is shown above corresponding panel. Panel (B) Individual-level correlations between participant judgments and model predictions, with participant number on the *x*-axis and model correlation on the *y*-axis. Main model shown in red and alternate model shown in gray. Leftmost plot shows correlations from all trials, middle plot isolates trials where the pirate got the map, and rightmost plot isolates trials where the pirate skipped the map. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

our alternate models, which encoded simple heuristics that ignored expected costs. In particular, while our alternate models were able to predict qualitative patterns in participant judgments (e.g.: that an agent knows more about Field A than Field B), they failed to capture participants' ability to infer different degrees of knowledge (e.g.: how much more the agent knows about Field A than Field B), instead making coarse, categorical inferences that the agent knew "a lot" or "a little". However, the results of Experiment 2 also suggested that participants may strategically revert to a simpler heuristic in contexts where precise knowledge estimates seem less important, a possibility we discuss further below.

Our computational model followed the same principles that shape related models of Theory of Mind, where mental-state inference is structured around an assumption that agents act to maximize utilities—the difference between the costs that agents incur and the rewards they obtain (Gergely & Csibra, 2003; Jara-Ettinger et al., 2016; Jern et al., 2017; Lucas et al., 2014). Our model builds on these ideas, and extends them by explicitly modeling the idea that, by observing the apparent costs agents incur, we can recover the amount of knowledge they possess. The quantitative fit between our model and participants suggests that the mechanisms supporting inferences about specific epistemic states follow the same principles as the mechanisms supporting broader inferences about the magnitude of knowledge.

Related work has developed computational models that explain how people infer each other's beliefs about the world (Baker et al., 2017). These inferences, however, depend on access to a highly constrained set

of epistemic hypotheses, and to observable behavior that is diagnostic of the agent's epistemic state. While these inferences are undoubtedly critical for social interaction, many everyday social behaviors lack the information needed to make such precise and targeted epistemic inferences. We show that, in such situations, people can nonetheless derive quantitative estimates of how much knowledge someone might possess (or believe they can come to possess). This capacity might be particularly important in informal pedagogy, as it might help us identify agents who are knowledgeable, who we could subsequently seek out to learn from. These inferences, given that they require fewer observations, might also serve as a powerful attention cue. Imagine, for instance, being a competitor in a setting like Experiment 1. Quickly detecting that an agent is knowledgeable might prompt us to attend to them carefully as they take additional actions, so that we can further uncover what specific knowledge they have.

Following a large tradition in computational cognitive science, our model was designed to explain human behavior at a computational level of analysis (Marr, 1982). Models at the computational level typically remain agnostic about the underlying algorithmic implementation in the human mind. In our case, however, we believe there are strong reasons to suspect our model is not a plausible candidate for an algorithmic implementation. This is because our model makes two critical assumptions: First, observers must have access to a range of epistemic hypotheses that they can evaluate; second, they must have a way to quantify the amount of knowledge contained within each epistemic hypothesis.

While the first assumption may seem plausible in some situations, there are many cases where we cannot represent the internal structure of epistemic hypotheses, or have access to the hypothesis space. For instance, while we know that pilots can fly planes, most of us do not know how to represent what a pilot knows (unlike in our experiments, where we knew how to represent different possible knowledge states the agents might have). This suggests that some amorphous inferences cannot be supported by an algorithm that requires people to integrate many specific hypotheses about an agent's knowledge. Similarly, the second assumption (that it is possible to quantify the amount of knowledge in each hypothesis) was easy to formalize in our experimental contexts. But this is not always the case. In the same example about pilots, even when we build specific representations of knowledge, such as "the pilot knows how turn the autopilot on and off", it is difficult to gauge the amount of knowledge involved without having the knowledge ourselves. For instance, if it is just a simple button press, little knowledge is needed. But if using autopilot requires managing a wide range of other parameters, then a lot of knowledge is needed.

The fact that our model is an unlikely candidate for an algorithmic implementation makes people's results, in some sense, even more interesting. Somehow, participants in our task were able to generate estimates of knowledge that quantitatively resembled our normative model. This suggests that people have access to some approximations that manage to produce inferences that approximate normative inferences. Thus, our results are best thought of as establishing that people have a capacity to make quantitative and graded amorphous knowledge inferences, and opens questions for future research about how exactly people accomplish this.

Our results also leave open the possibility that participants used some intermediate strategy that is less complex than our main model but more sophisticated than the simple "knows a lot or knows very little" heuristic encoded in our alternate model. One possibility is that, instead of averaging expected search costs across all possible epistemic states, which is typically intractable, people may estimate the search costs from a small number of states sampled from some internal distribution, then apply the usual expected utility maximization reasoning to these estimated search costs. This type of sampling-based approximation is common in computer science, and recent work suggests that it may also be an essential strategy for bounded cognitive agents (Lieder & Griffiths, 2020). This could explain why participants were clearly sensitive to search costs in most trials, but failed to recognize how a denser field is less costly to search in other trials, as we found in our Experiment 1 results.

The results of Experiment 2 also suggested that participants were using different reasoning strategies depending on whether the pirate took a detour to obtain more information. When the pirate got the map, participants very closely matched our main model predictions, showing a fine-grained sensitivity to the difficulty of searching the island and the cost incurred to obtain more information. When the pirate skipped the map, however, participant judgments were less sensitive to these features, and significantly more correlated with our alternate model predictions. One possibility is that participants selectively deploy different inference strategies depending on the nature of the task and the perceived importance of making precise estimates. For example, participants may see skipping the map (i.e.: not incurring an extra cost) as the default behavior, signaling that the pirate "knows enough" without requiring additional numerical precision. When the pirate obtains the map, however, this might trigger additional computation to determine the exact size of the knowledge deficit that would justify incurring the extra cost. While this may reflect two fundamentally different strategies for these two types of cases, it may also reflect that participants use the same approximation strategy described above (estimating expected search costs by sampling a few epistemic states), but devote fewer mental resources (i.e.: draw fewer samples) to estimating the map's value when the pirate skips it. Whether and when participants do in fact

strategically deploy different inference strategies remains a question for future work.

Our results also leave an empirical question open: although our focus was on amorphous knowledge inferences, we do not know if people also spontaneously attempted to make specific epistemic inferences. Although it is impossible to infer exactly what the agent knew, some context might reveal partial information. For instance, in Experiment 2, if a ship bypasses the island port and travels far away to collect a map, people might think that the pirates were confident that the treasure would not be close to the port. Furthermore, inferring specific epistemic states is not necessarily inconsistent with our main hypothesis: if, as suggested earlier, people are approximating the computations in our main model by sampling a handful of epistemic states (rather than exhaustively averaging over all of them), then we could think of this behavior (inferring a specific epistemic state) as a single-sample approximation of the main model. Prior research suggests that single-sample approximations can be rational in certain circumstances (Vul et al., 2014), which opens the possibility that people make approximate amorphous knowledge inferences by sampling individual epistemic states. This points to future research in which people's knowledge judgments are modeled as a hierarchical two-tiered inference where we use observable action to simultaneously make broad epistemic inferences and specific targeted inferences when possible.

Overall, our work sheds light on a common everyday epistemic inference: the ability to infer how much others know or believe they can learn, even when there is insufficient information to infer the exact contents of their knowledge. This work highlights a space of inferences that have been historically understudied in Theory of Mind, but that might be equally important. The capacity to build quick, high-level snapshots of what is in other minds might be one of the most important representations that direct our decisions over whom to attend to, seek information from, and trust.

#### CRediT authorship contribution statement

**Rosie Aboody:** Writing – original draft, Investigation, Visualization, Conceptualization, Writing – review & editing, Methodology, Funding acquisition, Data curation, Formal analysis. **Isaac Davis:** Writing – review & editing, Conceptualization, Writing – original draft, Methodology, Software, Formal analysis. **Yarrow Dunham:** Supervision, Writing – review & editing, Conceptualization. **Julian Jara-Ettinger:** Writing – review & editing, Project administration, Conceptualization, Writing – original draft, Resources, Funding acquisition, Supervision, Methodology.

#### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

#### Acknowledgments

We thank the members of the Yale Computational Social Cognition Lab for helpful conversations and advice, and we thank Emily Gerdin for her assistance with rubber duck debugging. This work was supported by National Science Foundation award BSC-2045778, and also by Yale's Franke Program in Science and the Humanities, via a Franke Interdisciplinary Graduate Award to RA.

#### Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.cognition.2025.106236>.

## Data availability

All stimuli, data, model, and analysis code is available at: <https://osf.io/pjy6x>.

## References

- 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), 1–10.
- Csibra, G. (2003). Teleological and referential understanding of action in infancy. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 358(1431), 447–458.
- Davis, I., Carlson, R., Dunham, Y., & Jara-Ettinger, J. (2023). Identifying social partners through indirect prosociality: a computational account. *Cognition*, 240, Article 105580.
- Gergely, G., & Csibra, G. (2003). Teleological reasoning in infancy: The naive theory of rational action. *Trends in Cognitive Sciences*, 7(7), 287–292.
- Jara-Ettinger, J., Gweon, H., Schulz, L. E., & Tenenbaum, J. B. (2016). The naïve utility calculus: Computational principles underlying commonsense psychology. *Trends in Cognitive Sciences*, 20(8), 589–604.
- Jara-Ettinger, J., Schulz, L. E., & Tenenbaum, J. B. (2020). The naive utility calculus as a unified, quantitative framework for action understanding. *Cognitive Psychology*.
- Jern, A., & Kemp, C. (2014). Reasoning about social choices and social relationships. *Proceedings of the 36th Annual Meeting of the Cognitive Science Society*.
- Jern, A., Lucas, C. G., & Kemp, C. (2017). People learn other people's preferences through inverse decision-making. *Cognition*, 168, 46–64.
- Koenig, M. A., Cole, C. A., Meyer, M., Ridge, K. E., Kushnir, T., & Gelman, S. A. (2015). Reasoning about knowledge: Children's evaluations of generality and verifiability. *Cognitive Psychology*, 83, 22–39.
- Koenig, M. A., & Harris, P. L. (2005). Preschoolers mistrust ignorant and inaccurate speakers. *Child Development*, 76(6), 1261–1277.
- Kominsky, J. F., Langthorne, P., & Keil, F. C. (2016). The better part of not knowing: Virtuous ignorance. *Developmental Psychology*, 52(1), 31.
- Landrum, A. R., & Mills, C. M. (2015). Developing expectations regarding the boundaries of expertise. *Cognition*.
- 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.
- Liu, S., Ullman, T. D., Tenenbaum, J. B., & Spelke, E. S. (2017). Ten-month-old infants infer the value of goals from the costs of actions. *Science*, 358(6366), 1038–1041.
- Lucas, C. G., Griffiths, T. L., Xu, F., Fawcett, C., Gopnik, A., Kushnir, T., Markson, L., & Hu, J. (2014). The child as econometrician: A rational model of preference understanding in children. *PLoS One*, 9(3), Article e92160.
- Lutz, D. J., & Keil, F. C. (2002). Early understanding of the division of cognitive labor. *Child Development*, 73(4), 1073–1084.
- Marr, D. (1982). *Vision: A computational investigation into the human representation and processing of visual information*. MIT Press.
- Ronfard, S., & Corriveau, K. H. (2016). Teaching and preschoolers' ability to infer knowledge from mistakes. *Journal of Experimental Child Psychology*, 150, 87–98.
- Ullman, T., Baker, C., Macindoe, O., Evans, O., Goodman, N., & Tenenbaum, J. (2009). Help or hinder: Bayesian models of social goal inference. *Advances in Neural Information Processing Systems*, 22.
- Vul, E., Goodman, N., Griffiths, T. L., & Tenenbaum, J. B. (2014). One and done? Optimal decisions from very few samples. *Cognitive Science*, 38(4), 599–637.
- Woodward, A. L. (1998). Infants selectively encode the goal object of an actor's reach. *Cognition*, 69(1), 1–34.