

Citation: Marti, L., Mollica, F., Piantadosi, S., & Kidd, C. (2018). Certainty is Primarily Determined by Past Performance During Concept Learning. *Open Mind: Discoveries in Cognitive Science*, 2(2), 47–60. [https://doi.org/10.1162/opmi\\_a\\_00017](https://doi.org/10.1162/opmi_a_00017)

DOI:  
[https://doi.org/10.1162/opmi\\_a\\_00017](https://doi.org/10.1162/opmi_a_00017)

Supplemental Materials:  
[www.mitpressjournals.org/doi/suppl/10.1162/opmi\\_a\\_00017](http://www.mitpressjournals.org/doi/suppl/10.1162/opmi_a_00017)

Received: 9 February 2017  
Accepted: 17 April 2018

Competing Interests: The authors declare that they have no competing interests.

Corresponding Author:  
Louis Marti  
[LMarti13@gmail.com](mailto:LMarti13@gmail.com)

Copyright: © 2018  
Massachusetts Institute of Technology  
Published under a Creative Commons  
Attribution 4.0 International  
(CC BY 4.0) license



The MIT Press

# Certainty Is Primarily Determined by Past Performance During Concept Learning

Louis Marti<sup>1,2</sup>, Francis Mollica<sup>1</sup>, Steven Piantadosi<sup>1,2</sup>, and Celeste Kidd<sup>1,2</sup>

<sup>1</sup>Brain and Cognitive Sciences, University of Rochester, Rochester

<sup>2</sup>Psychology, University of California, Berkeley

**Keywords:** certainty, confidence, metacognition, learning, concepts

## ABSTRACT

Prior research has yielded mixed findings on whether learners' certainty reflects veridical probabilities from observed evidence. We compared predictions from an idealized model of learning to humans' subjective reports of certainty during a Boolean concept-learning task in order to examine subjective certainty over the course of abstract, logical concept learning. Our analysis evaluated theoretically motivated potential predictors of certainty to determine how well each predicted participants' subjective reports of certainty. Regression analyses that controlled for individual differences demonstrated that despite learning curves tracking the ideal learning models, reported certainty was best explained by performance rather than measures derived from a learning model. In particular, participants' confidence was driven primarily by how well they observed themselves doing, not by idealized statistical inferences made from the data they observed.

## INTRODUCTION

Daily life requires making judgments about the world based on inconclusive evidence. These judgments are intrinsically coupled to people's subjective *certainty*, a metacognitive assessment of how accurate judgments are. While it is clear certainty impacts behavior, we do not fully understand how subjective certainty is linked to objective, veridical measures of certainty or probability. For example, people presented with disconfirming evidence can become even more entrenched in their original beliefs. Tormala, Clarkson, and Henderson (2011) and Tormala and Petty (2004) found that when people were confronted with messages that they perceived to be strong (e.g., from an expert) but contradicted their existing beliefs, their belief certainty *increased* instead of decreased. Similarly, the Dunning-Kruger effect—by which unskilled people overestimate their abilities and highly competent people underestimate them—also provides evidence of a miscalibration (Kruger & Dunning, 1999). Confidence is also influenced by social factors. Specifically, individuals calibrate their confidence to the opinions of others, irrespective of the accuracy of those opinions (Yaniv, Choshen-Hillel, & Milyavsky, 2009). Tsai, Klayman, and Hastie (2008) found that presenting individuals with more information raised their confidence irrespective of whether accuracy increased. Miscalibration is also present during “wisdom of the crowds” tasks. When questions require specialized information, individuals are equally as confident regardless of accuracy. This applies to both answers to questions and predictions about the accuracy of others (Prelec, Seung, & McCoy, 2017). Additionally, confidence in a memory has no relationship to whether or not the memory actually occurred (Loftus, Donders, Hoffman, & Schooler, 1989; McDermott & Roediger, 1998). Finally, simply taking prescription stimulants (e.g., Adderall, Ritalin) increases individuals' senses of certainty (Smith & Farah, 2011).

Studies examining perceptual phenomena, however, imply a tight link between certainty and reality. Individuals calculate their own subjective measure of visual uncertainty, which has been found to predict objective uncertainty (Barthelmé & Mamassian, 2009). Others have found correlates for subjective certainty such as reaction time, stimuli difficulty, and other properties of the data (Drugowitsch, Moreno-Bote, & Pouget, 2014; Kepecs, Uchida, Zariwala, & Mainen, 2008; Kiani, Corthell, & Shadlen, 2014). More evidence demonstrating the linkage between perceptual certainty and reality was presented when Sanders, Hangya, and Kepecs (2016) described a computational model that predicted certainty in auditory and numerical discrimination tasks.

Thus, while our certainty might be a useful guide with regard to perceptual decisions, such as trying to locate a friend yelling for help in the middle of the woods, it may be misleading in higher-level domains, such as deciding whether to see a chiropractor versus a medical doctor. However, no experiment has evaluated quantitatively measured changes in certainty during learning in tasks outside of perception. In ordinary life, evidence accumulation is likely to be less like perceptual learning and more like tasks for which learners must acquire abstract information about more complex latent variables—like rules, theories, or structures. Here, we examine certainty during learning using an abstract learning task with an infinite hypothesis space of logical rules. We present three experiments that used a Boolean concept-learning task to measure how certain learners should have been, given the strength of the observed evidence. With a potentially overwhelming hypothesis space, is a person's subjective certainty driven by veridical probabilities, or by something else?

Historically, Boolean concept-learning tasks have been used to study concept acquisition because they allowed researchers to examine the mechanisms of learning abstract rules while focusing on a manageable, simplified space of hypotheses (Bruner & Austin, 1986; Feldman, 2000; Goodman, Tenenbaum, Feldman, & Griffiths, 2008; Shepard, Hovland, & Jenkins, 1961). Experiment 1 compared measures from an idealized learning model to measures derived from participants' behavior to determine which best matched participants' ratings of certainty. Results suggest that the most important predictor of certainty is people's recent feedback/accuracy, not measures of, for example, entropy derived from the model. Furthermore, a logistic regression with the best predictors demonstrates that most of them provide unique contributions to certainty, implicating many factors in subjective judgments. Experiment 2 tested these predictors when participants were not given feedback. These results show that when feedback is removed, model predictors perform no better than in Experiment 1. Experiment 3 examined participants' certainty about individual trials rather than the overall concept. Similar to Experiment 1, in Experiment 3 people primarily relied on recently observed feedback. Our results show that participants used their overall and recent accuracy—not measured or derived from rule-learning models—to construct their own certainty.

## **EXPERIMENT 1**

### ***Motivation***

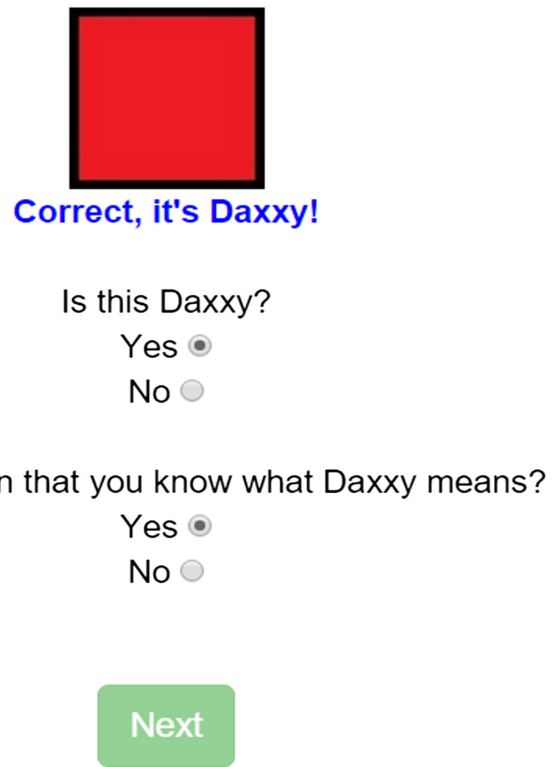
The aim of Experiment 1 was to measure subjective certainty of participants during concept learning and attempt to predict it using plausible model-based and behavioral predictors. In this experiment, certainty judgments were about what underlying concept (rule) generated the data they saw, as opposed to their certainty about the correct answer for any given trial (see Experiment 3).

### Methods

We tested 552 participants recruited via Amazon Mechanical Turk in a standard Boolean concept-learning task during which we measured their knowledge of a hidden concept (via *yes* or *no* responses) and their certainty throughout the learning process (see Figure 1 and Table 1). In this experiment, participants were shown positive and negative examples of a target concept “daxxy,” where membership was determined by a latent rule on a small set of feature dimensions (e.g., color, shape, size), following experimental work by Shepard et al. (1961) and Feldman (2000). The latent rules participants were required to learn varied across a variety of logical forms. After responding to each item, participants were provided feedback and then rated their certainty on what the word “daxxy” meant. For our analyses we considered and compared several different models of what might drive uncertainty (see Table 2). These predictors can be classified into two broad categories. Model-based predictors were calculated using our ideal learning model, while behavioral predictors were calculated using the behavioral data (see Appendix A in the Supplemental Materials [Martí, Mollica, Piantadosi, & Kidd, 2018] for additional method details).

### Results

We first visualize plots of participants’ certainty and accuracy for each concept in order to show (a) whether certainty and accuracy improved over the course of the experiment, (b) whether theoretically harder concepts (according to Feldman, 2000) were, in fact, more difficult for participants, and (c) whether participants’ certainty correlated with their accuracy in general.



**Figure 1.** In Experiment 1, participants saw 24 trials (as above), randomized between conditions. Feedback was displayed after responding.

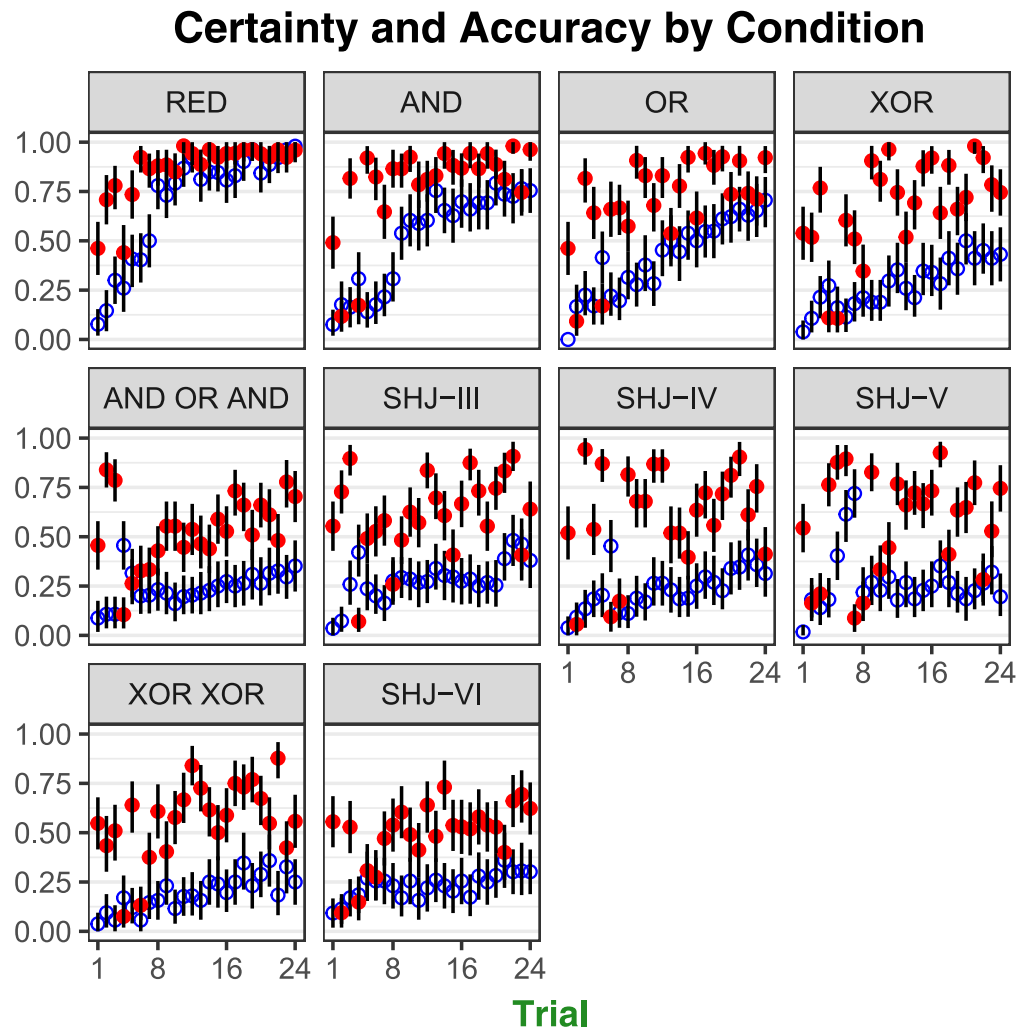
**Table 1.** Concepts presented to participants. Concepts 1 and 5–9 are the Shepard, Hovland, and Jenkins family consisting of three features and four positive examples.

Concept		
1	SHJ-I <sub>3[4]</sub>	red
2	AND	red $\wedge$ small
3	OR	red $\vee$ small
4	XOR	red $\oplus$ small
5	SHJ-II <sub>3[4]</sub>	(red $\wedge$ small) $\vee$ (green $\wedge$ large)
6	SHJ-III <sub>3[4]</sub>	(green $\wedge$ large $\wedge$ triangle) $\vee$ (green $\wedge$ large $\wedge$ square) $\vee$ (green $\wedge$ small $\wedge$ triangle) $\vee$ (red $\wedge$ large $\wedge$ square)
7	SHJ-IV <sub>3[4]</sub>	(green $\wedge$ large $\wedge$ triangle) $\vee$ (green $\wedge$ large $\wedge$ square) $\vee$ (green $\wedge$ small $\wedge$ triangle) $\vee$ (red $\wedge$ large $\wedge$ triangle)
8	SHJ-V <sub>3[4]</sub>	(green $\wedge$ large $\wedge$ triangle) $\vee$ (green $\wedge$ large $\wedge$ square) $\vee$ (green $\wedge$ small $\wedge$ triangle) $\vee$ (red $\wedge$ small $\wedge$ square)
9	SHJ-VI <sub>3[4]</sub>	(green $\wedge$ large $\wedge$ triangle) $\vee$ (green $\wedge$ small $\wedge$ square) $\vee$ (red $\wedge$ large $\wedge$ square) $\vee$ (red $\wedge$ small $\wedge$ triangle)
10	XOR XOR	red $\oplus$ small $\oplus$ square

Figure 2 shows participants’ certainty and accuracy (*y*-axis) over trials of the experiment (*x*-axis). The accuracy curves indicate participants learned the concepts in some conditions but not others. This is beneficial to our analysis as it allows us to analyze conditions and trials in which participants should have had high uncertainty. Overall, participant certainty was inversely proportional to concept difficulty. Participant certainty generally increased, but only reached high values in conditions in which they also achieved high accuracy. The increasing trend of certainty in conditions for which accuracy did not go above 50% may be reflective of overconfidence. It is also important to note that even though participants received exhaustive evidence, there were still multiple logical rules that were both equivalent and correct. Despite this, participants still became certain over time.

**Table 2.** Certainty predictors (behavioral predictors in gray).

Predictor	Description
Trial	Number of trials seen so far
Total Accuracy	Total performance thus far
Local Accuracy	Performance on previous <i>N</i> trials ( <i>N</i> = 2, 3, 4, 5)
Local Accuracy Current	Performance on previous <i>N</i> trials ( <i>N</i> = 2, 3, 4, 5) and a guess on the current trial
Current Accuracy	Performance on the current trial
Entropy	Model uncertainty over hypotheses regarding what the concept is
Domain Entropy	Model uncertainty over which objects belong to the concept
Change in Entropy	Entropy change from the previous trial
Change in Domain Entropy	Domain entropy change from the previous trial
Cross Entropy	How much beliefs about hypotheses have changed since the previous trial
Domain Cross Entropy	How much beliefs about which objects belong to the concept have changed since the previous trial
MAP	The probability of the best hypothesis
Maximum Likelihood	The probability of the best hypothesis ignoring the prior probability
Response Probability	The probability of the participant’s response given the model predictions



**Figure 2.** Mean certainty (hollow blue circles) and mean accuracy (filled red circles) across concepts for Experiment 1. Chance is 50% across all conditions if guesses are made randomly.

We will first consider our predictors as separate models in order to determine which best predict certainty. Subsequently we will build a model using the best predictors of each type in order to determine the unique contributions of each predictor.

We assessed our predictors with generalized logistic mixed-effect models fit by maximum likelihood with random subject and condition effects.<sup>1</sup> First, this analysis shows model accuracy significantly predicts behavioral accuracy ( $R^2 = .50$ ,  $\beta = .748$ ,  $z = 30.423$ ,  $p < .001$ ; Figure 3), meaning that overall performance can be reasonably well predicted by the learning model.

Figure 4 then shows mean certainty responses for each trial and condition ( $y$ -axis) over several different key predictors of certainty ( $x$ -axis). A perfect model here would have data points lying along the line  $y = x$  with a high  $R^2$  and very little residual variance. **Local**

<sup>1</sup> We also analyzed our data on an individual level in order to ensure our findings were not due to averaging effects (Estes & Todd Maddox, 2005). See Table A.1 in the Supplemental Materials (Martí et al., 2018).

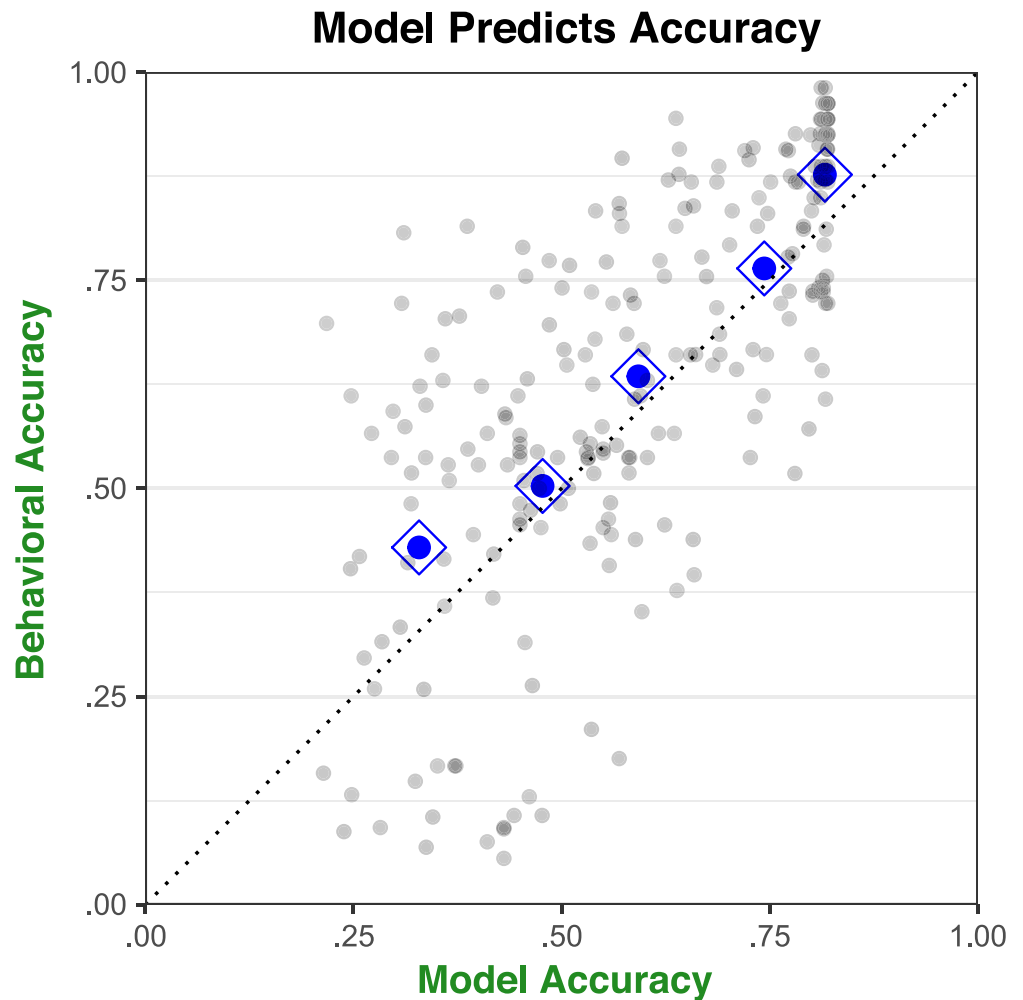


Figure 3. Model vs. behavioral accuracy for Experiment 1.

**Accuracy 5 Back**, the accuracy averaged over the past 5 items, has a high  $R^2$ , meaning that individuals with low local accuracy were uncertain and individuals with high local accuracy were highly certain. Likewise, **Domain Entropy** also has a high  $R^2$  and is very ordered compared to the other model predictors (see Figure A.1 in the Supplemental Materials [Martí et al., 2018] for additional predictor visualizations).

Table A.2 in the Supplemental Materials (Martí et al., 2018) shows the full model results, giving the performance of each model in predicting certainty ratings.<sup>2</sup> These have been sorted by Akaike information criterion (AIC), which quantifies the fit of each model penalizing its number of free parameters (closer to  $-\infty$  is better). The AIC score is derived from a generalized logistic mixed effect model fit by maximum likelihood with random subject and condition effects. This table also provides an  $R^2$  measure, calculated using the Pearson correlation between the means of each response and predictor for each trial and condition (this ignores variance from participants). As this table makes clear, the behavioral predictors tend to outperform the model predictors, at times by a substantial amount. The best predictor, **Local Accuracy 5 Back** accounts for 58% of the variance. Additionally, **Local Accuracy** models

<sup>2</sup> See Table A.3 in the Supplemental Materials (Martí et al., 2018) for simplified grammar predictors.

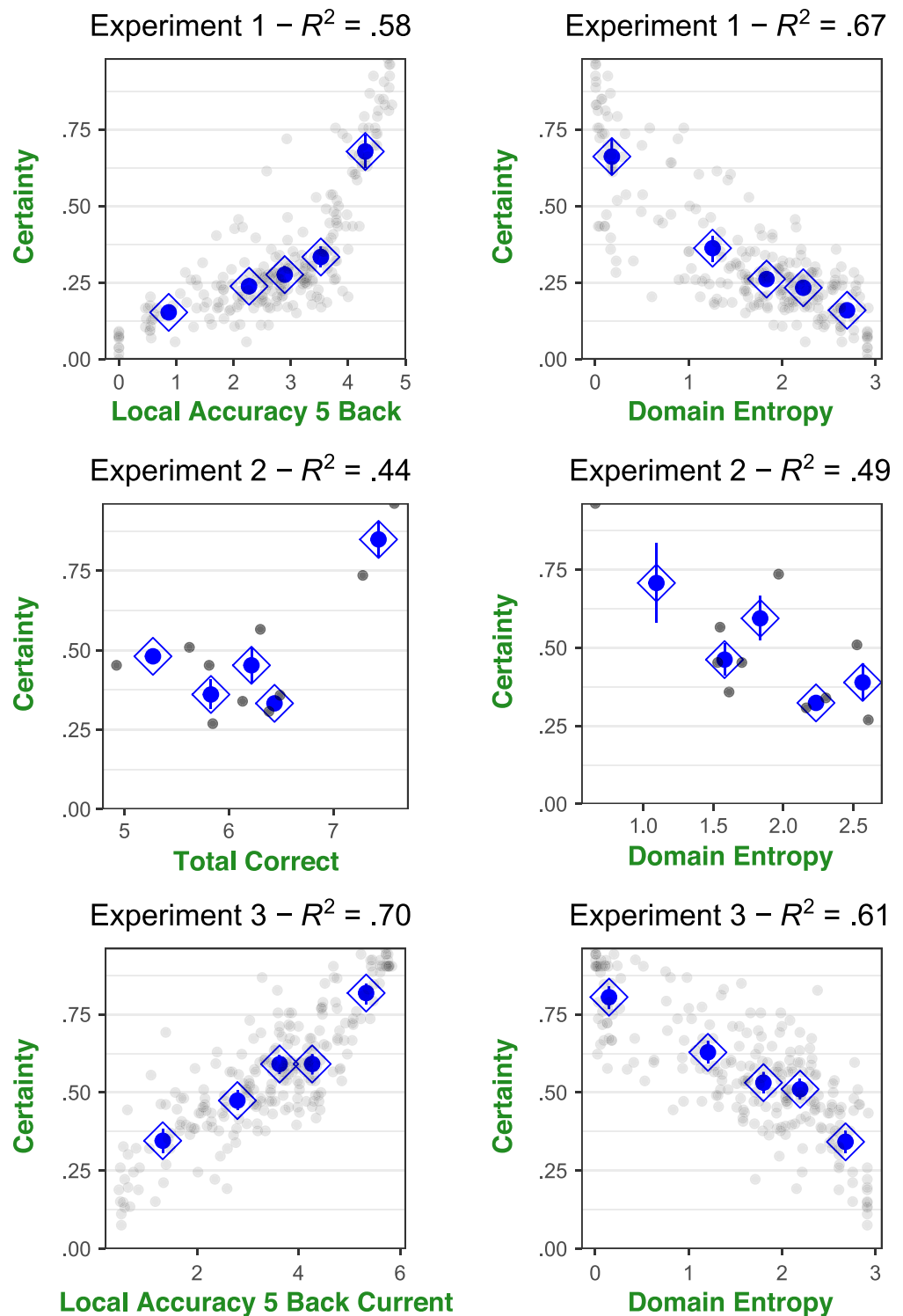


Figure 4. Key model fits for Experiments 1-3, showing mean participant responses for each concept and trial (gray) and binned model means in each of five quantiles (blue) for certainty rating ( $y$ -axis) as a function of model ( $x$ -axis). Diagonal lines with low variance correspond to models which accurately capture human behavior.



**Table 3.** Regression for best predictors (standardized) in Experiment 1 (behavioral predictors in gray).

Predictor	Beta	Standard Error	z Value	p
Intercept	−0.82	0.02	−37.61	< .001
Local Accuracy 5 Back	0.69	0.04	19.82	< .001
Log Trial	−0.60	0.04	−13.93	< .001
Total Correct	0.54	0.04	12.00	< .001
Domain Entropy	−0.34	0.06	−5.91	< .001
Entropy	−0.10	0.05	−1.93	.054
Log Maximum Likelihood	−0.04	0.04	−1.11	.269

outperform most of the other alternatives, a pattern that is robust to the way in which local accuracy is quantified (e.g., the number back that were counted or whether the current trial is included). The quantitatively best **Local Accuracy** model tracks accuracy over the past five trials. One possible explanation for this is that participants were simply basing their certainty on recent performance. The high performance of both **Local Accuracy** and **Total Correct** implies that people’s certainty is largely influenced by their own perception of how well they were doing on the task.

Strikingly, the lackluster performance of the majority of ideal learner models suggests that subjective certainty is not calibrated to the ideal learner. This is consistent with the theory that learners were likely not maintaining more than one hypothesis—perhaps they stored a sample from the posterior, but did not have access to the full posterior distribution. Strikingly, the idealized model of entropy over hypotheses—what might have corresponded to our best a priori guess for what certainty should reflect—performs especially poorly, worse than many behavioral and other model-based predictors. Such a failure of metacognition is consistent with the poor performance of **Current Accuracy**, a measure of whether or not the participant got the current trial correct. Subjective certainty does not accurately predict accuracy on the current trial, or vice versa.

Our first analysis treated each predictor separately and found the best, but what if multiple predictors were jointly allowed to predict certainty? To answer this, we created a model using the top three behavioral predictors and the top three model predictors in order to determine the unique contributions of each (see Table 3).<sup>3,4</sup> As the table makes clear, all behavioral predictors, along with **Domain Entropy**, make significant, unique contributions to certainty. Conversely, **Entropy** and **Log Maximum Likelihood** were not significant when controlling for the other predictors, demonstrating they provide no unique contributions to certainty. In alignment with the results of our AIC analysis, the (normalized) beta weights, which quantify the strength of each predictors’ influence, reveal that the behavioral predictors have the largest influence.

**Discussion**

Our results showed that an ideal learning model predicts learners’ accuracy in our task. These results hold regardless of whether certainty is measured on a binary, or a continuous scale (see Experiment 4 in Supplemental Materials Appendix D [Martí et al., 2018]). A plausible

<sup>3</sup> This regression was moderately sensitive to which predictors were included, likely due to some degree of multicollinearity.

<sup>4</sup> It was not possible to use random slopes (Barr, Levy, Scheepers, & Tily, 2013) in this regression due to a lack of convergence.



hypothesis would then be that the predictors derived from our ideal learning model would also be related to learners' certainty, perhaps to a large degree. Instead, we found that **Local Accuracy** and **Total Correct** are most predictive of people's certainty, outperforming our other predictors by predicting as much as 58% of the possible variance. In fact, overwhelmingly, the behavioral predictors performed better than the model predictors. **Domain Entropy** performs well and even has the highest  $R^2$  value, however it is important to emphasize that the  $R^2$  values did not take into account the subject and condition used in the mixed effect model. When these effects are controlled, we find that Domain Entropy has less of an influence than behavioral predictors, although its contribution to certainty is still nonzero. Performance of the predictors in a model that controls these effects should be a more reliable guide to each predictor's effect. Overall, the results suggested that participants primarily used the feedback on each trial in order to guide their senses of uncertainty about the concept.

## EXPERIMENT 2

### Motivation

Experiment 1 leaves open the possibility that both **Local Accuracy** and model-based predictors influence behavior, but that feedback overshadowed other predictors, perhaps because feedback was a quick and reliable cue. Experiment 2 tested this by removing feedback and thus removing it as a cue. We accomplished this by providing participants with only a *single* trial.

The critical question is whether the model-based predictors will become *more* predictive of responses compared to Experiment 1. If so, the cues to certainty may be strategically chosen based on what is informative, with participants able to use model-based measures when information about performance is absent. Alternatively, if the model-based predictors do not improve relative to Experiment 1, that would suggest that factors like **Local Accuracy** may be *the* driving force in metacognitive certainty and absent these predictors, people do not fall back on other systems.

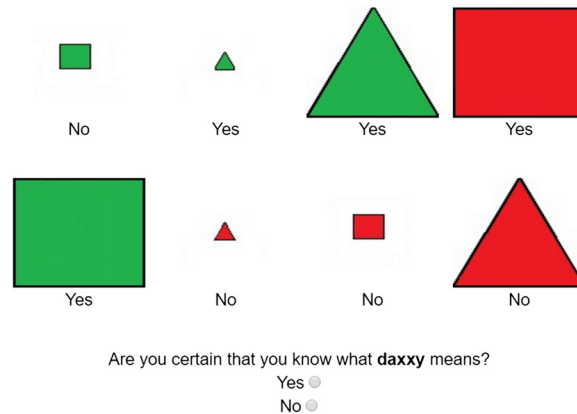
### Methods

Like Experiment 1, Experiment 2 presented participants with the task of discovering a hidden Boolean rule (see Figure 5 and Figure 6). We tested 577 participants via Amazon Mechanical Turk on a single-trial version of the same task used in Experiment 1, using the same set of concepts. The experimental trial tested participants on a single concept and displayed all eight images seen in a block of Experiment 1 simultaneously, each labeled with a *yes* or *no* to indicate whether it was part of the concept (see Figure 5). The participant answered whether they were certain what the concept was. They then saw the same set of eight images (randomized by condition) and were asked to label each as being a part of the concept (see Figure 6). (See Supplemental Materials Appendix B [Martí et al., 2018] for further details.)

### Results

Unlike Experiment 1, accuracy was high across most conditions, with average accuracy ranging from 62% to 95% across conditions (see Figure B.1 in the Supplemental Materials [Martí et al., 2018] for details). This was likely due to participants viewing the data simultaneously and testing them immediately afterward. Such a format would make it much easier to determine the concept and lead to reduced memory demands compared to Experiment 1. Despite this, subjective certainty was similar to Experiment 1 in that it related inversely to concept difficulty. Thus, since information regarding the underlying concept was still encoded and used in calculating their certainty, task differences did not seem to influence their certainty.

Look at each image below and try to figure out what makes something **daxxy** (yes) or not (no)



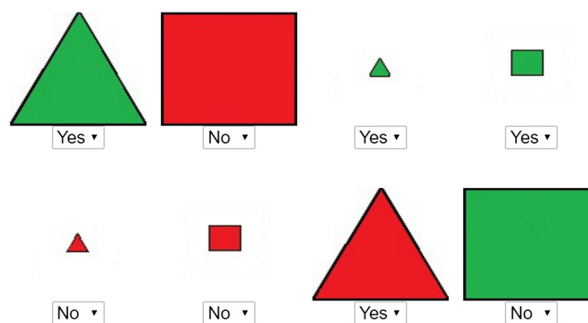
**Figure 5.** In Experiment 2, participants saw a single trial (as above), randomized between conditions.

For Experiment 2, we assessed our predictors with generalized logistic mixed-effect models fit by maximum likelihood with random condition effects. Unlike Experiment 1, the model fit for accuracy in Experiment 2 is not significant ( $R^2 = .02$ ,  $\beta = -.049$ ,  $z = -1.114$ ,  $p = .265$ ; see Figure B.2 in the Supplemental Materials [Martí et al., 2018]). This is likely due to data sparsity, although it is possible that participants did not learn these concepts as well due to the presentation format. In evaluating predictors of certainty Figure B.3 and Table B.1 in the Supplemental Materials (Martí et al., 2018) make clear that the results are similar to Experiment 1, with the best-performing predictors being behavioral measures. In this case, the only behavioral predictor, **Total Correct** is also the best predictor of certainty. Likewise, while **Domain Entropy** is the best performing model predictor, it is not as good as **Total Correct**. This is strong evidence that removing feedback had little to no effect on participants' propensity to avoid model-based predictors when constructing their own subjective certainty.

### Discussion

Our results demonstrate that feedback is not overriding model-based predictors when participants evaluate subjective certainty. When feedback is removed, participants still primarily used

Now we'd like to see which of these you think is **daxxy**. Select yes/no for each



**Figure 6.** In Experiment 2, after responding regarding their certainty, participants labeled each stimulus to assess their accuracy.

a behavioral predictor of overall accuracy in evaluating their own certainty. This could plausibly be because behavioral predictors provide a low-cost and rapid way of calculating certainty while model-based predictors are nonobvious and require more complex calculations.

### EXPERIMENT 3

#### *Motivation*

Both Experiment 1 and Experiment 2 asked about participants' certainty about a target *concept* that was underlying all of the observed data ("Are you certain you know what Daxxy means?"). However, word meanings are highly context dependent. A participant may be highly certain they know the meaning of "daxxy" within the confines of the experiment, but highly uncertain in general. Additionally, other work on metacognition has examined participants' certainty about their current *response*, where model-based effects can sometimes be seen. Experiment 3 examined trial-based certainty measures using the same setup of logical rules used in Experiments 1 and 2. If we find behavioral predictors no longer predict certainty but model-based predictors do, this would provide strong evidence that trial-certainty and concept-certainty are informed by two distinct processes.

#### *Methods*

Experiment 3 was a variant of Experiment 1 in which instead of asking "Are you certain that you know what Daxxy means?" we asked "Are you certain you're right?" after each response. We tested 536 participants on Amazon Mechanical Turk, using otherwise identical methods to Experiment 1 (see Supplemental Materials Appendix C [Martí et al., 2018] for further details).

#### *Results*

Unsurprisingly, participant accuracies were similar to Experiment 1, replicating the general observed trends (see Figure C.1 in the Supplemental Materials [Martí et al., 2018] for details). Importantly however, certainty in Experiment 3 seems to much more closely track accuracy on each trial, meaning that it is likely veridically reflecting participants' knowledge of each item response (as opposed to the meaning of "daxxy"). We assessed our predictors with generalized logistic mixed-effect models fit by maximum likelihood with random subject and condition effects. Like Experiment 1, the model fit between behavioral and model accuracy in Experiment 3 is reliable ( $R^2 = .50$ ,  $\beta = .808$ ,  $z = 31.529$ ,  $p < .001$ ; see Figure C.2 in the Supplemental Materials [Martí et al., 2018]).

Behavioral predictors once again overwhelmingly outperform the model-based predictors. Similar to Experiment 1, **Local Accuracy 5 Back Current** is the best predictor at 70% of variance explained, and the best model-based predictor is again **Domain Entropy**, which accounts for 61% of the variance (for details, see Figure C.3 and Table C.1 in the Supplemental Materials [Martí et al., 2018]).

#### *Discussion*

Experiment 3 provides strong evidence that participants primarily relied on local accuracy for their trial-based certainty just as they did for concept-based certainty. This reflects the fact that trial-based certainty, while more independent than concept-based certainty per trial, was

still influenced by performance and feedback on previous trials. Like Experiment 1, participants did not seem to be using most model-based predictors in their certainty calculations, despite behaving in line with model predictions with regard to accuracy. These results are seemingly in conflict with the Sanders et al. (2016) model, which they demonstrated to be a good predictor of participant certainty. One possibility is that these differences were the result of cross-trial learning in our task required. Neither Sanders et al. (2016) tasks required such cross-trial learning.

## GENERAL DISCUSSION

In conjunction with past research, our results paint a picture of how subjective certainty is derived for high-level logical domains like Boolean concept learning. It appears that certainty estimation primarily makes use of behavioral and overt task features, but that some model predictors are also relevant. In contrast, perceptual certainty and certainty involving one's memory of a fact (such as asking which country has a higher population; Sanders et al., 2016) seem to default to using predictors derived from ideal learning models.

In Experiments 1 and 3, **Local Accuracy** and **Total Correct** were very successful predictors of certainty. This means that participants seemed to primarily be basing their certainty on their past performance—*inferring certainty from their own behavior and feedback*. One view is that certainty's function is as a guide to inform our beliefs and decisions. If certainty was fulfilling this function, one might expect **Current Accuracy** to be an excellent predictor. Instead, we find it is an extremely poor predictor, implying that people's sense of certainty in these tasks is not likely to be a useful or important cause of behavior and is not calibrated well to their future performance. This is also in line with past research showing that some people's certainty is not based solely on their perceived probability of being correct, but also on the inverse variance of the data (Navajas et al., 2017). This general pattern is not unlike findings from metacognitive studies showing that often people do not understand—or perhaps even remember—the causes of their own behavior (Johansson, Hall, Sikström, & Olsson, 2005; Nisbett & Wilson, 1977). People do not directly observe their own cognitive processes and are often blind to their internal dynamics. This appears to be true in the case of subjective certainty reports when feedback is present and learning is taking place. In these cases, people do not appear to reflect an awareness of how much certainty they *should* have.

Past studies in memory have found that initial eyewitness confidence reliably predicted eyewitness accuracy, however, confidence judgments after memory “contamination” has occurred were no longer reliable (Wixted, Mickes, Clark, Gronlund, & Roediger, 2015). Given our results, a possible explanation for this is that the feedback in our experiments played the same role as the memory contamination in the eyewitness studies. In other words, recent feedback heavily influences certainty, and if that feedback is unreliable, it could lead to false memories.

It should be noted that one possible reason the behavioral predictors outperform the model predictors is that the behavioral predictors will vary with participants' mental states and thus with the natural idiosyncrasies within, although this effect may be mitigated by our use of mixed-effect models. For example, individual differences in attention that influence performance at the subject level could be captured by the behavioral predictors, but not the model-based predictors, which are functions only of the observed data. Though difficult to quantitatively evaluate, this difference may in part explain why the behavioral predictors are dominant in capturing performance, and this possible mechanism is consistent with the idea

that certainty is primarily derived from observing our own behavior and secondarily by the properties of the data.

Our analyses also help inform us about which factors *do not* drive certainty during learning, and several are surprising. One reasonable theory posits that participants could base their certainty off of their confidence in the Maximum a Posteriori (MAP) hypothesis under consideration. Since the MAP predictors do not perform well, it is unlikely that learners' certainty relies on internal estimates of the probabilities of the most likely hypothesis.

## CONCLUSION

Our findings suggest that although several types of predictors make unique contributions to certainty, the primary predictors of certainty are from observations of people's own behavior and performance, not from measures derived from an idealized learning model. Although learning patterns follow an idealized mathematical model, subjective certainty is only secondarily influenced by that model regardless of whether or not participants were able to observe how well they were doing. This is likely due to the underlying process of hypothesis formation and revision, as well as the way in which probabilities are handled beyond that which an ideal learner provides. These results also provide counterintuitive insight into why humans become certain. Certainty about a latent, abstract concept does not seem to be determined by the same mechanisms that drive learning. Instead, a large component of certainty could reflect factors that are largely removed from the veridical probabilities that any given hypothesis is correct.

## ACKNOWLEDGMENTS

We thank the Jacobs Foundation, the Google Faculty Research Awards Program, and the National Science Foundation Research Traineeship Program (Grant 1449828) for the funding to complete this work. We also thank members of the Kidd Lab and the Computation and Language Lab for providing valuable feedback.

## FUNDING INFORMATION

Celeste Kidd, Jacobs Foundation (DE); Celeste Kidd, Google (<http://dx.doi.org/10.13039/100006785>); Louis Martí, National Science Foundation (<http://dx.doi.org/10.13039/100000001>), Award ID: 1449828.

## AUTHOR CONTRIBUTIONS

LM: Conceptualization: Lead; Data curation: Lead; Formal analysis: Lead; Investigation: Lead; Methodology: Lead; Project administration: Lead; Software: Lead; Validation: Lead; Visualization: Lead; Writing—original draft: Lead; Writing—review & editing: Lead. FM: Formal analysis: Supporting; Methodology: Supporting; Writing—review & editing: Supporting. SP: Conceptualization: Supporting; Formal analysis: Supporting; Methodology: Supporting; Resources: Supporting; Supervision: Equal; Writing—review & editing: Supporting; Validation: Supporting. CK: Formal analysis: Supporting; Funding acquisition: Lead; Methodology: Supporting; Resources: Conceptualization: Supporting; Lead; Supervision: Equal; Writing—review & editing: Supporting; Validation: Supporting.

## REFERENCES

Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, *68*, 255–278.

Barthelmé, S., & Mamassian, P. (2009). Evaluation of objective uncertainty in the visual system. *PLOS Computational Biology*, *5*(9), e1000504.

- Bruner, J. S., & Austin, G. A. (1986). *A study of thinking*. Piscataway, NJ: Transaction.
- Drugowitsch, J., Moreno-Bote, R., & Pouget, A. (2014). Relation between belief and performance in perceptual decision making. *PLOS ONE*, *9*(5), e96511.
- Estes, W. K., & Todd Maddox, W. (2005). Risks of drawing inferences about cognitive processes from model fits to individual versus average performance. *Psychonomic Bulletin & Review*, *12*, 403–408.
- Feldman, J. (2000). Minimization of Boolean complexity in human concept learning. *Nature*, *407*, 630–633.
- Goodman, N. D., Tenenbaum, J. B., Feldman, J., & Griffiths, T. L. (2008). A rational analysis of rule-based concept learning. *Cognitive Science*, *32*, 108–154.
- Johansson, P., Hall, L., Sikström, S., & Olsson, A. (2005). Failure to detect mismatches between intention and outcome in a simple decision task. *Science*, *310*, 116–119.
- Kepecs, A., Uchida, N., Zariwala, H. A., & Mainen, Z. F. (2008). Neural correlates, computation and behavioural impact of decision confidence. *Nature*, *455*, 227–231.
- Kiani, R., Corthell, L., & Shadlen, M. N. (2014). Choice certainty is informed by both evidence and decision time. *Neuron*, *84*, 1329–1342.
- 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*, 1121–1134.
- Loftus, E. F., Donders, K., Hoffman, H. G., & Schooler, J. W. (1989). Creating new memories that are quickly accessed and confidently held. *Memory & Cognition*, *17*, 607–616.
- Martí, L., Mollica, F., Piantadosi, S., & Kidd, C. (2018). Supplemental materials for "Certainty is primarily determined by past performance during concept learning." *Open Mind: Discoveries in Cognitive Science*, *2*(2), 47–60. doi:10.1162/opmi\_a\_000017
- McDermott, K. B., & Roediger, H. L. (1998). Attempting to avoid illusory memories: Robust false recognition of associates persists under conditions of explicit warnings and immediate testing. *Journal of Memory and Language*, *39*, 508–520.
- Navajas, J., Hindocha, C., Foda, H., Keramati, M., Latham, P. E., & Bahrami, B. (2017). The idiosyncratic nature of confidence. *Nature Human Behaviour*, *1*(11), 810–818.
- Nisbett, R. E., & Wilson, T. D. (1977). Telling more than we can know: Verbal reports on mental processes. *Psychological Review*, *84*, 231–259.
- Prelec, D., Seung, H. S., & McCoy, J. (2017). A solution to the single-question crowd wisdom problem. *Nature*, *541*, 532–535.
- Sanders, J. I., Hangya, B., & Kepecs, A. (2016). Signatures of a statistical computation in the human sense of confidence. *Neuron*, *90*, 499–506.
- Shepard, R. N., Hovland, C. I., & Jenkins, H. M. (1961). Learning and memorization of classifications. *Psychological Monographs: General and Applied*, *75*(13), 1–42.
- Smith, M. E., & Farah, M. J. (2011). Are prescription stimulants "smart pills"? The epidemiology and cognitive neuroscience of prescription stimulant use by normal healthy individuals. *Psychological Bulletin*, *137*, 717–741.
- Tormala, Z. L., Clarkson, J. J., & Henderson, M. D. (2011). Does fast or slow evaluation foster greater certainty? *Personality and Social Psychology Bulletin*, *37*, 422–434.
- Tormala, Z. L., & Petty, R. E. (2004). Source credibility and attitude certainty: A metacognitive analysis of resistance to persuasion. *Journal of Consumer Psychology*, *14*, 427–442.
- Tsai, C. I., Klayman, J., & Hastie, R. (2008). Effects of amount of information on judgment accuracy and confidence. *Organizational Behavior and Human Decision Processes*, *107*(2), 97–105.
- Wixted, J. T., Mickes, L., Clark, S. E., Gronlund, S. D., & Roediger, H. L., III. (2015). Initial eyewitness confidence reliably predicts eyewitness identification accuracy. *American Psychologist*, *70*, 515–526.
- Yaniv, I., Choshen-Hillel, S., & Milyavsky, M. (2009). Spurious consensus and opinion revision: Why might people be more confident in their less accurate judgments? *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *35*, 558–563.