

The role of decision confidence in advice-taking and trust formation

Niccolò Pescetelli^{1,2} and Nick Yeung¹

¹Department of Experimental Psychology, University of Oxford, UK

²Max Planck Institute for Human Development, Berlin, Germany

July 9, 2020

Abstract

In a world where ideas flow freely across multiple platforms, people must often rely on others’ advice and opinions without an objective standard to judge whether this information is accurate. The present study explores the hypothesis that an individual’s internal decision confidence can be used as a signal to learn the accuracy of others’ advice, even in the absence of feedback. According to this “agreement-in-confidence” hypothesis, people can learn about an advisor’s accuracy across multiple interactions according to whether the advice offered agrees with their own initial opinions, weighted by the confidence with which these initial opinions are held. We test this hypothesis using a judge-advisor system paradigm to precisely manipulate the profiles of virtual advisors in a perceptual decision making task. We find that when advisors’ and participants’ judgments are independent, people can correctly learn advisors’ features, like their accuracy and calibration, whether or not objective feedback is available. However, when their judgments (and thus errors) are correlated—as is the case in many real social contexts—predictable distortions in trust can be observed between feedback and feedback-free scenarios. Using agent-

based simulations, we explore implications of these individual-level heuristics for network-level patterns of trust and belief formation.

1 Introduction

We rely on advice in many everyday contexts, from finance and politics to education and health, but often lack immediate feedback or other objective standards with which to judge the accuracy of that advice. Yet in these contexts we must learn to distinguish good from bad advisors and consequently who to listen to. How people do this, and how reliably they do so, are open questions that have received relatively little systematic study to date (see (Weiss & Shanteau, 2003) for exceptions). The present research addresses this issue.

Although trust is a multidimensional construct, in the current work we are interested in a specific aspect of trust pertaining to the accuracy and competence of the trustee (Mayer, Davis, & Schoorman, 1995). In relation to this, we ask three related questions. First, we ask whether people can learn their advisors' competence through experience even in domains where feedback is unreliable, costly or completely absent, and in the absence of contextual cues (e.g., reputation). Second, we ask what heuristics people use to form competence representations in these contexts, and under what circumstances these heuristics are useful vs. maladaptive. Third, we explore how heuristics used by individual decision makers can drive emerging patterns of trust and influence in larger groups. We find that people can discern the usefulness of advice without the benefit of external feedback, and that simple mechanisms can explain the observed behavior.

When available, external feedback can guide learning about others through reinforcement learning mechanisms that have been thoroughly described in previous research (Behrens, Hunt, Woolrich, & Rushworth, 2008; Guggenmos, Wilbertz, Hebart, & Sterzer, 2016; Sutton & Barto, 1998). But a similar kind of social inference is likely to be necessary even when such objective external feedback is not readily available. Our hypothesis is that a solution to this seemingly computationally intractable problem lies in the use

of two readily available pieces of evidence in social decision making contexts — namely, an advisor’s agreement rate with one’s own beliefs, and one’s own internal confidence in those beliefs: If for a given decision we are certain we are correct, then we can equally be certain that anyone who disagrees with us is wrong, and accordingly down-weight their opinion in the future. If on the contrary we make a choice with less confidence, we should still down-weight our trust in that advisor in future interactions, but now to a lesser extent. Thus, people can overcome the absence of objective feedback by accumulating over time the trial-by-trial co-variation between their internal decision confidence and the actual state of the environment (Guggenmos & Sterzer, 2017; Guggenmos et al., 2016; Pescetelli, Rees, & Bahrami, 2016). We call this strategy the *agreement-in-confidence* heuristic.

Our work builds on previous research into group decisions and information integration processes (Bonaccio & Dalal, 2006; Rader, Larrick, & Soll, 2017; Sniezek & Buckley, 1989; Yaniv & Kleinberger, 2000). It is already well-known that confidence of both judge and advisors act as a weight in opinion aggregation: confident judges are less likely to ask for advice and are less influenced by advice they receive (Pescetelli, Hauperich, & Yeung, 2019, in prep.), (Tost, Gino, & Larrick, 2012); meanwhile, confident people are trusted more and are more influential within groups and juries (Penrod & Cutler, 1995; Roediger III, Wixted, & Desoto, 2012; Swol & Sniezek, 2005; Zarnoth & Sniezek, 1997), irrespective of true accuracy (Hertz, Romand-Monnier, Kyriakopoulou, & Bahrami, 2016; Mahmoodi et al., 2015). According to a confidence heuristic (Price & Stone, 2004; Pulford, Colman, Buabang, & Krockow, 2018) an advisor’s confidence signals their likely accuracy. This heuristic is normatively justified to the extent that confidence predicts objective accuracy (Bahrami et al., 2010; Henmon, 1911; Koriat, 2012), and reflects a subjective probabilistic estimate of decision accuracy (Aitchison, Bang, Bahrami, & Latham, 2015; Fleming & Daw, 2017; Meyniel, Sigman, & Mainen, 2015; Pouget, Drugowitsch, & Kepecs, 2016).

Importantly, we propose that the role of confidence goes beyond serving as an external social signal of *advice* accuracy (Bahrami et al., 2010; Price & Stone, 2004), and additionally serves as a learning signal by which an advisee can evaluate the accuracy of the *advisor* themselves, thus enabling formation of stable representations of advisor competence. In contrast with other strategies based on classification variability (Weiss & Shanteau, 2003), which require repeated observations, the agreement-in-confidence heuristic enables people to learn an advisor’s accuracy even in one-shot interactions as a basis for weighting their advice in future interactions. This strategy is consistent with various social psychological phenomena, including the “false consensus” effect, naïve realism and social judgment theory’s “latitude of acceptance”, which all indicate a tendency for people to discount disagreeing opinions, under-weight advice as a function of distance from one’s own opinion, and consider one’s own opinions as more objective or frequent than others’ (Ecken & Pibernik, 2016; Liberman, Minson, Bryan, & Ross, 2012; Minson, Liberman, & Ross, 2011; Ross, Greene, & House, 1977; Schultze, Gerlach, & Rittich, 2018; Sherif, Sherif, & Nebergall, 1965; Soll & Larrick, 2009; Yaniv, 2004). Importantly, however, our approach differs in suggesting that these phenomena are part of a normatively justified strategy that enables people to discern advisors’ features without the benefit of feedback: When a judge and advisor are independent, their rate of agreement varies as a simple monotonic function of their respective accuracies, so that agreement rate can be used to infer an advisor’s accuracy. However, a judge using this strategy will always underestimate the accuracy of the advisor, unless they themselves are perfectly accurate (Figure 1A). Using the agreement-in-confidence heuristic, a judge can have more nuanced and accurate assessments of advisors, because they can weight their learning about advisors according to their certainty in their initial view. Here we test the hypothesis that people use this agreement-in-confidence heuristic to learn about the accuracy of their advisors. Moreover, we extend these ideas to identify the boundary conditions that determine whether use of this strategy is adaptive or maladaptive depending on features of the environment, as predicted by adaptive rationality theories (Gigerenzer & Selten, 2002; Simon, 1972).

1.1 Overview of research

In a series of experiments and agent-based modelling simulations, we explore the following three key hypotheses: (H1) People can learn about the competence of their advisors even in the absence of feedback or contextual cues. (H2) This learning depends on an agreement-in-confidence heuristic that is normatively prescribed but can lead to systematic biases in learning. (H3) Use of the agreement-in-confidence heuristic by individual decision makers can influence the evolution of trust and belief within a social network.

We test the first two hypotheses—that people can learn the accuracy of advisors even in the absence of feedback, and do so based on the use of internal metacognitive information—using a judge-advisor system (JAS) task (Bonaccio & Dalal, 2006). Participants (judges) performed a series of simple perceptual judgments, first giving an initial response and associated confidence, then revising decision and confidence after having been shown the opinion of different virtual advisors with predetermined informational profiles. Crucially, we manipulated across participants the presence of objective trial-by-trial feedback. The group receiving trial-by-trial feedback provides a baseline where we expected participants to correctly judge advisor accuracy (Behrens et al., 2008). The behavior of interest is whether, in the absence of feedback, participants still learn to trust advisors differentially as a function of their objective accuracy. Two measures of perceived competence were recorded to allow for dissociations between explicit and implicit behaviors: explicit numerical ratings of advisor competence vs. measurement of the influence of their advice on participants’ decisions (Bonaccio & Dalal, 2006). We label the influence measure as “implicit” because participants were not explicitly prompted to report it.

Experiment 1 tested whether people learn crucial characteristics of advisors, like their accuracy and calibration, in the absence of feedback (H1). When feedback is readily available, both accuracy (Behrens et al., 2008) and calibration (Tenney, MacCoun, Spellman, & Hastie, 2007) have been shown to be valued advice features affecting both trustworthiness and influence. We replicate these findings, and extend them to show that people

are capable of learning these advisor features even without objective feedback.

In Experiment 1, advisors’ judgements were independent from the participant’s judgements. Experiments 2 and 3 extended Experiment 1 to explore crucial boundary conditions of consensus and confidence-based estimates of advisor competence that are a key implication of the agreement-in-confidence heuristic (H2). When judge and advisor opinions are correlated rather than independent, agreement-based heuristics will systematically overestimate the accuracy of advisors. It is plausible that such correlations are common in real world scenarios where people share the same information, for example via news sources or social groupings, or approach questions with similar biases (Del Vicario, Bessi, et al., 2016; Sunstein, 2001; Tversky & Kahneman, 1974). To mimic these scenarios of correlated opinions, we break the coupling between agreement and accuracy by creating dependence between participants’ initial judgments and the advice they receive. We show that this leads to predictable distortions in rated competence and influence, which differ from those seen when feedback is provided.

These distortions are potentially relevant to understand decision-making in many real-world environments—from online debates to information consumption—that are characterized by multiple judges and correlated information among them (Del Vicario, Bessi, et al., 2016; Kao & Couzin, 2014; Kao, Miller, Torney, Hartnett, & Couzin, 2014; Sunstein, 2001; Yaniv, Choshen-Hillel, & Milyavsky, 2009). In the final section, we use agent-based modelling to generalize the simple mechanisms of trust formation we identify empirically to consider their impact in larger groups of interacting agents. In our simulations, we show that, under realistic assumptions, the empirically-identified heuristics lead to distinct patterns of trust and belief formation at the collective level (H3). Specifically, our agent-based models suggest that Bayesian normative heuristics can lead to the emergence of clusters of individuals sharing similar biases that are persistent over time.

2 Experiment 1

Experiment 1 investigated whether people can learn about advisors’ features —specifically, their accuracy and calibration — in environments that lack objective feedback. Four virtual advisors with differing profiles were designed, who differed in accuracy and calibration (Fleming & Lau, 2014). Participants repeatedly experienced each advisor to give them the opportunity to learn about the quality of their advice. The crucial question of interest was whether the perceived competence of advisors would be sensitive to their objective accuracy and calibration, even if participants did not have access to feedback on each trial. By including a second group of participants, who did have access to trial-by-trial feedback, we could also assess the effect of using objective vs. subjective learning signals on evaluations of advisor competence.

2.1 Methods

2.1.1 Participants

There were 46 participants (26 females, age = 23 ± 0.45), half of whom were pseudo-randomly assigned to the Feedback condition and the other half to the No-Feedback condition.

2.1.2 Paradigm

The perceptual task (Boldt & Yeung, 2015) required participants to judge which of two briefly presented boxes contained more dots (Figure 1B). One box contains more dots, specifically $ndots = 200 + d$, compared to the other with $ndots = 200 - d$, with dots pseudo-randomly assigned to locations in a 20 x 20 grid anew on each trial and with equal trial numbers with left and right box as the correct answer. By manipulating the d parameter we titrated the difficulty of the task (Treutwein, 1995) to ensure similar overall accuracy across participants (nominal accuracy rate = 70.7%).

Participants registered their response and confidence judgment, unspeeded, by mouse-click on a semi-continuous scale in 10 steps, ranging from “100% sure left” to “100% sure

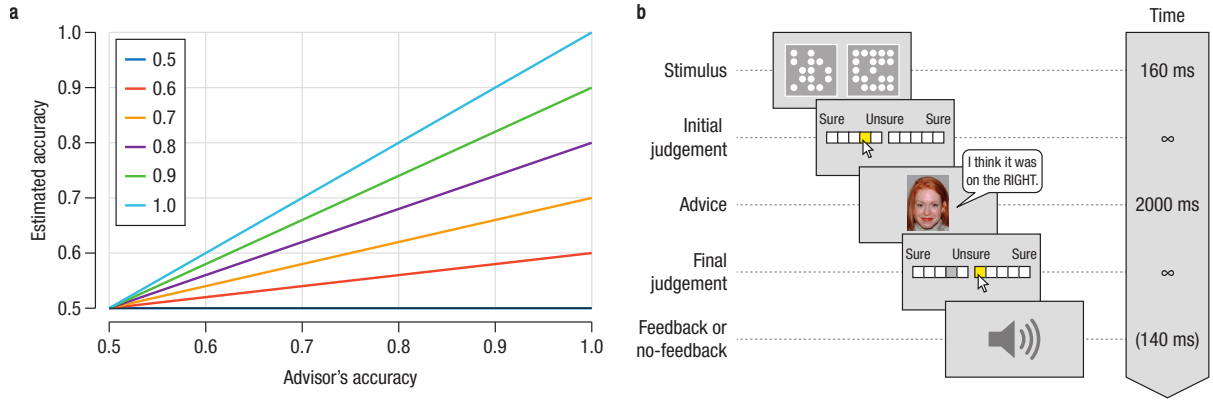


Figure 1: (a): A simple model that estimates accuracy from pure agreement rate will underestimate the accuracy of any advisor unless the model itself is always correct in its judgments: Estimated accuracy = $a \times b + (1 - a) \times (1 - b)$, where a is the objective accuracy of the judge's decisions and b is the advisor's accuracy; i.e., their likelihood of agreeing on the correct answer plus their likelihood of agreeing on an incorrect answer. (b): Schematic illustration of Experiment 1 paradigm. The computerised judge-advisor system (Sniezek & Buckley, 1989) involved on each trial a perceptual decision, advice from one of four task virtual advisers, followed by a final decision.

right". Text landmarks signalling 10% increases aided the interpretation of the scale. The middle point of the scale (50% or total uncertainty) was removed and a gap appeared instead, meaning that participants had to commit to one interval (2-alternative forced-choice). After confirming their response with the spacebar, one of four different advisors appeared centrally as a head-shot picture. Advice was provided in the form of spoken sentences that expressed a binary level of confidence (low vs. high) and either agreement or disagreement with the participant's judgment (see Supplementary Information for further details).

Participants were then given the opportunity to update their decision and confidence level using the same interface and input method as used in the pre-advice period. In the Feedback condition only, after the final decision was confirmed, a high frequency tone indicated whenever the participant's final decision was incorrect. In the No-Feedback condition, a new trial started immediately after participants had confirmed their final answer.

At the end of each block, all participants saw a summary of their post-advice percentage accuracy. Because advisors appeared equally often and in randomised order within blocks, this feedback could not favor one advisor over the others. Participants performed

500 trials across 10 experimental blocks. Prior to these, two initial blocks with a fifth advisor served as practice and were removed from all the analyses. On each experimental block, each advisor appeared ten times. Ten randomly selected trials within each block were presented with a black silent silhouette and a post-advice decision was not required (null trials), to motivate participants to provide meaningful answers in their pre-advice answers on each trial. After every two experimental blocks, participants answered a brief questionnaire about their explicit opinions about the four advisers. Four questions asked participants to directly rate on a scale from 1 (Not at all) to 50 (Extremely) how much they thought each adviser was accurate (Q1), confident (Q2), trustworthy (Q3) and influential on their own choices (Q4) (see Supplementary Information for complete description).

2.1.3 Manipulation

We orthogonally manipulated the average accuracy of the four advisors and their confidence-to-accuracy calibration. The two accurate advisors gave correct answers on 80% of trials, whereas the two inaccurate advisors gave correct answers 60% of trials. Crossed with this factor, the two calibrated advisors were always correct when expressing answers with low uncertainty (“I’m sure”) and less accurate when uncertain (“I think”), whereas the two uncalibrated advisors expressed uncertainty independently from objective accuracy. This led to the profiles shown in Table 1. Note that all advisors were equally often confident vs. unconfident across trials, to avoid participants simply trusting the advisor who was the most confident *on average* when objective feedback is not available (Sah, Moore, & Maccoun, 2013).

2.1.4 Exclusion criteria

An exclusion criterion was set *a priori* for staircase convergence. Participants who showed progressively increasing thresholds (i.e., increasing dot difference d across the experiment) were to be eliminated as this indicated that they were randomly guessing. None of the participants had to be removed when this criterion was applied to our sample. At the

Events count	Advisors			
	Accurate Calibrated	Accurate Uncalibrated	Inaccurate Calibrated	Inaccurate Uncalibrated
Incorrect Confident	0	1	0	2
Incorrect Unconfident	2	1	4	2
Correct Unconfident	3	4	1	3
Correct Confident	5	4	5	3

Table 1: Experiment 1 - Advisors profile. Values in the central section represent the number of times each event occurred over the course of a 50-trial block (10 null trials, with no advisor, are not shown). Calibration (metacognitive sensitivity) of each advisor and their informativeness are reported in Supplementary Information.

end of the experiment the average difficulty parameter d across participants (pooled data) was 9.6 ± 2.81 .

2.2 Results

The key set of analyses assessed whether participants were sensitive to advisors' features (accuracy and calibration) in the absence of feedback (i.e., the between-participants manipulation). These questions were investigated through the analysis of both explicit ratings of perceived competence and implicit influence measure, with data from all blocks collapsed given that preliminary analyses indicated no notable effects of time. A summary table is provided in the Appendix, summarizing all ANOVA results of Rated competence and Influence, for all three experiments.

2.2.1 Competence ratings

Participants provided explicit ratings of advisor competence after every second experimental block, along four dimensions: accuracy, confidence, trustworthiness and influence. An initial rating was provided before the start to assess baseline perceived competence, e.g., due to advisor appearance. Baseline ratings were then subtracted from subsequent

ones to remove these effects. Ratings were then converted into a unitary measure of perceived competence via principal components analysis (see Supplementary Information). We used PCA over simple averaging because (a) we were agnostic on which dimensions of advisor competence were made salient by our manipulation; (b) we did not want to average over distinct constructs (e.g., accuracy and trust).

A mixed-design ANOVA on the resulting rated competence scores, with factors of Feedback (between-participants) and advisor Accuracy and Calibration (both within-participants), revealed significant main effects for advisor Accuracy ($F(1, 44) = 9.68, p = .003, \eta_G^2 = 0.079$) and advisor Calibration ($F(1, 44) = 12.32, p = .001, \eta_G^2 = 0.076$), but not Feedback ($F < 1$). Participants gave higher competence ratings for accurate over inaccurate advisors, and for calibrated over uncalibrated advisors. No interaction term reached significance ($F(1, 44) < 1.9, p > .16$). Importantly, neither within-participants manipulation interacted with Feedback, suggesting that participants were sensitive to the accuracy of the advice and that this sensitivity did not vary consistently according to the presence or absence of feedback (Figure 2A).

We next ran planned 2-way ANOVAs separately for each feedback group to assess whether the overall patterns described above were also reliable in each group. In the Feedback group this analysis revealed reliable main effects of both advisor Accuracy ($F(1, 22) = 8.26, p = .008, \eta_G^2 = .09$) and advisor Calibration ($F(1, 22) = 8.71, p = .007, \eta_G^2 = 0.12$), but no reliable interaction ($F(1, 22) = 1.24, p = .27, \eta_G^2 = .02$). In the No Feedback group, although comparable numerical trends were apparent, neither main effect of Accuracy ($F(1, 22) = 3.42, p = .07, \eta_G^2 = 0.06$) and Calibration ($F(1, 22) = 4.03, p = .05, \eta_G^2 = 0.04$) reached statistical significance. The interaction term was not significant ($F < 1$).

2.2.2 Influence

Influence was quantified as the signed difference between post and pre-advice confidence (Equation 3, Supplementary Information), and represents the shift in a participant's expressed judgement observed after social information. We replicated the results sep-

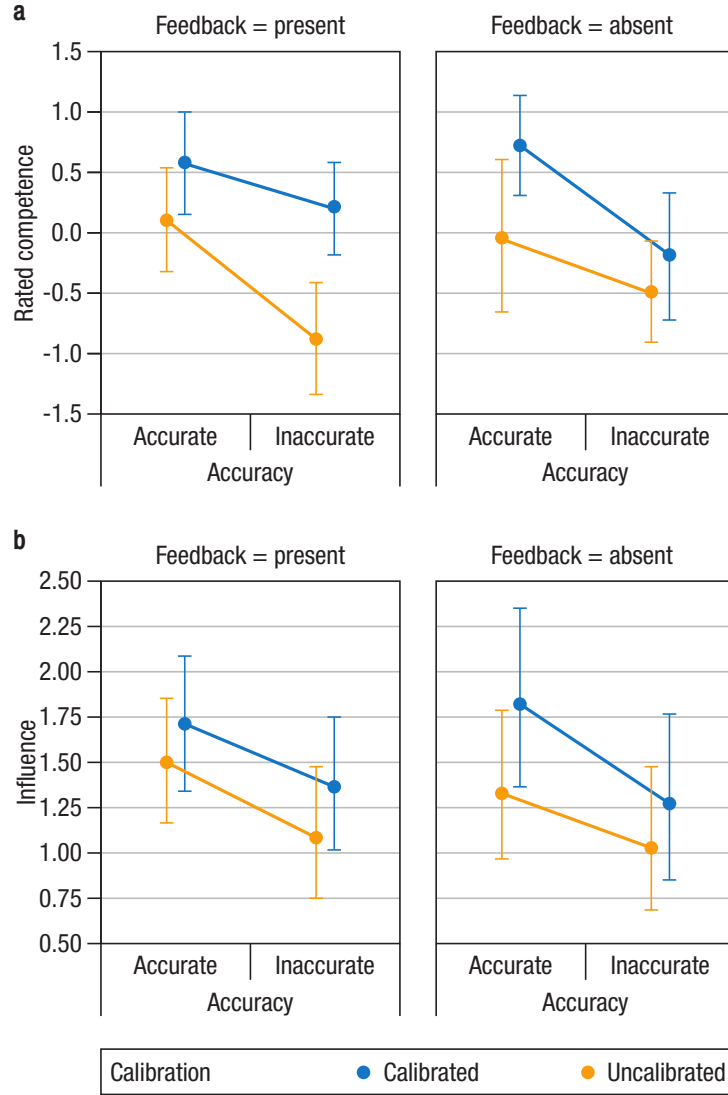


Figure 2: Experiment 1 - rated competence and influence patterns for humans and simulations. (a) Average competence ratings in the two feedback groups as a function of advisor accuracy and calibration. Error bars represent 95% bootstrap confidence intervals. (b) The effect of advisor Accuracy, advice Confidence and Calibration, and Feedback group on the influence measure, averaging across the two levels of advisor confidence. Error bars represent 95% bootstrap confidence intervals.

arately for agreement and disagreement trials (SI, §2). Influence results (Figure 2B) were analysed using a mixed-design ANOVA that included the same factors of Feedback presence (between-participants), advisor Accuracy and Calibration (both within-participants) as above, together with the additional within-participants factor of the Confidence expressed by the advisor on a given trial. This analysis revealed significant main effects of advisor Accuracy ($F(1,44) = 14.80, p < .001, \eta_G^2 = 0.02$) and Calibration ($F(1,44) = 15.84, p < .001, \eta_G^2 = 0.01$), mirroring the pattern of results

seen for rated competence. A reliable main effect of advisor Confidence ($F(1, 44) = 55.82, p < .001, \eta_G^2 = .12$) unsurprisingly indicated that more confidently expressed advice had greater influence. There was a significant interaction between Calibration and advisor Confidence ($F(1, 44) = 9.62, p = .003, \eta_G^2 = 0.004$) indicating that the effect of confidently (vs. unconfidently) expressed advice was greater for calibrated advisors than uncalibrated advisors. The analysis also revealed a significant 3-way interaction between advisor Accuracy, Calibration and Confidence ($F(1, 44) = 4.75, p = .03, \eta_G^2 = 8.5e - 04$), indicating that the two-way interaction between Calibration and Confidence was larger for inaccurate advisors than accurate ones. Importantly, there was no reliable main effect of Feedback ($F < 1$), nor any reliable interaction between Feedback and any other main effects (all $F_s(1, 44) < 1.1, ps > .29$) suggesting again that participants' sensitivity to advisor accuracy, here expressed in terms of the influence of their advice, did not depend significantly on the provision of trial-by-trial feedback.

Although we observed no reliable effects of feedback, we ran planned follow-up ANOVAs for each feedback group separately to assess whether the overall patterns described above were also reliable in each group. These analyses revealed significant effects in both groups of advisor Accuracy ($F_s(1, 22) > 6.87, ps < .05, \eta_G^2 > 0.02$) and Calibration ($F_s(1, 22) > 6.35, ps < .01, \eta_G^2 > .01$), and advice Confidence ($F_s(1, 22) = 22.95, ps < .001, \eta_G^2 > 0.07$). For the Feedback group, we also found a reliable two-way interaction between Calibration and Confidence ($F(1, 22) = 7.37, p = .01, \eta_G^2 = .008$) and a reliable three-way interaction between Accuracy, Calibration and Confidence ($F(1, 22) = 8.32, p = .008, \eta_G^2 = .003$), indicating that the influence difference between calibrated and uncalibrated advisors was mainly shown for highly confident advice compared to uncertain advice, and more so for inaccurate advisors than accurate ones. Similar patterns were seen numerically in the No Feedback group, but the effects were not statistically reliable ($F_s(1, 22) < 2.59, ps > .12$).

2.3 Discussion

Overall, participants in Experiment 1 gave higher ratings of competence for, and were more influenced by, advisors who were characterised by high accuracy rates and high calibration. The finding that advisor confidence and calibration affect perceptions of advice is broadly consistent with previous findings (Price & Stone, 2004; Sah et al., 2013; Sniezek & Van Swol, 2001; Tenney et al., 2007). However, the present results extend this work to show, in a highly-controlled perceptual task, that participants distinguish advice along these dimensions even in the absence of objective feedback. Thus, although rated competence and influence were not identical across groups, no consistent (statistically reliable) differences were observed, and overall participants in the No Feedback group were able to learn distinguishing characteristics of advisors (accuracy and calibration).

These empirical results are consistent with the hypothesis that, in the absence of feedback, people make use of internal signals to evaluate the quality of advice they receive. In algorithmic simulations based on straightforward implementations of two related learning strategies—estimating advice accuracy according to simple agreement, or agreement-in-confidence—we show how this differentiation of advice quality can be achieved (Supplementary Information §3). As such, the combined empirical and simulation results indicate the normative value of these heuristic strategies in enabling people to discern the usefulness of advice in feedback-poor environments: When advice is independent from a judge’s initial opinion, agreement and confidence covary with accuracy and thus are useful cues to integrate over time so to learn about the competence of an advisor. Extending these ideas, Experiments 2 and 3 aimed to test the agreement-in-confidence hypothesis more directly, and to explore crucial limitations in the use of agreement and confidence in estimating advisor accuracy, by exposing participants to advice that was not independent of their own initial decisions.

3 Experiment 2

Experiment 1 showed that people are able to detect subtle advisor differences even in the absence of feedback. According to our hypothesis, they achieve this differentiation of advisor competence using an agreement-in-confidence heuristic, whereby feedback is replaced by the interaction of past agreement with the advisor and internal metacognitive signals: Specifically, to the extent that an observer’s internal confidence is calibrated (i.e., predictive of their objective accuracy), they can validly learn whether about advisor’s accuracy (vs. inaccuracy) by integrating agreement (vs. disagreement) with their own confidently held opinions.

Giving weight to advice that agrees with one’s own view seems normatively appropriate from a Bayesian standpoint (Dawes, 1989; Krueger & Clement, 1994), but depends critically on the assumption that observers are independent. However, people’s judgments are rarely independent and distorted only by random noise (Koriat, 2012; Krause, Ruxton, & Krause, 2010): Dependence between individuals’ opinions can arise from use of similar cognitive heuristics that lead to similar reasoning errors or information sampling (Tversky & Kahneman, 1974; Vandormael, Herce Castañón, Balaguer, Li, & Summerfield, 2017), from being exposed to similar signals (Kao & Couzin, 2014; Kao et al., 2014) or belonging to the same social clique (Jamieson & Cappella, 2008; Jasny, Waggle, & Fisher, 2015; Sunstein, 2001). People tend not to take into account advisors’ judgments interdependence (Yaniv et al., 2009). As a result, crowds are known to be susceptible to error cascades (Le Bon, 1895; Mackay, 1841), economic bubbles (De Martino, O’Doherty, Ray, Bossaerts, & Camerer, 2013), polarisation (Myers & Lamm, 1976) and *groupthink* (Janis, 1972; Turner & Pratkanis, 1998). Moreover, in the context of advice, advisors’ opinions and suggestions might be deliberately framed in relation to the advisee. For any or all of these reasons, advice may not be independent of a decision maker’s judgment in many realistic scenarios. Heavily weighing agreeing advisors might in these cases be maladaptive. As a first exploration of scenarios with correlated advice, we thus investigated whether people rely on accurate advisors or instead advisors who tend to agree with their own initial judgment.

In this experiment, advisor accuracy was manipulated so that two advisors were on average highly accurate (around 80% accuracy) and two advisors were on average relatively inaccurate (around 60% accuracy). Orthogonally, advisor agreement rate with the participant’s initial judgment was manipulated to create two advisors who agreed with the participant frequently (around 80% of trials) and two advisors who tended to have a lower agreement rate with the participant (around 60%). We conceive of an advisor with low accuracy but a high rate of agreement as someone who shares biases with the participant and so makes similar (correlated) mistakes. The accurate but disagreeing advisor conversely represents an advisor who uses different information and therefore tends to be correct when the participant makes mistakes, and vice versa. Of interest was the impact on the perceived competence and influence of each advisor, as a function of the advisors’ accuracy and agreement rates and, separately, the participants’ access to objective feedback.

3.1 Methods

3.1.1 Participants

The experiment included 46 participants, equally divided between the two feedback groups (37 females in total, 18 of whom were in the Feedback group, mean age: 21.63 ± 3.02).

3.1.2 Paradigm

The overall design was very similar to Experiment 1, with advice provided by four virtual advisors, characterized by distinct informational profiles, appearing in the context of the same dot-count perceptual decision task. Advice was always presented in the form of a binary left/right judgment (i.e., with no accompanying indication of high vs. low advice confidence), which could agree or disagree with participants’ original judgments. Participants completed ten blocks of 44 trials each. The presentation of the four advisors was randomly shuffled across trials within-block, with each advisor appearing exactly ten times. Four additional trials served as null trials (as above). Explicit ratings of advisor

competence along four dimensions were again collected at the end of every second block and aggregated as before: accuracy (Q1), likeability (Q2), trustworthiness (Q3) and influence (Q4).

3.1.3 Manipulation

To disentangle advisor agreement rate and accuracy, the probability of agreement conditional on the participant’s choice accuracy was manipulated. Through the staircase procedure it was expected that all participants would converge to an accuracy level of about 70%. This enabled us to manipulate advisor accuracy and agreement rate separately, by pre-determining the probability of agreement differently according to whether the participant’s initial decision was correct vs. incorrect. Both accuracy and agreement were manipulated to have two levels (high=80% and low=60%). This gave rise to the four advisor profiles defined in Table 2. Probabilities are expressed as a fraction over the number of participants’ expected correct (7) and incorrect (3) judgments, over the number of encounters with one advisor during one block (10). Of key interest was the separate impact of advisor accuracy and agreement rate on explicit ratings of competence and implicit measures of advisor influence, separately for the Feedback and No Feedback groups.

3.1.4 Exclusion criteria

The first two experimental blocks were removed from the analysis to allow the staircase procedure to fully adapt to each individual’s threshold. This was necessary given that our manipulation was heavily dependent on the expected accuracy rate of the participants. A further exclusion criterion was set to exclude all participants whose threshold never converged, which suggests a random response strategy. None of the participants had to be removed on the basis of this criterion. The perceptual task difficulty d (dot difference between boxes) after staircasing was 9.93 ± 2.96 (pooled data).

	Advisors			
	High Accuracy High Agreement	High Acc. Low Agr.	Low Acc. High Agr.	Low Acc. Low Agr.
$p(Agr Cor_s)$	6.5/7	5.5/7	5.5/7	4.5/7
$p(Agr Inc_s)$	1.5/3	0.5/3	2.5/3	1.5/3
Expected Acc. rate	80%	80%	60%	60%
Expected Agr. rate	80%	60%	80%	60%

Table 2: Experiment 2 advisors’ profiles. Expected accuracy and agreement rates of different advisors are disentangled by manipulating the probability of the advice agreeing with the participant, conditional on the participant’s accuracy. In the table, probabilities are expressed as a fraction of the number of participants’ expected correct (7) and incorrect (3) judgments, during the number of encounters with one advisor (10) in a single experimental block. We report each advisor information value in supplementary information.

3.2 Results

As for Experiment 1, separate analyses were conducted on rated competence and influence measures. Of interest was whether the manipulated within-participants factors (advisor agreement rate and accuracy) varied in impact across two groups of participants, with differential access to objective trial-by-trial feedback.

3.2.1 Competence ratings

A mixed-design ANOVA was run on competence ratings with feedback group as a between-participants factor and advisor accuracy (low vs. high) and agreement rate (low vs. high) as within-participants factors. This analysis revealed significant main effects for both accuracy ($F(1, 44) = 8.36, p = .005, \eta_G^2 = .06$) and agreement rate ($F(1, 44) = 22.52, p < .001, \eta_G^2 = .1$), but not for feedback ($F < 1$). A significant interaction between feedback and accuracy ($F(1, 44) = 8.41, p = .005, \eta_G^2 = .06$) indicated much greater impact of advisor accuracy on participants’ rated competence when feedback was available than when it was absent (Figure 3A). There was no reliable interaction between feedback and agreement rate ($F(1, 44) = 1.88, p = .17, \eta_G^2 = .01$) or between agreement and accuracy ($F < 1$), nor a significant three-way interaction ($F < 1$).

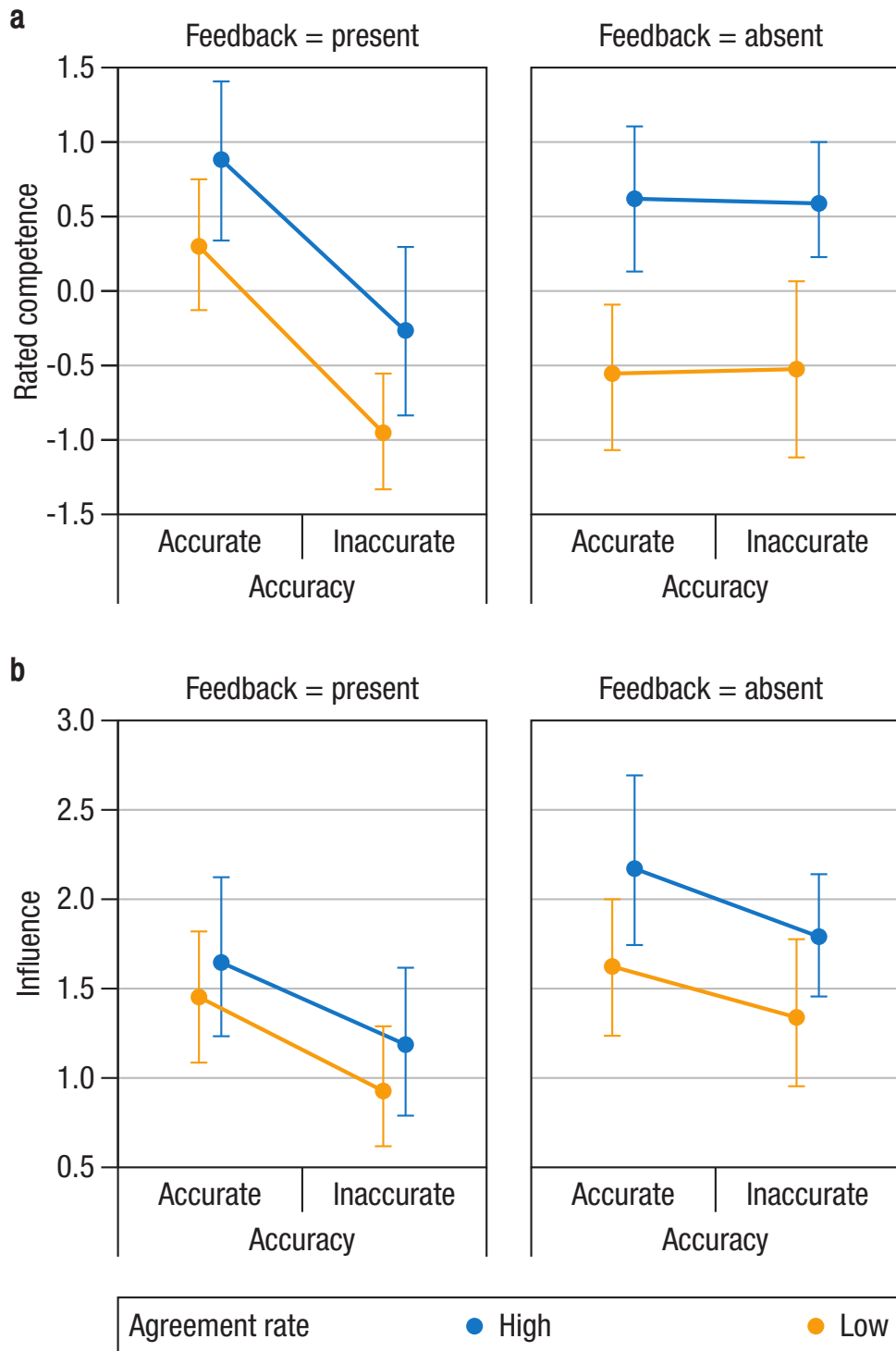


Figure 3: Experiment 2 - rated competence ratings and influence for human participants and simulations. (a) Average competence ratings in the two feedback groups, divided by advisor accuracy rate and agreement rate. Error bars represent 95% bootstrap confidence intervals. (b) Influence that advice had on participants' opinions in the two feedback groups and divided by advisor accuracy rate and agreement rate. Error bars represent 95% bootstrap confidence intervals.

Separate planned ANOVAs were conducted for the two groups separately. Analysis of competence ratings from the Feedback group revealed significant effects of both advisor's Accuracy ($F(1, 22) = 21.36, p < .001, \eta_G^2 = .20$) and Agreement rate ($F(1, 22) = 5.62, p = .02, \eta_G^2 = .06$), but no significant interaction ($F < 1$). The effect of accuracy was much stronger than the one observed for agreement as indicated by the generalised eta squared values (Bakeman, 2005). Nonetheless, agreement rate had a significant effect on rated competence even when controlling for objective accuracy. Analysis of competence ratings from the No-Feedback group found that Agreement rate ($F(1, 22) = 18.91, p < .001, \eta_G^2 = .18$) but not Accuracy ($F < 1$) affected participants' explicit competence ratings, with no significant interaction ($F < 1$). These findings suggest that when feedback was not directly available to estimate partners' accuracy, participants rated agreeing advisors as more competent than disagreeing ones, but did not consistently differentiate accurate vs. inaccurate advisors.

3.2.2 Influence

The next analysis focused on influence, our implicit measure of perceived competence, using the same mixed-design ANOVA as above. Table 2 in the Appendix summarises the results of this ANOVA and corresponding analyses in Experiments 2 and 3. This analysis revealed significant main effects of advisor Accuracy ($F(1, 44) = 14.79, p < .001, \eta_G^2 = .04$) and Agreement rate ($F(1, 44) = 13.91, p < .001, \eta_G^2 = .03$), with no significant interactions, including between Accuracy and Feedback ($F < 1$) and between Agreement rate and Feedback ($F(1, 44) = 2.01, p = .16, \eta_G^2 = .005$). Thus, participants were more influenced by accurate advisors and by advisors characterized by high agreement rates with their own judgments, and these effects did not vary consistently as a function of whether or not trial-by-trial feedback was available. We replicate the results for agreement and disagreement confidence changes separately (SI §2).

Planned ANOVAs on the influence measure for each group separately revealed, for the Feedback group, significant effects of advisor Accuracy ($F(1, 22) = 12.71, p = .001, \eta_G^2 = .06$) and Agreement rate ($F(1, 22) = 5.81, p = .02, \eta_G^2 = .01$), with no significant in-

teraction ($F < 1$). As would be expected, when feedback was available, participants were more influenced by accurate advisors over inaccurate ones. More surprisingly, they were also more influenced by advisors who more often agreed with their own judgment, even though feedback was available that indicated equivalent accuracy rates for pairs of advisors characterized by different agreement rates. Analysis of advisor influence in the No-Feedback group revealed a similar pattern, with significant main effects for both advisor Agreement rate ($F(1, 22) = 8.60, p = .007, \eta_G^2 = .06$) and Accuracy ($F(1, 22) = 4.09, p = .05, \eta_G^2 = .02$), although the effect size of Agreement was the greater of the two. No significant interaction was observed ($F < 1$). Thus, in the No Feedback group, the results suggest a dissociation between what people reported in explicit competence ratings, with higher competence reported for agreeing advisors regardless of accuracy, and what was apparent in advisors' influence on participants' decisions, which showed effects of both accuracy and agreement rate.

3.3 Discussion

The results of this experiment reveal a strong effect of advisor agreement on perceived competence, regardless of the presence or absence of objective feedback and of the objective accuracy of advice: Participants showed greater reliance on advisors who more often agreed with their own initial judgments, both in their explicit ratings of competence and in the degree to which the advisors influenced their final decisions. This impact of agreement is in line with established findings in social psychology — e.g., phenomena of naïve realism and latitude of acceptance — that people tend to discard disagreeing opinions and tend to see their own subjective views as being more objective than others' (Liberman et al., 2012; Minson et al., 2011; Sherif et al., 1965; Shultz, Katz, & Lepper, 2001). Our findings indicate that these effects extend to perceived competence that is learnt through repeated interaction with different advisors, and moreover that it persists even when participants have access to objective feedback. The agreement effect apparent in the Feedback group— not predicted by a simple *Accuracy* heuristic that learns advisor accuracy based on observed feedback (Supplementary Information §3)— suggests an

intrinsic value of agreement when learning about advisor competence.

When feedback was unavailable, participants' reliance on advisors was dominantly determined by advisors' agreement rate, with much weaker effects of the objective accuracy of their advice. Nevertheless, at least for the implicit measure of advisor influence, we found that participants were still able to distinguish more accurate from less accurate advisors. Importantly, going beyond previous findings, this result cannot be explained in terms of participants simply downweighting advice on a given trial when the advice disagrees with their initial opinion (Lieberman et al., 2012; Sherif et al., 1965). Rather, we observe patterns of trust and influence that reflect learning across aggregated sets of trials, such that that advisors who more regularly disagreed with a participants' choices were less influential, even on those occasions where their advice happened to agree with the participants' view. Conversely, we find that advisors who more regularly agreed with the participant were more influential, even on those occasions where their opinion diverged from the participants' initial choice.

The subtle but consistent effect of advisor accuracy—even when feedback was absent and agreement rates were matched—is predicted by our agreement-in-confidence heuristic, which learns about the trust of an advisor by accumulating agreement instances weighed by internal decision confidence (Supplementary Information, Figure S6). The absence of this effect in explicit competence ratings is the only instance we found in our experiments of a systematic divergence between explicit ratings of competence and implicit measures of influence, which generally produced very similar patterns of results. The dominance of agreement over accuracy in explicit ratings might have resulted in a halo-dumping effect (Clark & Lawless, 1994) whereby, when prompted to discriminate among advisors, participants used only the most accessible dimension (i.e., agreement rate). In Experiment 3, we therefore matched all advisors in terms of agreement and accuracy rates, while at the same time varying the amount of shared information and bias between the advisors and the participants, to provide a direct test of the hypothesis that internal confidence judgments are used when making inferences about the competence of others.

4 Experiment 3

Experiment 3 aimed to provide stronger evidence that subjective confidence contributes to the formation of judgments about advice accuracy. Here, we manipulated the probability that advice would agree with a participant’s initial judgment, conditional on their initial confidence in that judgment. There were three advisors: (1) an unbiased advisor who tended to agree with the participant’s choice about 70% of the time, independent of participants’ confidence; (2) a “bias-sharing” advisor who more often agreed with the participant’s initial choice when the participant expressed high confidence in this choice; and (3) an “anti-bias” advisor who more often agreed with the participant when the participant was unsure in their initial decision. Crucially, overall agreement rate and accuracy was identical across the three advisors. The labeling of the advisors in terms of bias does not reflect an actual bias manipulation, but rather reflects our aim to capture a property that may hold in many real-world situations, where the individuals share biases in their opinions and choices (e.g., reflecting their shared reliance on common sources of information).

4.1 Methods

4.1.1 Participants

50 participants were tested and divided in the two experimental groups. Due to participants failing to attend or complete sessions, numbers across groups were unbalanced with 24 participants in the No-Feedback group and 26 in the Feedback group.

4.1.2 Paradigm

The experiment consisted of 12 experimental blocks of 35 trials each, with each of the three advisors seen on 10 trials and with 5 null trials. The perceptual task was the same as for Experiments 1 and 2, except that a different confidence rating scale was used because the 10-point scale used in Experiments 1 and 2 would not allow us to distinguish fine gradations in confidence needed here. Instead, participants rated their confidence on

a 100-point scale (50 points per interval, left vs. right). Two blocks of 25 trials served as the practice blocks and used a fourth practice advisor. Questionnaires were again administered every two blocks.

4.1.3 Manipulation

The advice profiles of the three advisors were manipulated so that advisors were matched for accuracy (70%) and their overall agreement rate (70%) with participants' initial decisions. The pattern of agreement was manipulated, however, such that the three different advisors' likelihood of agreement varied according to the participant's initial confidence (Table 3 and Figure 4). To this end, the distribution of the participant's pre-advice confidence judgments was divided into three confidence bins: the low, middle and high confidence bins, comprising 30%, 40% and 30% of trials, respectively. On trials in which the participant's initial judgment was correct, the three advisors had different agreement patterns across these bins. An unbiased advisor had a probability of agreement of 70% independent of the participant's confidence. A bias-sharing advisor had an 80% probability of agreeing when the participant was highly confident and 60% when the participant expressed low confidence in their initial decision. An anti-bias advisor, conversely, had 60% probability of agreement when the participant was highly confident and 80% when s/he was uncertain. All three advisors had equal chance of agreement when the participant's decision was correct and pre-advice confidence fell in the middle bin (70%). Likewise, all advisors had a 30% agreement rate independent of the participant's confidence when the participant's initial decision was incorrect. This ensured that all advisors were matched across all trials in terms of average agreement rate ($0.7 * 0.7 + 0.3 * 0.3 = 0.58$) and accuracy ($0.7 * 0.7 + 0.3 * 0.7 = 0.7$).

By limiting analyses to trials within the intermediate confidence bin, we could compare advisors on trials that were matched for confidence and a priori likelihood of advice agreement. The confidence reference distribution used to assign trials to the low, middle and high confidence bins was first set up on the basis of each participant's confidence ratings in the first two practice blocks. The reference distribution was updated after each

block to reflect the distribution of confidence judgments provided during the previous two blocks, to allow for possible shifts of confidence during the course of the experiment.

	Advisors		
	Bias-sharing	Unbiased	Anti-bias
$p(\text{Agree} \text{Correct}^s, \text{Confidence}_{\text{low}}^s)$	60%	70%	80%
$p(\text{Agree} \text{Correct}^s, \text{Confidence}_{\text{mid}}^s)$	70%	70%	70%
$p(\text{Agree} \text{Correct}^s, \text{Confidence}_{\text{high}}^s)$	80%	70%	60%
$p(\text{Agree} \text{Incorrect}^s)$	30%	30%	30%

Table 3: Experiment 3 advisors’ profiles. Agreement probability of different advisors is manipulated conditional on the participant’s pre-advice confidence and accuracy. This manipulation allowed to create three different advisors who were matched in terms of agreement rate and accuracy, but who differed in terms of information value (see Supplementary Information §1.2 for details).

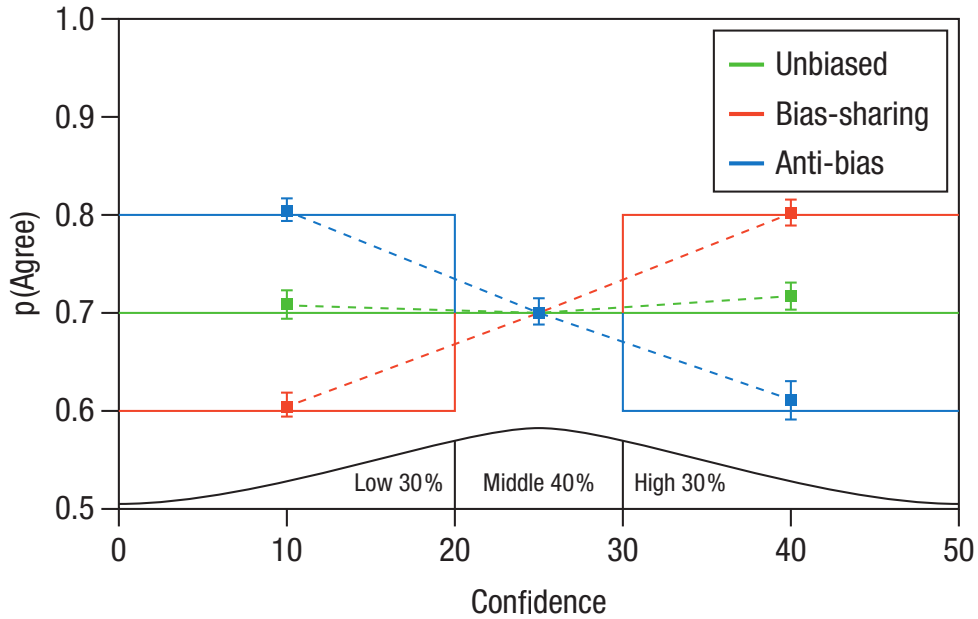


Figure 4: Experimental manipulation of Experiment 3. The probability of advisor’s agreement conditional on participant’s accuracy and confidence was manipulated so that the three advisors differed in their pattern of agreement (i.e., bias) despite being equal on average agreement rate and accuracy rate. Continuous lines represent expected agreement rates, dashed lines represent empirical data pooled across the two feedback groups.

Following Experiment 2, we expected different patterns of results to emerge from feedback and feedback-free conditions. If people use subjective confidence to learn about advisors’ accuracy (*Confidence* model in Supplementary information §3), we should observe that competence ratings and influence will favor the bias-sharing advisor over the

other advisors when feedback is unavailable. This is because the participant will experience more high-confidence agreements and fewer low-confidence agreements with the bias-sharing advisor compared to the other two. In contrast, heuristics using simple agreement counts or objective feedback (*Consensus* and *Accuracy* models in Supplementary Information §3) would not distinguish among advisors, given that advisors are matched in terms of objective accuracy and agreement rates.

4.1.4 Exclusion criteria

An exclusion criterion based on staircase convergence was set so to exclude all participants who showed random guessing. Application of this criterion resulted in the exclusion of one participant from the Feedback group and one participant from the No-Feedback group, leaving a total of 25 and 23 participants in these groups, respectively. Average difficulty parameter d was 9.98 ± 2.82 (pooled data).

4.2 Results

Degrees of freedom were corrected for violations of sphericity according to the Greenhouse-Geisser procedure, with epsilon values reported as appropriate.

4.2.1 Competence ratings

A mixed-design ANOVA on competence ratings from the end-of-block questionnaires revealed no significant main effect of Advisor Type ($F(2, 92) = 1.66, p = .19, \eta_G^2 = .03, \epsilon = 0.99$) or Feedback ($F < 1$), but a significant interaction between these factors ($F(1, 92) = 6.64, p = .002, \eta_G^2 = .12, \epsilon = 0.99$). Figure 5A shows how advisor perceived competence varied according to the presence vs. absence of feedback.

Planned follow-up one-way ANOVAs for each group separately revealed, for the Feedback group, a significant effect of Advisor Type ($F(2, 48) = 4.90, p = .01, \eta_G^2 = .16$), with rated competence being highest for the anti-bias advisor, intermediate for the unbiased advisor, and lowest for the bias-sharing advisor. Pairwise comparisons indicated that the bias-sharing advisor was perceived significantly less accurate than the anti-bias

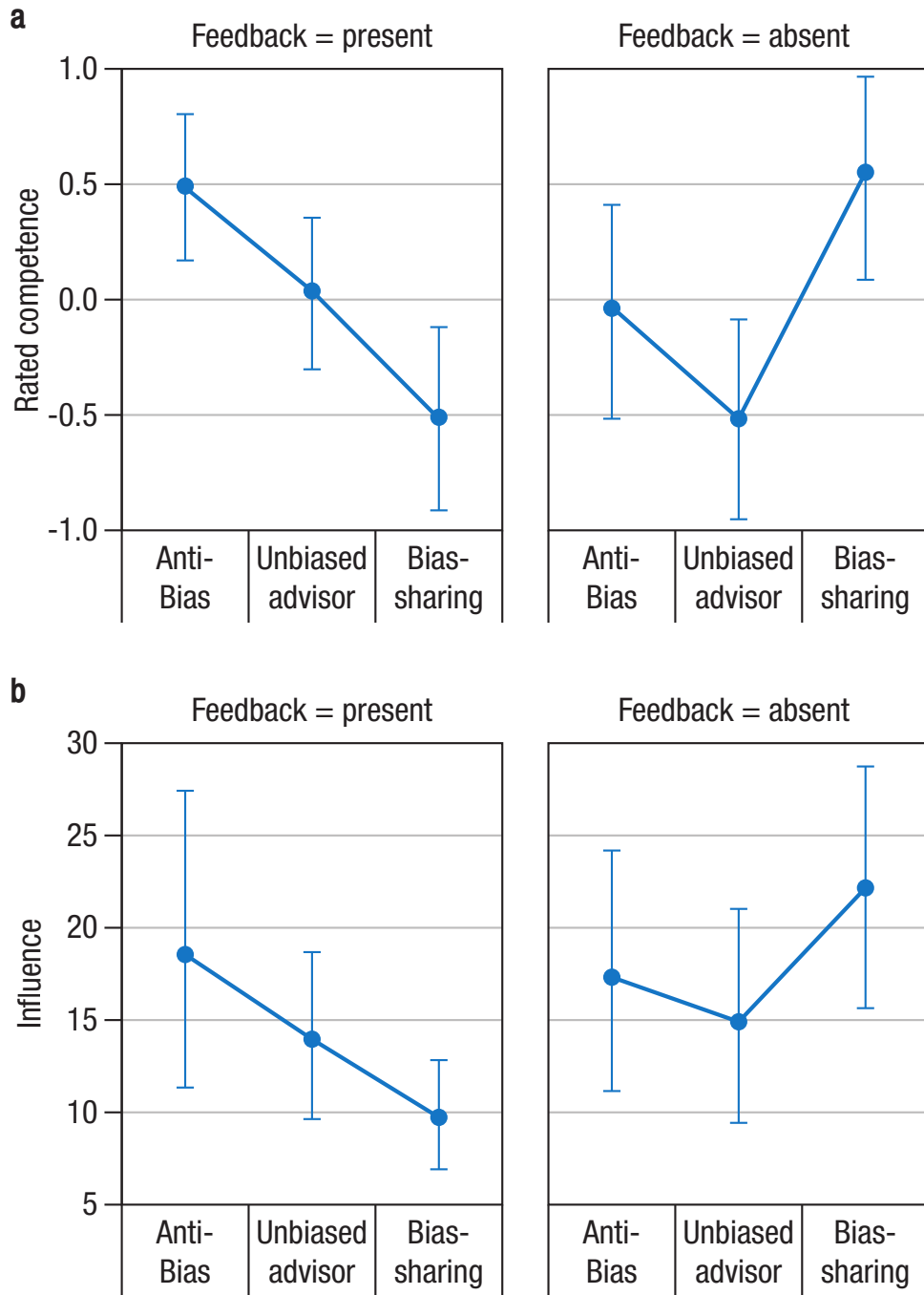


Figure 5: Experiment 3 - competence ratings and influence of human participants and simulations. (a) The effect of advisor type on competence ratings, separately for the Feedback and No-Feedback condition. Error bars represent 95% bootstrap confidence intervals. (b) The effect of advisor type on the influence measure, divided by the two feedback conditions. Error bars represent 95% bootstrap confidence intervals.

598 advisor ($t(24) = 3.09, p = .004, d = 1.07$), with no reliable differences observed otherwise
599 ($ts(24) < 1.60, p > .12, d < 0.56$). This pattern, with a high level of reliance on the
600 anti-bias advisor and lower in the bias-sharing advisor is in accordance with the pattern
601 of information gain of each advisor (Table S1), rather than advisors' objective accuracy.

Similar results are reported below for influence, suggesting this pattern is robust to different measures of trust.

A corresponding analysis on competence ratings from the No-Feedback group also revealed a significant difference among Advisors ($F(2, 44) = 3.56, p = .03, \eta_G^2 = .13$), but the pattern observed was very different. As predicted, rated competence was highest for the bias-sharing advisor when feedback was absent, numerically so compared with the anti-bias advisor ($t(22) = 1.48, p = .07, d = 0.53$, one-tail), and significantly so compared with the unbiased advisor ($t(22) = 2.81, p = .005, d = 0.98$, one-tail). The anti-bias and the unbiased advisors did not significantly differ from each other ($p > .1, d = 0.41$), but rated competence was, unexpectedly, numerically higher for the anti-bias observer than the unbiased advisor.

4.2.2 Influence

Similar patterns of trust were apparent as measured implicitly via advisor influence (confidence change following advice). A mixed-design ANOVA on measured influence revealed no significant main effect of Feedback ($F(1, 46) = 1.36, p = .24, \eta_G^2 = .01$) nor Advisor ($F(2, 92) = 1.12, p = .33, \eta_G^2 = .009, \epsilon = 0.77$), but a significant interaction between the two ($F(2, 92) = 4.80, p = .01, \eta_G^2 = .03, \epsilon = 0.77$) (Figure 5B).

When looking at influence in the Feedback condition only, a one-way ANOVA revealed a marginally significant effect of Advisor Type ($F(2, 48) = 2.88, p = .06, \eta_G^2 = .05$). Planned comparisons (two-tailed t-tests) showed that the bias-sharing advisor was less influential than both the anti-bias advisor ($t(24) = 1.98, p = .05, d = 0.54$) and the unbiased advisor ($t(24) = 2.26, p = .03, d = 0.44$). The anti-bias advisor was numerically more influential than the bias-sharing advisor, but this difference was not reliable ($t(24) = 1.10, p = .28, d = 0.26$). Similar numerical trends were observed in both agreement and disagreement trials separately (SI §2). The effects remained significant when considering disagreement trials only.

In the No-Feedback condition, there was a significant effect of Advisor ($F(2, 44) = 3.25, p = .04, \eta_G^2 = .03$). Planned comparisons showed that the bias-sharing advisor was

significantly more influential than the unbiased advisor ($t(22) = 2.63, p = .007, d = 0.46$, one-tail) and numerically more influential than the anti-bias advisor ($t(22) = 1.46, p = .07, d = 0.29$, one-tail). No significant difference was found between the unbiased and the anti-bias advisors ($p > .1, d = 0.16$), but the direction of the difference was the same as for rated competence, with the anti-bias advisor somewhat more influential on participants' decisions and confidence than the unbiased advisor. Thus, these results seem to suggest that the bias-sharing advisor was more influential than the other two when trial-by-trial feedback was not available.

In summary, as in Experiment 2, we find that participants' perception of advisor accuracy varies systematically according to whether trial-by-trial feedback is present or absent when advice is correlated with their own initial decisions. In particular, advisors' perceived competence follows their relative informativeness when feedback is provided, but when feedback is absent participants tend to rely on advisors who share their judgment biases (i.e., who agree with their confidently held judgments). Of three heuristic models fitted to these empirical data—using accuracy, agreement and agreement-in-confidence learning signals respectively—only the agreement-in-confidence heuristic was able to produce diversified beliefs about the advisors (Supplementary Information §3).

4.3 Experiment discussion

Experiment 3 showed again that systematic differences in perceived competence emerge (both in humans and models) as a function of feedback when advisors' judgements are non-independent from the advisee's, and further demonstrated a strong influence of our agreement-in-confidence manipulation on trust. Here, the presence of feedback partly reversed the pattern of competence ratings and influence measure that was observed when trial-by-trial feedback was unavailable: When objective feedback was provided, people perceived as more competent, and were more influenced by, the advisor who more frequently agreed with them when they themselves were unsure (vs. less frequently when they were sure), and showed less trust in an advisor who tended to agree with them

in decisions in which they were already sure. Evidently, though advisor accuracy and agreement are critical determinants of perceived competence (as evidenced in Experiment 2), participants remain sensitive to other dimensions of advice. Here, they appeared sensitive to the informational value (see Supplementary Information §1.2) of the advisor and not only to accuracy *per se*: They relied more on advisors whose judgments were less redundant with their own.

On the contrary, when objective trial-by-trial feedback was removed, the pattern of results partly reversed such that competence ratings and influence were greatest for the advisor who tended to agree with participants' confidently-made judgments. These advisors were more influential also on those occasions where they disagreed with participants' initial judgements, indicating that participants were actively learning about their advisors' overall competence, rather than simply discounting disagreeing advice (Supplementary information §2). Surprisingly, participants did not perceive as least accurate the anti-bias advisor who agreed with them more frequently when their initial judgment was made with low confidence. If anything, trust was lower in the unbiased advisor. The difference was not reliable and hence must be interpreted with caution, but we note a possible link to an existing proposal suggesting that scarcely differentiated judgments across different observations tend to be indicative of lower expertise (Weiss & Shanteau, 2003). Importantly, regardless of the explanation of this unexpected result, the findings of Experiment 3 demonstrate participants' sensitivity to advisors in relation to their own confidence in their judgments, and learn differentiated patterns of competence accordingly.

5 Network-level impacts of individuals' trust strategies: an agent-based simulation

The experimental and simulation results described above demonstrate that the absence of trial-level objective feedback does not preclude agents from inferring the competence of other agents. However, systematic deviations in perceptions of advisor competence can

be observed between feedback and feedback-free scenarios when judges' and advisors' opinions are correlated (e.g., due to shared biases or common sources of information). To extend this work, we used agent-based modeling to explore how these effects might play out in more complex, multi-actor situations.

If patterns of trust vary depending on feedback availability, we might expect different macro-level patterns of trust in networks where feedback is available vs. difficult to obtain. Furthermore, the presence of bi-directional information channels between agents (as opposed to judge-advisor systems) may lead to clustering of people sharing similar biases, increasing their polarization (i.e., confidence) and tendency to show herding behavior (Janis, 1972; Turner & Pratkanis, 1998). These effects parallel important observations about social networks, where echo-chambers and recommendation systems can produce clusters of individuals with similar characteristics that are impenetrable to external information (Del Vicario, Bessi, et al., 2016; Jasny et al., 2015; Pariser, 2011; Sunstein, 2001). Consuming within-cluster information more than between-cluster information can in turn lead to spurious consensus effect (Bessi, 2016; Del Vicario, Vivaldo, et al., 2016; Yaniv et al., 2009). Similarly, people in a group are known to have better access to shared (rather than private) information (Lightle, Kagel, & Arkes, 2009; Stasser & Titus, 2003). The present ideas provide a normative account of these otherwise suboptimal effects: According to our findings, in the absence of an objective standard, it may be adaptive to prefer advisors whose opinions agree with our own confidently-held views. However, this strategy can lead to systematic biases in trust when opinions and judgments are non-independent (a common feature of real social networks).

To explore the implications of these ideas, we manipulated agents' access to metacognitive signals in an agent-based simulation, and assessed how feedback availability interacts with cognitive mechanisms to produce emerging network macro-structures (Couzin, Krause, James, Ruxton, & Franks, 2002; Epstein, 2013). We expected that (a) when judgments in the population are independent, agreement-based heuristics are useful in reliably approximating true underlying expertise; (b) when judgments are correlated—e.g., due to the presence of shared biases in the population—agents sharing similar biases will tend

713 to cluster together due to the use of a confidence-agreement heuristic in feedback-poor
714 (but not feedback-rich) scenarios; and (c) the use of internal representations of others’
715 competence to discount advice will be beneficial to agents’ performance when judgments
716 are independent but not when they are correlated.

717 5.1 Model description

718 Our agent-based models simulated the development of trust among agents as they make
719 decisions and revise these decisions on the basis of learning the opinions of others. The
720 decisions are simulated as a series of simple binary choices—is a stimulus from category
721 A or B?—as a generalised case of the specific dot discrimination task used in our exper-
722 iments. Of interest was how perceived competence (“trust”) among agents was affected
723 by the relative quality of their decisions, the availability of feedback and, crucially, the
724 degree to which agents shared biases in their decision processes. Our models formalised
725 these biases as differences in prior (base rate) expectations about the likelihood of A
726 or B being the correct answer, which via Bayesian belief updating (SI §5) will bias the
727 interpretation of incoming information in the decision process (as an analogue of what
728 might be expected to happen in everyday decisions as a function of our political leanings,
729 news readings, musical preferences, etc.).

730 After making a judgment, agents selected one other agent to interact with either at
731 random (random sampling) or proportionally to their trust (biased sampling). Agents
732 then updated their initial judgment, either without discounting advice or by discounting
733 the advisor’s judgment proportionally to their current level of trust in the advisor. After
734 updating their judgments, each agent updated its trust judgments based on the available
735 information about other agents (feedback or *estimated* partner’s accuracy based on the
736 agreement-in-confidence heuristics described above). We assessed the effect of feedback
737 by parametrically manipulating its availability.

738 When learnt trust depends on agreement-based heuristics, levels of trust should track
739 true accuracy when judgments are independent but generate clustering of populations—
740 namely high within-group trust and low between-group trust—when judgment correla-

tions emerge within such populations. We test how network clustering is shaped by the presence or absence of objective feedback and show that bias-specific segregation arises only in the absence of feedback.

Finally, once bias-specific segregation is established, we ask whether such clustering remains stable. In particular, after 500 iterations we allow agents to dynamically change their original bias as a function of experience (Akaishi, Umeda, Nagase, & Sakai, 2014; Zylberberg, Wolpert, & Shadlen, 2018). In the present context, if an agent systematically reports “A” but receives negative feedback, they should reduce their bias by decreasing their prior probability $p(A)$. Similarly, when feedback is absent, an agent who systematically reports “A” but finds themselves, after interacting with other agents, believing that B s are more frequent than expected, should reduce their bias towards A s. Conversely, bias should get stronger if the social contexts reinforces it (although see (Bail et al., 2018)) (see Supplementary Information §5 for details on network clustering computation and bias updating rules).

5.2 Results

We first tested the hypothesis that agreement-based heuristics are adaptive in situations of independent judgments because they reliably track others’ expertise without the need to rely on external forms of feedback. To this end, we set all agents’ initial bias $p(A)$ to 0.50 and the probability of feedback to 0. We then created two sub-populations of agents that orthogonally varied in their overall judgment accuracy (modeled as the degree of perceptual noise, with decreasing accuracy as a function of increasing noise) and calculated the average trust toward a target sub-population (here, Population 2). The results show that average trust in Population 2 agents correctly tracks their underlying perceptual noise: Trust in Population 2 is highest when their perceptual noise is low and decreases as their noise increases (y-gradient in Figure 6A). Interestingly, the effect is not linear and we observe an interaction (x-gradient in Figure 6A) between the accuracy of Population 1 and the accuracy of Population 2 (no interaction is observed in Population 2’s perception of itself, Figure S8), mirroring the simple numerical analysis in Figure

1B: In the absence of feedback, poor performers in Population 1 fail to distinguish the accuracy of others, evident as a narrower range of trust values in rightward pixels in Figure 6A [cf. (Kruger & Dunning, 1999)]. Nevertheless, overall, we see that trust in Population 2 systematically tracks its level of performance, even in the absence of objective feedback.

Actively inferring the competence of others is beneficial for agents. Population 1 agents generally improve their accuracy by receiving advice from Population 2 members (most pixels in Figure 6B are light-colored), but this benefit becomes an accuracy cost if Population 2 agents are much worse at the task than Population 1 agents (dark pixels in the lower left of Figure 6B; cf. Bahrami et al., 2010). However, this cost of receiving bad advice is mitigated when advice is discounted according to trust learnt through confidence-weighted agreement (Figure 6C, light-colored pixels in the lower left corner). Thus, learning differentiated patterns of trust enables agents to benefit from advice when it is useful, and downweight or ignore it otherwise.

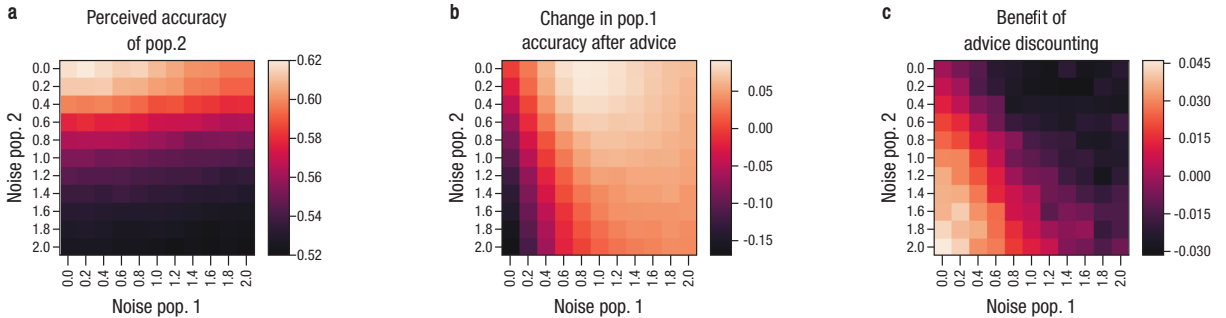


Figure 6: A: average trust shown in agents belonging to population 2. Trust in these agents increases as the noise affecting them decreases. Interestingly, trust formation is affected by perceptual noise of the observed agent(s) as well as perceptual noise of the observing agent(s) [cf.(Kruger & Dunning, 1999)]. B: impact of advice for Population 1’s accuracy. Pixels’ color represent the difference between post and pre-advice accuracy. C: the use of a trust-based advice discounting strategy increases post-interaction accuracy of accurate agents’ but not inaccurate ones, particularly when interacting with low accuracy advisors (bottom left corner).

Our second hypothesis was that when the true state of the world is difficult to discern (e.g., when objective feedback is rare or absent, signal strength is weak or perceptual noise is large), bias-specific segregation can arise. We therefore ran a simulation in which we varied the bias $p(A)$ shown by two subpopulations of equal size and equal average accuracy (equivalent average perceptual noise). The two populations differed in

787 their base rate estimates of the relative likelihood of the two outcomes (A and B), with
 788 Population 1 biased towards judging events as “B”s and Population 2 biased towards
 789 judging events as “A”s (ie., $p(A)$ drawn in the range $[0, .5]$ and $[.5, 1]$ respectively). On
 790 each iteration, agents interacted with another randomly selected agent and updated their
 791 trust via a simple delta rule (Supplementary Information, equation S25). Figure 7 shows
 792 the average clustering—quantified as the degree to which agents learn greater trust in
 793 others who share their initial biases than others who have different initial biases—after
 794 1000 iterations. The data are plotted as a function of signal strength (low in panel A to
 795 high in panel B), perceptual noise (variation along the x-axis of each panel), and feedback
 796 probability (y-axis).

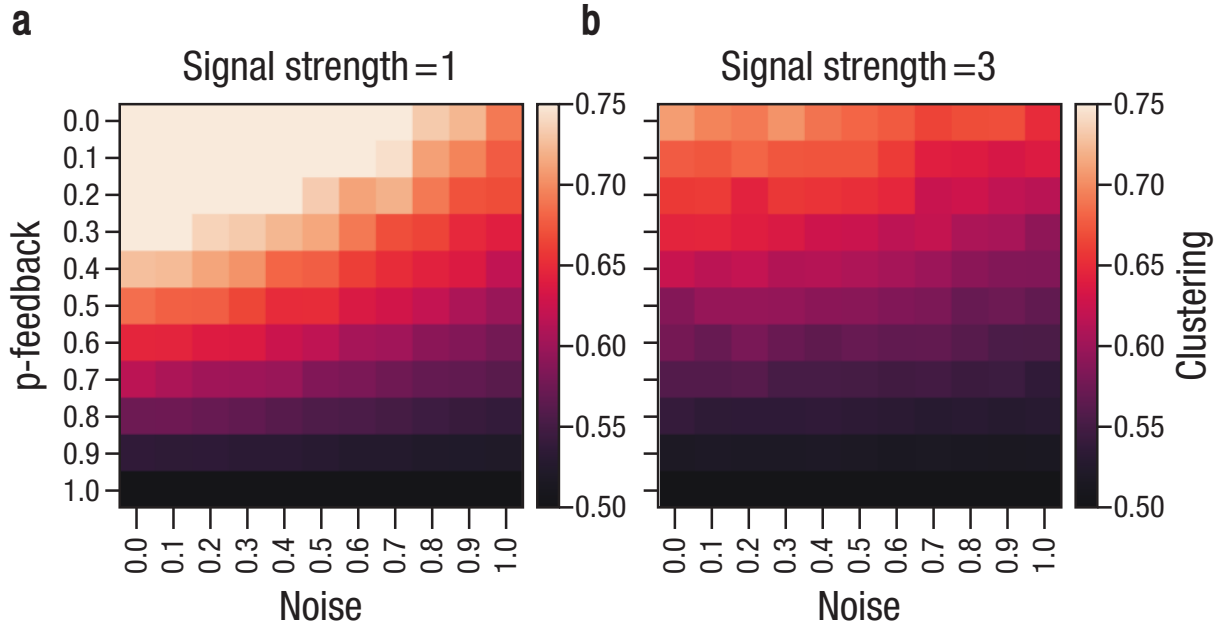


Figure 7: Clustering as a function of signal strength, probability of feedback and noise. Each pixel in each panel indicates the average clustering value (over 20 simulations), computed as described in the Supplementary Information, after 1000 iterations of decision, advice and update in a network of agents with fixed stimulus strength, perceptual noise and feedback probability parameters. A value of 0.5 indicates equal trust in in-group and out-group members. Values greater than 0.5 indicate greater trust in in-group than out-group members. Right and left columns distinguish simulations with and without bias update respectively. The two panels represent increasing signal strengths.

797 Average clustering decreased as the signal strength increased (Figure 7A-B panels),
 798 as perceptual noise increased (gradient over x-axis), and as the probability of receiving
 799 objective feedback increased (gradient over y-axis). In these cases, agents' decisions are

dominated by external evidence rather than their prior expectations (Equation S20) and within-group and between-group agreement rates are similar. On the contrary, when feedback availability, noise or signal strengths decrease, decisions and thus agreement-based trust will more strongly be influenced by prior expectations, as predicted by Bayes theorem. This finding suggests that simple rules of trust update can perform very differently depending on what learning signal is used. In the absence of objective feedback, the circular nature of using one’s own belief to estimate others’ accuracy produces higher clustering and segregation, particularly when decisions are often ambiguous. The presence of noise reduces the risk of getting stuck on local minima, particularly when correlations exist between judges (Couzin et al., 2011; Kao & Couzin, 2014; Shirado & Christakis, 2017).

Our third and final hypothesis was that segregation between in-group and out-group members, once established, can resist modification and in turn shape the evolution of shared beliefs. In this simulation, after 500 iterations we allowed agents to update their bias $p(A)$ via delta rule, whereby their experience of the relative prevalence of A and B outcomes—as determined by their posterior opinion, opinions of partners, and objective feedback, when available—led to modification of their estimates of $p(A)$. Of interest was the way in which the biases of the two groups evolved across interactions as a function of feedback availability, perceptual noise, and information sampling strategy.

When biased agents (differing in prior $p(A)$ and marked by different colors in Figure 8) select partners at random (columns A and B), or when objective feedback is available (saturated lines), decision biases rapidly diminish such that agents from both populations converge on accurate estimates of the relative likelihood of A vs. B outcomes (i.e., $p(A) = 0.5$), independently of noise (Figure 8A,C). Columns B and D plot the distribution of biases across members of each group after the final iteration. However, when feedback is less available (low saturation lines) and trust drives agents’ selection of their advisors (Figure 8C-D), population biases are persistent in a manner that exhibits sensitivity to both information selection strategy and levels of perceptual noise. Cluster segregation is alleviated by greater perceptual noise (S9-10) but magnified further if agents discount

advice received based on trust (Figure S11). Thus, social information can stabilise (and even increase) initial group biases insofar as agents rely on this social information in lieu of objective feedback. Under these circumstances, a positive feedback loop can develop, whereby agents who share an initial bias (of, say, believing that $p(A) > 0.5$) will tend to ask each other for advice more often, thereby compounding each other’s biased beliefs over time, thereby increasing their tendency to seek advice from like-minded agents, thereby compounding their bias, and so on across iterations.

In this way, our agent-based modelling simulations demonstrate how sub-optimal behaviors at the network-level—such as self-sustaining biases and even extremes of polarization—can emerge from learning strategies that are normatively sensible at the level of individual agents: Using decision confidence as a proxy for objective accuracy can allow independent agents to learn about the validity of social information sources when feedback is infrequent or absent. Moreover, using social information can in turn allow agents to learn about the world under these conditions of diminished feedback. However, these learning strategies exhibit important limitations when agents’ initial decisions are not independent but rather exhibit patterns of shared vs. distinct biases across individuals within a population: They lead to predictable patterns of clustering and trust when agents’ agreement rates are inflated by their shared biases or shared access to correlated information. Moreover, bias updating and partner-selection strategies also affect network structure by changing the segregation of individuals sharing the same bias, leading to stabilisation of networks with systematically biased or even increasingly polarized beliefs about the world (in our simulations, the relative likelihood of A vs. B being the correct decision).

6 General Discussion

The present research explored how people learn the accuracy of others’ advice in contexts where objective external feedback is not readily available or costly to acquire. In these scenarios, we hypothesize that metacognitive confidence provides a useful proxy for objective feedback. This is due to the fact that in many tasks, confidence provides a

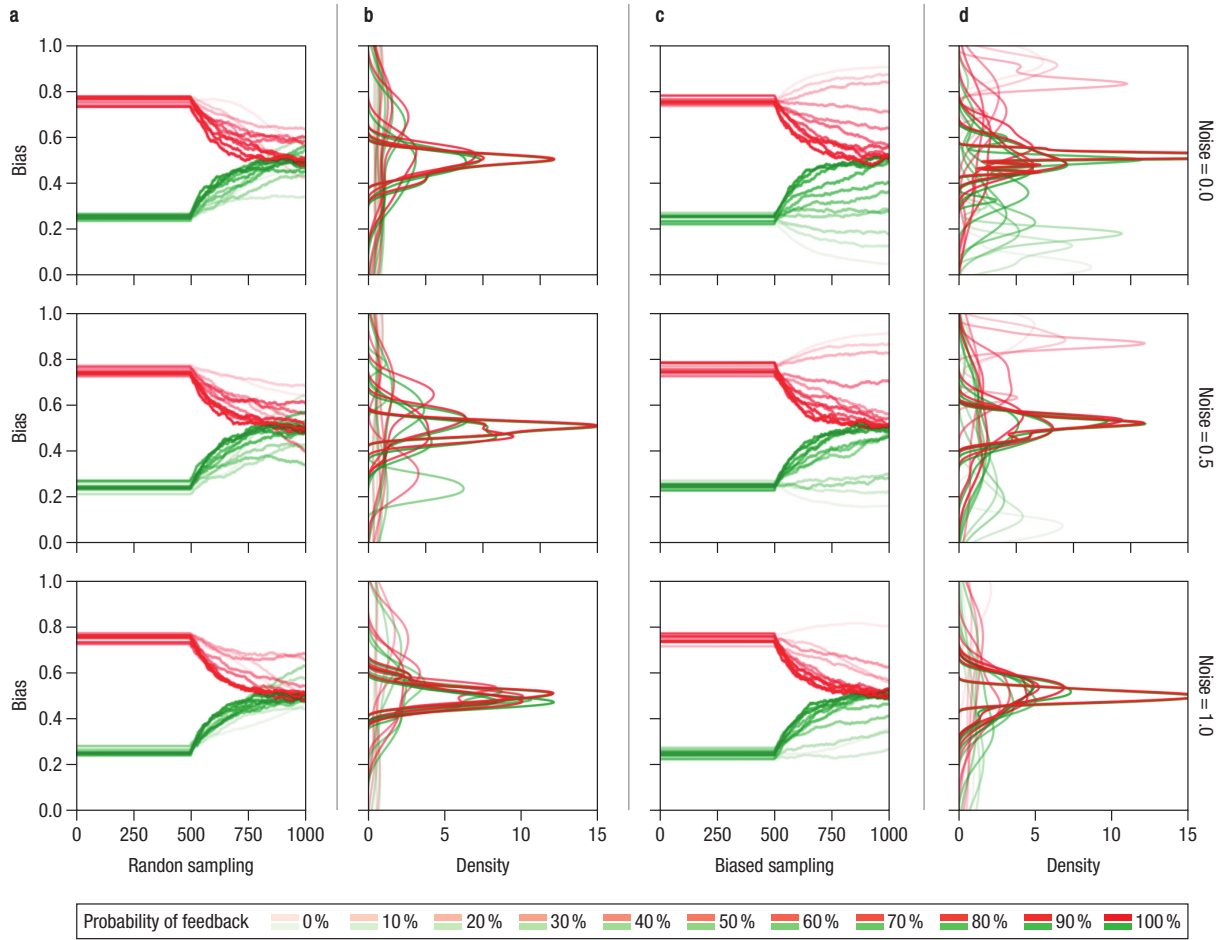


Figure 8: Bias as a function of time, feedback probability and agents’ perceptual noise. Rows represent increasing levels of noise. Lines alpha values (saturation) represent probability of objective feedback. Columns A and B show agents’ bias evolution over time and final bias distribution when agents choose their partners at random. Columns C and D show the same graphs when agents choose their partner proportionally to their trust.

finely calibrated estimate of underlying accuracy (Fleming, 2016; Henmon, 1911; Koriati, 2012; Pescetelli et al., 2016). Thus, people can apply this internal confidence estimate to received advice in order to estimate its accuracy: as $p(\text{correct})$ in the case of agreement, and $1-p(\text{correct})$ in the case of disagreement. This signal, accumulated over time, can support learning of reliable estimates of others’ competence. Experiment 1 empirical data and models show the power of this agreement-in-confidence heuristic in identifying subtle differences in advisor quality when advisors provide new independent information: Participants (and simple models) learned to distinguish advisors of differing accuracy and confidence calibration even in the absence of objective feedback, and did so to a similar extent to participants (and simple models) who had access to objective feedback

after every decision. Thus, we propose that people are able to follow internal signals (e.g., decision confidence) in a comparable manner to external signals—such as rewards or feedback—to learn about the accuracy of advice and advisors according to associative learning principles (Behrens et al., 2008; Sutton & Barto, 1998).

Experiments 2 and 3 tested a key prediction (and a crucial limitation of) the agreement-in-confidence heuristic, in common with an established research tradition exploring heuristics and biases in human judgments (Gigerenzer, 2008; Tversky & Kahneman, 1974): According to this tradition, our limited cognitive capacities mean that we must often adapt to use approximations or “short-cuts” to find good-enough solutions to otherwise intractable problems (Tversky & Kahneman, 1983). However, use of such heuristics can lead to systematic errors when their assumptions are violated. In the current work, the agreement-in-confidence heuristic provides a solution to the challenging problem of estimating the accuracy of information in the absence of an objective standard. This solution works well when agreement correlates with accuracy, as is the case where initial judgment and advice are independent (Experiment 1). However, the process goes astray when the independence between judgment and advice is broken, as in Experiments 2 and 3 where advice was contingent on participants’ initial decision and expressed confidence. Adverse patterns of competence ratings and influence emerged when objective feedback was unavailable. Our findings are in agreement with studies showing that people often tend to ignore the correlation structure of advisors and heavily rely on agreement among information sources (Yaniv et al., 2009).

These findings extend our understanding of advice-taking and trust formation in three main ways. First, they extend our understanding of the uses of metacognitive signals in learning, by showing that learning takes place also in the absence of objective feedback or contextual cues. Previous studies have recognized the adaptive value of metacognitive monitoring in a range of cognitive tasks, like uncertainty monitoring (Yeung & Summerfield, 2012), information seeking (Desender, Boldt, & Yeung, 2018), and cognitive control (Botvinick, Braver, Barch, Carter, & Cohen, 2001). The present findings demonstrate that confidence is not only the end-product of information flow from the external stimulus

to a perceptual inference, but it can feed back to help making inferences about external events, thus complementing objective learning signals (Behrens et al., 2008). Evidence seems to support this notion, as confidence signals have been shown to parallel reward prediction error signals in feedback-free situations (Guggenmos et al., 2016; Zylberberg et al., 2018).

Second, our findings extend understanding of metacognition in the social domain. Specifically, our results go beyond previous evidence that people downweight advice that disagrees with their initial opinion (Lieberman et al., 2012; Yaniv, 2004), to explore how aggregated experiences of agreement and disagreement shape people’s learned evaluations of overall advisor competence. Thus, our effects are observed in terms of patterns of trust and influence that reflect learning across repeated interactions, such that that advisors who more regularly disagree with a participants’ confidently-made choices are less influential, even on those trials where their advice happens to agree with the participants’ view. Conversely, we find that advisors who more regularly agree with the participant, particularly when the participant is confident in their initial choice, are more influential even on those trials where their opinion diverges from the participants’ (Supplementary results, §2).

A long tradition in social psychology has investigated the importance of agreement and confidence in advice taking (Bonaccio & Dalal, 2006; Ecken & Pibernik, 2016; Swol & Snizek, 2005), persuasion (Price & Stone, 2004), influence (Lieberman et al., 2012; Rader et al., 2017; Sah et al., 2013), and group dynamics (Sherif et al., 1965; Stasser & Davis, 1981), and its relevance in applied fields such as judicial systems and organizations (Bovens & Hartmann, 2004; Roediger III et al., 2012; Schum & Martin, 1982; Stasser & Davis, 1981). This literature consistently demonstrates that people tend to pay an accuracy cost for not using advice enough or using it in a self-serving manner (Minson et al., 2011; Soll & Mannes, 2011), for example by discounting advice proportionally to the conflict with one’s own opinion. Here, we suggest that these phenomena might emerge as by-products of a normatively-justified heuristic that tries to overcome the computational intractability of learning without objective feedback by approximating it with *subjective*

924 probabilistic estimates.

925 Our third and final contribution is in identifying the boundary conditions of this
926 best-response strategy. Other studies in the social domain have focused on the way
927 that advisor’s confidence provides a useful signal that can benefit group decision making
928 (Bahrami et al., 2010; Koriat, 2012; Sorkin, Hays, & West, 2001) and social coordination
929 (Shea et al., 2014). Most previous literature has focused on situations where observers are
930 independent. Here, we investigate environments characterized by correlated information
931 among judges, where the heuristics that people use to gauge the usefulness of social
932 signals turn out to be maladaptive. Far from being fringe cases, these environments may
933 be common in many real world scenarios, both in physical and digital space (Del Vicario,
934 Vivaldo, et al., 2016; Kao et al., 2014). Correlation among judges emerge from sharing
935 similar social or information cliques, which in turn lead to sharing similar biases. Our
936 findings (Experiment 2-3, and agent-based simulation) show the micro- and macro-scale
937 effects of using confidence and agreement in feedback-poor information environments.

938 Our results were generally in line with the idea that we trust advisors according to their
939 objective accuracy when feedback is present, and trust agreeing-in-confidence advisors
940 when it is absent. However, there were some notable exceptions. First, advisor agreement
941 rate had an effect over and above accuracy even when feedback was available on every trial
942 in Experiment 2, suggesting that people value agreement even when redundant. Second,
943 people seem to value information over pure accuracy, as shown by participants in the
944 Feedback group of Experiment 3, who preferred the advisor who agreed with them more
945 frequently when they themselves were correct but low in confidence (the anti-bias advisor)
946 over equally accurate but less informative advisors. Finally, participants in Experiment
947 3 did not prefer the neutral to anti-bias advisor when feedback was absent. We explain
948 this result in light of the fact that, without feedback, expertise can still be estimated by
949 comparing the classification variability observed for similar stimuli with the classification
950 variability observed across different stimuli (Weiss & Shanteau, 2003).

Our agent-based models allowed us to explore how these cognitive mechanisms might scale up in larger population interactions. In the models, we observed that when feedback was unavailable and agents' judgments were independent, agreement-based trust formation strategies helped agents trust more accurate advisors. However, when individuals covaried in their signals (e.g., by sharing systematic biases), individuals segregated according to their initial biases. Agents within a homogeneous population were more likely to trust and influence each other than agents belonging to different populations. The polarization of each cluster's average bias was reduced by the presence of random noise, increased signal strength or the presence of objective feedback.

These results provide a potential new understanding of echo-chambers and assortative mixing online (Bollen, Gonçalves, Ruan, & Mao, 2011; Sunstein, 2001), as a by-product of an otherwise adaptive mechanism—use of an agreement-in-confidence heuristic to estimate advice accuracy—when it is difficult to objectively assess the validity of others' opinions, and when these opinions might themselves be subject to systematic shared biases. In these scenarios, random fluctuations in initial bias for one option within a population of equally accurate individuals can spiral out to form densely connected communities, thus effectively modifying the network structure. The results show how simulated agents endowed with realistic cognitive mechanisms can shed light on the emergence of complex patterns at the population level (Epstein, 2013). More broadly, the present research suggests the value of identifying strategies used by individual decision makers—who are likely to rely on imperfect heuristics given necessary limits in the information they access and the cognitive resources they have available to process that information—and exploring how these strategies can influence behavior at the group and network level. The debate around the false consensus effect (Dawes, 1989; Krueger & Clement, 1994; Ross et al., 1977) is a clear example of how understanding biases in terms of their adaptive potential can fostered fruitful debate in psychology. In a similar vein, we hope our work can further understanding of phenomena like echo-chambers and assortative mixing, to design better platforms for democratic debate that are mindful of the heuristics they might elicit in their users.

7 Conclusions

The current work aimed to show that confidence is a valuable attribute of someone’s judgment in social decision-making. It helps others discriminate when one is more likely to be correct (and thus value their contribution) but also helps a decision-maker to make consistent judgments about others irrespective of feedback availability. However, this potentially adaptive solution to an intractable problem of learning in the absence of feedback can backfire when judgments from different observers are not independent, leading to systematic biases in trust and influence.

8 Context of the Research

This research was conducted as part of the D.Phil. research of the first author (N.P.). It was inspired by research on the computational mechanisms of social interaction, e.g., by Bahador Bahrami (with whom the first author completed his MSc dissertation), and the confluence of these ideas with research in N.Y.’s lab concerning the nature and functional role of confidence judgements. Although confidence has for a long time been studied in terms of its role in social and organisational decision making, new developments in cognitive neuroscience place metacognitive signals within the wider context of cognitive control and uncertainty monitoring. Working within this framework, the present research adds to growing evidence — from N.Y.’s lab and elsewhere — indicating that metacognitive signals play an important role in adaptive behavior, both at the level of individual decision makers and at the level of groups of interacting individuals. This research also sits well within the research program of N.P., whose research focuses on social learning and decision-making by human and algorithmic agents. This work is being developed, although in different directions, by both authors in their respective groups: the ACC lab, led by N.Y. at the University of Oxford, and the Hybrid Collective Intelligence group, led by N.P. at the Max Planck Institute for Human Development. The former is investigating how decision confidence drives the collection of new evidence, while the latter develops the current findings to hybrid human-machine collective decisions.

(Pescetelli & Yeung, 2019)

References

- Aitchison, L., Bang, D., Bahrami, B., & Latham, P. (2015). Doubly Bayesian Analysis of Confidence in Perceptual Decision-Making. *PLoS Computational Biology*, 11(10), 1.
- Akaishi, R., Umeda, K., Nagase, A., & Sakai, K. (2014, 1). Autonomous mechanism of internal choice estimate underlies decision inertia. *Neuron*, 81(1), 195–206. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/24333055> doi: 10.1016/j.neuron.2013.10.018
- Bahrami, B., Olsen, K., Latham, P. E., Roepstorff, A., Rees, G., & Frith, C. D. (2010, 8). Optimally interacting minds. *Science (New York, N.Y.)*, 329(5995), 1081–5. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/20798320> doi: 10.1126/science.1185718
- Bail, C. A., Argyle, L. P., Brown, T. W., Bumpus, J. P., Chen, H., Hunzaker, M. B. F., ... Volfovsky, A. (2018, 9). Exposure to opposing views on social media can increase political polarization. *Proceedings of the National Academy of Sciences*, 115(37), 9216–9221. Retrieved from <http://www.pnas.org/lookup/doi/10.1073/pnas.1804840115> doi: 10.1073/pnas.1804840115
- Bakeman, R. (2005). Recommended effect size statistics for repeated measures designs. *Behavior Research Methods*, 37(3), 379–384.
- Behrens, T. E. J., Hunt, L. T., Woolrich, M. W., & Rushworth, M. F. S. (2008, 11). Associative learning of social value. *Nature*, 456(7219), 245–9. Retrieved from <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2605577&tool=pmcentrez&rendertype=abstract> doi: 10.1038/nature07538
- Bessi, A. (2016, 12). Personality traits and echo chambers on facebook. *Computers in Human Behavior*, 65, 319–324. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/S0747563216305817> doi: 10.1016/j.chb.2016.08.016

- Boldt, A., & Yeung, N. (2015, 2). Shared Neural Markers of Decision Confidence and Error Detection. *Journal of Neuroscience*, 35(8), 3478–3484. Retrieved from <http://www.jneurosci.org/cgi/doi/10.1523/JNEUROSCI.0797-14.2015> doi: 10.1523/JNEUROSCI.0797-14.2015
- Bollen, J., Gonçalves, B., Ruan, G., & Mao, H. (2011, 7). Happiness Is Assortative in Online Social Networks. *Artificial Life*, 17(3), 237–251. Retrieved from http://www.mitpressjournals.org/doi/10.1162/artl_a-00034 doi: 10.1162/artl_a-00034
- Bonaccio, S., & Dalal, R. S. (2006, 11). Advice taking and decision-making: An integrative literature review, and implications for the organizational sciences. *Organizational Behavior and Human Decision Processes*, 101(2), 127–151. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/S0749597806000719> doi: 10.1016/j.obhdp.2006.07.001
- Botvinick, M. M., Braver, T. S., Barch, D. M., Carter, C. S., & Cohen, J. D. (2001). Conflict monitoring and cognitive control. *Psychological Review*, 108(3), 624–652. Retrieved from <http://doi.apa.org/getdoi.cfm?doi=10.1037/0033-295X.108.3.624> doi: 10.1037/0033-295X.108.3.624
- Bovens, L., & Hartmann, S. (2004). *Bayesian Epistemology*. Oxford: Oxford University Press.
- Clark, C. C., & Lawless, H. T. (1994). Limiting response alternatives in time-intensity scaling: an examination of the halo-dumping effect. *Chemical Senses*, 19(6), 583–594. Retrieved from <https://ezproxy-prd.bodleian.ox.ac.uk:5876/chemse/article/19/6/583/article> doi: 10.1093/chemse/19.6.583
- Couzin, I. D., Ioannou, C. C., Demirel, G., Gross, T., Torney, C. J., Hartnett, A., ... Leonard, N. E. (2011). Uninformed Individuals Promote Democratic Consensus in Animal Groups. *Science (New York, N.Y.)*, 334(6062), 1578–1580. doi: 10.1126/science.1210280
- Couzin, I. D., Krause, J., James, R., Ruxton, G. D., & Franks, N. R. (2002, 9). Collective Memory and Spatial Sorting in Animal Groups. *Journal of Theoreti-*

cal Biology, 218(1), 1–11. Retrieved from <https://linkinghub.elsevier.com/retrieve/pii/S0022519302930651> doi: 10.1006/jtbi.2002.3065

Dawes, R. M. (1989, 1). Statistical criteria for establishing a truly false consensus effect. *Journal of Experimental Social Psychology*, 25(1), 1–17. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/002210318990036X> doi: 10.1016/0022-1031(89)90036-X

Del Vicario, M., Bessi, A., Zollo, F., Petroni, F., Scala, A., Caldarelli, G., ... Quattrociocchi, W. (2016, 1). The spreading of misinformation online. *Proceedings of the National Academy of Sciences*, 113(3), 554–559. Retrieved from <http://www.pnas.org/lookup/doi/10.1073/pnas.1517441113> doi: 10.1073/pnas.1517441113

Del Vicario, M., Vivaldo, G., Bessi, A., Zollo, F., Scala, A., Caldarelli, G., & Quattrociocchi, W. (2016, 12). Echo Chambers: Emotional Contagion and Group Polarization on Facebook. *Scientific Reports*, 6(1), 37825. Retrieved from <http://www.nature.com/articles/srep37825> doi: 10.1038/srep37825

De Martino, B., O'Doherty, J. P., Ray, D., Bossaerts, P., & Camerer, C. (2013, 9). In the Mind of the Market: Theory of Mind Biases Value Computation during Financial Bubbles. *Neuron*, 79(6), 1222–1231. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/S0896627313005680> doi: 10.1016/j.neuron.2013.07.003

Desender, K., Boldt, A., & Yeung, N. (2018, 5). Subjective Confidence Predicts Information Seeking in Decision Making. *Psychological Science*, 29(5), 761–778. Retrieved from <http://journals.sagepub.com/doi/10.1177/0956797617744771> doi: 10.1177/0956797617744771

Ecken, P., & Pibernik, R. (2016, 7). Hit or Miss: What Leads Experts to Take Advice for Long-Term Judgments? *Management Science*, 62(7), 2002–2021. Retrieved from <http://pubsonline.informs.org/doi/10.1287/mnsc.2015.2219> doi: 10.1287/mnsc.2015.2219

Epstein, J. M. (2013). *Agent_Zero: Toward Neurocognitive Foundations for Generative*

Social Science. Princeton and Oxford: Princeton University Press.

Fleming, S. M. (2016). Changing our minds about changes of mind. *eLife*, 5. doi: 10.7554/eLife.14790

Fleming, S. M., & Daw, N. D. (2017). Self-evaluation of decision performance: A general Bayesian framework for metacognitive computation. *Psychological review*, 124(1), 1–59. doi: 10.1002/sml.))

Fleming, S. M., & Lau, H. C. (2014, 7). How to measure metacognition. *Frontiers in Human Neuroscience*, 8. Retrieved from http://www.frontiersin.org/Human_Neuroscience/10.3389/fnhum.2014.00443/abstract doi: 10.3389/fnhum.2014.00443

Gigerenzer, G. (2008). Why Heuristics Work? *Perspectives on Psychological Science*, 3(1), 20–29.

Gigerenzer, G., & Selten, R. (2002). *Bounded Rationality: The Adaptive Toolbox*. Cambridge, MA: MIT Press.

Guggenmos, M., & Sterzer, P. (2017). A confidence-based reinforcement learning model for perceptual learning. *bioRxiv*.

Guggenmos, M., Wilbertz, G., Hebart, M. N., & Sterzer, P. (2016, 3). Mesolimbic confidence signals guide perceptual learning in the absence of external feedback. *eLife*, 5. Retrieved from <http://elifesciences.org/lookup/doi/10.7554/eLife.13388> doi: 10.7554/eLife.13388

Henmon, V. A. C. (1911). The relation of the time of a judgment to its accuracy. *Psychological Review*, 18(3), 186–201. Retrieved from <http://doi.apa.org/getdoi.cfm?doi=10.1037/h0074579> doi: 10.1037/h0074579

Hertz, U., Romand-Monnier, M., Kyriakopoulou, K., & Bahrami, B. (2016). Social influence protects collective decision making from equality bias. *Journal of Experimental Psychology: Human Perception and Performance*, 42(2), 164–172. Retrieved from <http://doi.apa.org/getdoi.cfm?doi=10.1037/xhp0000145> doi: 10.1037/xhp0000145

Jamieson, K. H., & Cappella, J. N. (2008). *Echo Chamber: Rush Limbaugh and the*

Conservative Media Establishment. Oxford: Oxford University Press. Retrieved from https://books.google.co.uk/books?id=1390a4M0sAgC&redir_esc=y

Janis, I. L. (1972). *Victims of groupthink*. Boston: Houghton Mifflin.

Jasny, L., Waggle, J., & Fisher, D. R. (2015, 5). An empirical examination of echo chambers in US climate policy networks. *Nature Climate Change*, 5(8), 782–786. Retrieved from <http://www.nature.com/doifinder/10.1038/nclimate2666> doi: 10.1038/nclimate2666

Kao, A. B., & Couzin, I. D. (2014, 4). Decision accuracy in complex environments is often maximized by small group sizes. *Proceedings of the Royal Society B: Biological Sciences*, 281(1784), 20133305–20133305. Retrieved from <http://rspb.royalsocietypublishing.org/cgi/doi/10.1098/rspb.2013.3305> doi: 10.1098/rspb.2013.3305

Kao, A. B., Miller, N., Torney, C., Hartnett, A., & Couzin, I. D. (2014, 8). Collective Learning and Optimal Consensus Decisions in Social Animal Groups. *PLoS Computational Biology*, 10(8), e1003762. Retrieved from <http://dx.plos.org/10.1371/journal.pcbi.1003762> doi: 10.1371/journal.pcbi.1003762

Koriat, A. (2012, 4). When are two heads better than one and why? *Science (New York, N.Y.)*, 336(6079), 360–2. Retrieved from <http://www.sciencemag.org/cgi/doi/10.1126/science.1216549> <http://www.ncbi.nlm.nih.gov/pubmed/22517862> doi: 10.1126/science.1216549

Krause, J., Ruxton, G. D., & Krause, S. (2010, 1). Swarm intelligence in animals and humans. *Trends in Ecology & Evolution*, 25(1), 28–34. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/S0169534709002298> doi: 10.1016/j.tree.2009.06.016

Krueger, J., & Clement, R. W. (1994). The truly false consensus effect: An ineradicable and egocentric bias in social perception. *Journal of Personality and Social Psychology*, 67(4), 596–610. Retrieved from <http://doi.apa.org/getdoi.cfm?doi=10.1037/0022-3514.67.4.596> doi: 10.1037/0022-3514.67.4.596

Kruger, J., & Dunning, D. (1999). Unskilled and Unaware of It: How Difficulties in

1150 Recognizing One's Own Incompetence Lead to Inflated Self-Assessments. *Journal*
1151 *of personality and social psychology*, 77(6), 1121–1134.

1152 Le Bon, G. (1895). *La psychologie des foules* (Félix Alca ed.). Paris: Ancienne Li-
1153 brairie Germer Bailliere et Cie. Retrieved from [http://socserv.mcmaster.ca/](http://socserv.mcmaster.ca/econ/ugcm/3113/lebon/Crowds.pdf)
1154 [econ/ugcm/3113/lebon/Crowds.pdf](http://socserv.mcmaster.ca/econ/ugcm/3113/lebon/Crowds.pdf)

1155 Liberman, V., Minson, J. A., Bryan, C. J., & Ross, L. (2012, 3). Naïve realism
1156 and capturing the wisdom of dyads. *Journal of Experimental Social Psychology*,
1157 48(2), 507–512. Retrieved from [https://linkinghub.elsevier.com/retrieve/](https://linkinghub.elsevier.com/retrieve/pii/S0022103111002599)
1158 [pii/S0022103111002599](https://linkinghub.elsevier.com/retrieve/pii/S0022103111002599) doi: 10.1016/j.jesp.2011.10.016

1159 Lightle, J. P., Kagel, J. H., & Arkes, H. R. (2009, 4). Information Exchange in Group
1160 Decision Making: The Hidden Profile Problem Reconsidered. *Management Science*,
1161 55(4), 568–581. Retrieved from [http://pubsonline.informs.org/doi/abs/10](http://pubsonline.informs.org/doi/abs/10.1287/mnsc.1080.0975)
1162 [.1287/mnsc.1080.0975](http://pubsonline.informs.org/doi/abs/10.1287/mnsc.1080.0975) doi: 10.1287/mnsc.1080.0975

1163 Mackay, C. (1841). *Extraordinary Popular Delusions and the Madness of Crowds*
1164 (Wordsworth ed.). Ware, UK: Wordsworth Edition Limited.

1165 Mahmoodi, A., Bang, D., Olsen, K., Zhao, Y. A., Shi, Z., Broberg, K., ... Bahrami, B.
1166 (2015). Equality bias impairs collective decision-making across cultures. *Proceedings*
1167 *of the National Academy of Sciences*, 201421692. Retrieved from [http://www.pnas](http://www.pnas.org/lookup/doi/10.1073/pnas.1421692112)
1168 [.org/lookup/doi/10.1073/pnas.1421692112](http://www.pnas.org/lookup/doi/10.1073/pnas.1421692112) doi: 10.1073/pnas.1421692112

1169 Mayer, R. C., Davis, J. H., & Schoorman, F. D. (1995, 7). An Integrative Model of
1170 Organizational Trust. *The Academy of Management Review*, 20(3), 709. Retrieved
1171 from <http://www.jstor.org/stable/258792?origin=crossref> doi: 10.2307/
1172 258792

1173 Meyniel, F., Sigman, M., & Mainen, Z. (2015). Confidence as Bayesian Probability:
1174 From Neural Origins to Behavior. *Neuron*, 88(1), 78–92. Retrieved from [http://](http://linkinghub.elsevier.com/retrieve/pii/S0896627315008284)
1175 linkinghub.elsevier.com/retrieve/pii/S0896627315008284 doi: 10.1016/j
1176 [.neuron.2015.09.039](http://linkinghub.elsevier.com/retrieve/pii/S0896627315008284)

1177 Minson, J. A., Liberman, V., & Ross, L. (2011, 10). Two to Tango. *Per-*
1178 *sonality and Social Psychology Bulletin*, 37(10), 1325–1338. Retrieved from

<http://journals.sagepub.com/doi/10.1177/0146167211410436> doi: 10.1177/
0146167211410436

Myers, D. G., & Lamm, H. (1976). The group polarization phenomenon. *Psychological Bulletin*, 83(4), 602–627. doi: 10.1037/0033-2909.83.4.602

Pariser, E. (2011). *The Filter Bubble: What The Internet Is Hiding From You*. London: Penguin.

Penrod, S., & Cutler, B. (1995). Witness confidence and witness accuracy: Assessing their forensic relation. *Psychology, Public Policy, & Law*, 1, 817–845.

Pescetelli, N., Hauperich, A.-K., & Yeung, N. (2019). *Confidence Drives Post-decisional Search of Information from Social Sources*.

Pescetelli, N., Rees, G., & Bahrami, B. (2016). The perceptual and social components of metacognition. *Journal of Experimental Psychology: General*, 145(8), 949–965. doi: 10.1037/xge0000180

Pescetelli, N., & Yeung, N. (2019). *The role of decision confidence in advice-taking and trust formation*. Retrieved from osf.io/phxcs

Pouget, A., Drugowitsch, J., & Kepecs, A. (2016, 2). Confidence and certainty: distinct probabilistic quantities for different goals. *Nature Neuroscience*, 19(3), 366–374. Retrieved from <http://www.nature.com/doifinder/10.1038/nn.4240> doi: 10.1038/nn.4240

Price, P. C., & Stone, E. R. (2004, 1). Intuitive evaluation of likelihood judgment producers: evidence for a confidence heuristic. *Journal of Behavioral Decision Making*, 17(1), 39–57. Retrieved from <http://doi.wiley.com/10.1002/bdm.460> doi: 10.1002/bdm.460

Pulford, B. D., Colman, A. M., Buabang, E. K., & Krockow, E. M. (2018, 10). The persuasive power of knowledge: Testing the confidence heuristic. *Journal of Experimental Psychology: General*, 147(10), 1431–1444. Retrieved from <http://doi.apa.org/getdoi.cfm?doi=10.1037/xge0000471> doi: 10.1037/xge0000471

Rader, C. A., Larrick, R. P., & Soll, J. B. (2017, 8). Advice as a form of social influence: Informational motives and the consequences for accuracy. *Social and Personality*

Psychology Compass, 11(8), e12329. Retrieved from <http://doi.wiley.com/10.1111/spc3.12329> doi: 10.1111/spc3.12329

Roediger III, H. L., Wixted, J. H., & Desoto, K. A. (2012, 7). The Curious Complexity between Confidence and Accuracy in Reports from Memory. In *Memory and law* (pp. 84–117). Oxford University Press. Retrieved from <http://www.oxfordscholarship.com/view/10.1093/acprof:oso/9780199920754.001.0001/acprof-9780199920754-chapter-4> doi: 10.1093/acprof:oso/9780199920754.003.0004

Ross, L., Greene, D., & House, P. (1977, 5). The false consensus effect: An egocentric bias in social perception and attribution processes. *Journal of Experimental Social Psychology*, 13(3), 279–301. Retrieved from <https://linkinghub.elsevier.com/retrieve/pii/002210317790049X> doi: 10.1016/0022-1031(77)90049-X

Sah, S., Moore, D. A., & Maccoun, R. J. (2013). Cheap talk and credibility: The consequences of confidence and accuracy on advisor credibility and persuasiveness. *Organizational Behavior and Human Decision Processes*, 121, 246–255. doi: 10.1016/j.obhdp.2013.02.001

Schultze, T., Gerlach, T. M., & Rittich, J. C. (2018, 7). Some People Heed Advice Less than Others: Agency (but Not Communion) Predicts Advice Taking. *Journal of Behavioral Decision Making*, 31(3), 430–445. Retrieved from <http://doi.wiley.com/10.1002/bdm.2065> doi: 10.1002/bdm.2065

Schum, D. A., & Martin, A. W. (1982). Formal and Empirical Research on Cascaded Inference in Jurisprudence. *Law & Society Review*, 17(1), 105–152. Retrieved from <https://www.jstor.org/stable/3053534?origin=crossref> doi: 10.2307/3053534

Shea, N., Boldt, A., Bang, D., Yeung, N., Heyes, C., & Frith, C. D. (2014, 2). Supra-personal cognitive control and metacognition. *Trends in Cognitive Sciences*. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/S1364661314000230> doi: 10.1016/j.tics.2014.01.006

Sherif, C., Sherif, M., & Nebergall, R. (1965). *Attitude and attitude change*. Philadelphia:

W.B. Saunders Company.

- Shirado, H., & Christakis, N. A. (2017, 5). Locally noisy autonomous agents improve global human coordination in network experiments. *Nature*, 545(7654), 370–374. Retrieved from <http://www.nature.com/doifinder/10.1038/nature22332> doi: 10.1038/nature22332
- Shultz, T. R., Katz, J. A., & Lepper, M. R. (2001). Clinging to beliefs: A constraint-satisfaction model. In *Proceedings of the twenty-third annual conference of the cognitive science society* (pp. 928–933). Mahwah, NJ: Erlbaum.
- Simon, H. A. (1972). Theories of bounded rationality. *Decision and organization*, 1(1), 161–176.
- Snizek, J. A., & Buckley, T. (1989). Social influence in the advisor-judge relationship. In *Annual meeting of the judgment and decision making society*. Atlanta, Georgia.
- Snizek, J. A., & Van Swol, L. M. (2001, 3). Trust, Confidence, and Expertise in a Judge-Advisor System. *Organizational behavior and human decision processes*, 84(2), 288–307. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/11277673> doi: 10.1006/obhd.2000.2926
- Soll, J. B., & Larrick, R. P. (2009, 5). Strategies for revising judgment: how (and how well) people use others' opinions. *Journal of experimental psychology: Learning, memory, and cognition*, 35(3), 780–805. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/19379049> doi: 10.1037/a0015145
- Soll, J. B., & Mannes, A. E. (2011, 1). Judgmental aggregation strategies depend on whether the self is involved. *International Journal of Forecasting*, 27(1), 81–102. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/S0169207010000877> doi: 10.1016/j.ijforecast.2010.05.003
- Sorkin, R. D., Hays, C. J., & West, R. (2001). Signal-detection analysis of group decision making. *Psychological Review*, 108(1), 183–203. Retrieved from <http://doi.apa.org/getdoi.cfm?doi=10.1037/0033-295X.108.1.183> doi: 10.1037/0033-295X.108.1.183
- Stasser, G., & Davis, J. H. (1981). Group decision making and social influence: A social

interaction sequence model. *Psychological Review*, 88(6), 523–551. Retrieved from
<http://content.apa.org/journals/rev/88/6/523> doi: 10.1037/0033-295X.88
.6.523

Stasser, G., & Titus, W. (2003). Hidden Profiles : A Brief History. *Psychological Inquiry*,
14(3), 304–313.

Sunstein, C. R. (2001). *Republic.com*. Princeton University Press.

Sutton, R. S., & Barto, A. G. (1998). *Reinforcement Learning: an introduction*. Cam-
bridge, MA: MIT Press.

Swol, L. M., & Sniezek, J. A. (2005). Factors affecting the acceptance of expert advice.
British journal of social psychology, 44(3), 443–461.

Tenney, E. R., MacCoun, R. J., Spellman, B. a., & Hastie, R. (2007, 1). Calibration
trumps confidence as a basis for witness credibility. *Psychological Science*, 18(1),
46–50. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/17362377> doi:
10.1111/j.1467-9280.2007.01847.x

Tost, L. P., Gino, F., & Larrick, R. P. (2012, 1). Power, competitiveness, and advice tak-
ing: Why the powerful dont listen. *Organizational Behavior and Human Decision*
Processes, 117(1), 53–65. Retrieved from [http://linkinghub.elsevier.com/](http://linkinghub.elsevier.com/retrieve/pii/S0749597811001233)
[retrieve/pii/S0749597811001233](http://linkinghub.elsevier.com/retrieve/pii/S0749597811001233) doi: 10.1016/j.obhdp.2011.10.001

Treutwein, B. (1995). Adaptive Psychophysical Procedures. *Vision Research*, 35(17),
2503–2522.

Turner, M. E., & Pratkanis, A. R. (1998, 2). Twenty-Five Years of Groupthink Theory
and Research: Lessons from the Evaluation of a Theory. *Organizational Behav-*
ior and Human Decision Processes, 73(2-3), 105–115. Retrieved from [http://](http://linkinghub.elsevier.com/retrieve/pii/S074959789892756X)
linkinghub.elsevier.com/retrieve/pii/S074959789892756X doi: 10.1006/
obhd.1998.2756

Tversky, A., & Kahneman, D. (1974). Judgment under Uncertainty: Heuristics and
Biases. *Science*, 185(4157), 1124–1131. Retrieved from [http://www.jstor.org/](http://www.jstor.org/stable/1738360)
[stable/1738360](http://www.jstor.org/stable/1738360)Copy

Tversky, A., & Kahneman, D. (1983). Extensional versus intuitive reasoning: The

conjunction fallacy in probability judgment. *Psychological Review*, 90(4), 293–315. Retrieved from <http://content.apa.org/journals/rev/90/4/293> doi: 10.1037/0033-295X.90.4.293

Vandormael, H., Herce Castañón, S., Balaguer, J., Li, V., & Summerfield, C. (2017, 3). Robust sampling of decision information during perceptual choice. *Proceedings of the National Academy of Sciences*, 114(10), 2771–2776. Retrieved from <http://www.pnas.org/lookup/doi/10.1073/pnas.1613950114> doi: 10.1073/pnas.1613950114

Weiss, D. J., & Shanteau, J. (2003, 1). Empirical assessment of expertise. *Human factors*, 45(1), 104–16. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/12916584>

Yaniv, I. (2004, 1). Receiving other peoples advice: Influence and benefit. *Organizational Behavior and Human Decision Processes*, 93(1), 1–13. Retrieved from <http://linkinghub.elsevier.com/retrieve/pii/S0749597803001018> doi: 10.1016/j.obhdp.2003.08.002

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(2), 558–563. doi: <http://dx.doi.org/10.1037/a0014589>

Yaniv, I., & Kleinberger, E. (2000, 11). Advice Taking in Decision Making: Egocentric Discounting and Reputation Formation. *Organizational behavior and human decision processes*, 83(2), 260–281. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/11056071> doi: 10.1006/obhd.2000.2909

Yeung, N., & Summerfield, C. (2012, 5). Metacognition in human decision-making: confidence and error monitoring. *Philosophical transactions of the Royal Society of London. Series B, Biological sciences*, 367(1594), 1310–21. Retrieved from <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3318764&tool=pmcentrez&rendertype=abstract> doi: 10.1098/rstb.2011.0416

Zarnoth, P., & Snizek, J. A. (1997, 7). The Social Influence of Confidence in

1324 Group Decision Making. *Journal of Experimental Social Psychology*, 33(4),
1325 345–366. Retrieved from [http://linkinghub.elsevier.com/retrieve/pii/](http://linkinghub.elsevier.com/retrieve/pii/S0022103197913263)
1326 S0022103197913263 doi: 10.1006/jesp.1997.1326

1327 Zylberberg, A., Wolpert, D. M., & Shadlen, M. N. (2018, 9). Counterfactual Rea-
1328 soning Underlies the Learning of Priors in Decision Making. *Neuron*, 99(5),
1329 1083–1097. Retrieved from [https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S0896627318306330)
1330 S0896627318306330 doi: 10.1016/j.neuron.2018.07.035

Supplementary Information

Niccolò Pescetelli^{1, 2} and Nick Yeung¹

¹Department of Experimental Psychology, University of Oxford

²Max Planck Institute for Human Development, Berlin, Germany

July 9, 2020

1 Supplementary methods

1.1 Experiment 1: Advisors description

After the participant confirmed their initial perceptual decision response with the spacebar, one of four different advisors appeared centrally as a head-shot picture. The advisors pictures were all Caucasian, smiling female characters (Tottenham et al., 2009), randomly assigned per participant to the four accuracy/calibration conditions described below. Advice was provided in the form of spoken sentences (2 s long), that expressed a binary level of confidence (low vs. high) and either agreement or disagreement with the participant’s judgment. Low confidence was expressed by the sentences “I think it was on the [LEFT/RIGHT]” and “It was on the [LEFT/RIGHT], I think”, with one of the two versions randomly assigned on every trial. Similarly, high confidence was expressed by the sentences “I’m sure it was on the [LEFT/RIGHT]!” and “It was on the [LEFT/RIGHT], I’m sure!”. The use of two inverted sentences for each confidence cue (“I’m sure” vs. “I think”) was to avoid over-repetition of a single sentence and to balance the differences in emphasis that the English language conveys when using the confidence cues at the beginning or end of the sentence. The selection of LEFT or RIGHT depended on the advisor’s choice and accuracy as described below. The spoken advice was pre-recorded from four female native English speakers, again randomised to conditions across participants.

Advisor calibration was defined as the strength of co-variation between confidence judgments and accuracy, and quantified as Type 2 A_{ROC} (A''_{ROC}), a method which that not make assumptions about the generative model of confidence (Fleming & Lau, 2014). Uncalibrated advisors both had an A''_{ROC} of 0.5, meaning that confidence was totally uninformative in predicting the advisor’s trial-level accuracy. Due to the experimental design—in which overall accuracy of advisors was fixed at 60% or 80%, and calibrated advisors were always correct when high in confidence—calibrated advisors differed in their metacognitive sensitivity according to this metric.

1.2 Advice Value as Information Gain

Experiment 1 We formalised an advisor’s informational value as the mean absolute information gained after each possible social encounter with a specific advisor. Informa-

tion gain is the difference between the posterior and prior probability of participant’s correct response:

$$IG = p(d = w|e) - p(d = w) \quad (1)$$

where posterior probability correct $p(d = w|e)$ represents the probability that the decision d is equal to the correct decision w , conditional on the specific social encounter e . Social encounter e represents one of the four possible events: the advisor (1) confidently disagrees, (2) unconfidently disagrees, (3) unconfidently agrees, (4) confidently agrees (where “confidently” and “unconfidently” refer to the level of confidence expressed by the advisor on that trial). Posterior probability $p(d = w|e)$ was computed using Bayes’ theorem and was proportional to participant’s prior probability correct $p(d = w)$ and the likelihood of the social event given participant’s accuracy $p(e|d = w)$. Given the staircase procedure, we used 70% as prior $p(d = w)$. The probability of agreement (or disagreement) conditional on correct response and the overall probability of agreement (or disagreement) were known by design. The mean absolute information gain so computed was lowest for the Inaccurate Uncalibrated advisor (0.08), intermediate for the Accurate Calibrated and Accurate Uncalibrated advisors (0.29 and 0.26 respectively) and highest for the Calibrated but Inaccurate advisor (0.38). This can be intuitively understood by looking at Table 1, in the main text. Although the Inaccurate Calibrated advisor’s accuracy rate is lower than the Accurate Calibrated advisor, outcomes can be better predicted by its judgment. In particular, its judgments correlate strongly positively when sure and strongly negatively when unsure with the correct answer. On the contrary when the Accurate Calibrated advisor is unsure there is a much higher uncertainty about the final outcome. We also computed an expected information gain IG_e for each advisor (Table 1 in the main text) by scaling IG by the overall probability of each event:

$$IG_e = IG * p(e) \quad (2)$$

where $p(e)$ is the overall probability of each social event (i.e., confident disagreement, unconfident disagreement, unconfident agreement, confident agreement). The expected information gain captures the idea that extremely informative but very unlikely events are not very valuable. IG_e values for each advisor were: Accurate Calibrated = .063, Accurate Uncalibrated = .063, Inaccurate Calibrated = .084, Inaccurate Uncalibrated = .021; suggesting that the Inaccurate Calibrated advisor’s advice was the most informative.

Experiment 2 Similar to Experiment 1, we used conditional probabilities and the participants’ expected accuracy to compute the informational value of each advisor. Advisors’ mean absolute information gain IG and expected information gain IG_e were computed as in the previous experiment. Contrary to Experiment 1, however, advisors did not express different levels of confidence. This created only two possible social situations e on each trial (instead of four as in the previous experiment), namely either agreement or disagreement. Information gain was highest for the accurate advisors (.28 and .27 for the high and low agreement advisors respectively) and the lowest for inaccurate advisors (.03 and .06 for the high and low agreement advisors respectively).

Experiment 3 The intuition that the anti-bias advisor was the most informative of the three advisors designed for Experiment 3 was confirmed using a numerical simulation.

We used numerical simulation rather than analytic calculations here because, in Experiment 3, the profile of the advisors cannot be calculated *a priori* but is dependent on the specific distribution of confidence of the participant. The simulations were based on an ideal Bayesian observer performing the task, with a Gaussian distribution of confidence centred on 25 and with a standard deviation of 10. For each initial confidence judgment, the information gained from observing agreement or disagreement was computed for each advisor as the difference between posterior confidence and prior confidence. Contrary to previous Experiments, prior probability correct was here defined on a trial level based on pre-advice confidence. An expected information gain was computed by multiplying the information gain so obtained by the normalisation term in the Bayes formula. This produced a curve of expected information gains after agreeing or disagreeing with each advisor over possible pre-advice confidence levels (Figure S1). The average area under the expected information gain curve was taken as an objective measure of advisor informativeness. Average areas under the curve were 14.74 for the unbiased advisor, 14.63 for the bias-sharing advisor and 15.78 for the anti-bias advisor. This procedure thus quantified and confirmed the intuition that the anti-bias advisor provided the most informative advice.

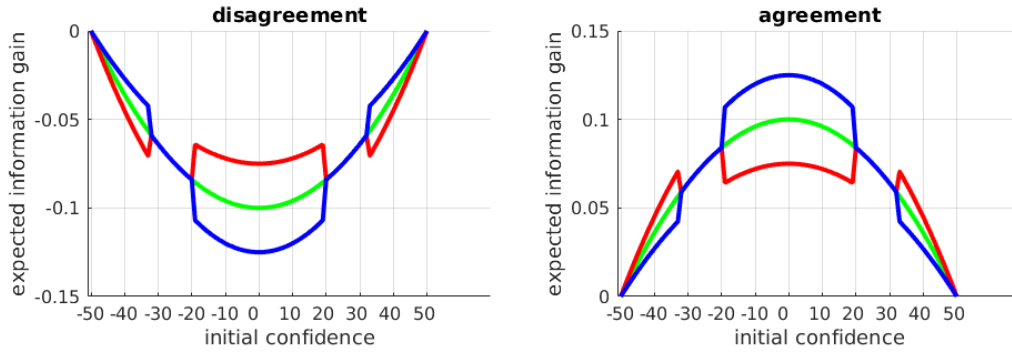


Figure S1: Information gained after agreeing and disagreeing with each advisor type (color code) for each initial subjective prior confidence. Information gain is scaled by the likelihood of agreement and disagreement events. Advice informativeness can be quantified by the area under the curve.

Experiment 1 Advisors				
	Accurate Calibrated	Accurate Uncalibrated	Inaccurate Calibrated	Inaccurate Uncalibrated
A''_{ROC}	.72	.5	0.84	.5
IG	0.29	0.26	0.38	0.08
IG _e	0.063	0.063	0.084	0.021
Experiment 2 Advisors				
	High Accuracy High Agreement	High Acc. Low Agr.	Low Acc. High Agr.	Low Acc. Low Agr.
IG	0.28	0.27	0.03	0.06
IG _e	0.09	0.13	0.01	0.03
Advisors				
	Bias-sharing	Unbiased	Anti-bias	
$AUC(IG_e)$	14.63	14.74	15.78	

Table S1: Experiment 1-3 information gain and expected information gain— IG and IG_e respectively—indicate average informational value of the advice, computed as information gain and expected information gain respectively. Experiment 1 also shows advisors’ calibration, measured as type II AROC.

1.3 Measures of interest

Two measures of estimated advice reliability were defined. The first was the explicit trust that participants expressed in the advisors as collected by the brief questionnaires presented to participants every two blocks. In Experiment 1, four questions asked participants to directly rate on a scale from 1 (“Not at all”) to 50 (“Extremely”) how much they thought each advisor was accurate (Q1), confident (Q2a), trustworthy (Q3) and influential on their own choices (Q4). In Experiments 2 and 3, due to the absence of a confidence judgment from advisors, question 2 was replaced with a question asking about how much participants liked each advisor (Q2b: likeability question). For all Experiments, the first questionnaire was presented immediately after the practice blocks but before any interaction with the advisors took place so to provide a baseline measure. Baseline ratings were subtracted from following ratings to account for confounding factors related to advisors’ appearance and inter-individual differences in the use of the scale. A principal component analysis (PCA) was performed for dimensionality reduction on normalised difference scores and the first component was taken as a unitary measure of expressed trust.

In Experiment 1, question loadings for the Feedback condition were 0.52 (Q1), 0.44 (Q2a), 0.50 (Q3), 0.52 (Q4); and for the No-Feedback condition were 0.51 (Q1), 0.39 (Q2), 0.54 (Q3), 0.52 (Q4).

In Experiment 2, question loadings for the Feedback condition were 0.53 (Q1), 0.39 (Q2b), 0.53 (Q3), 0.52 (Q4); and for the No-Feedback condition were 0.51 (Q1), 0.41 (Q2), 0.53 (Q3), 0.51 (Q4).

In Experiment 3, question loadings for the Feedback condition were 0.53 (Q1), 0.42 (Q2b), 0.53 (Q3), 0.50 (Q4); and for the No-Feedback condition were 0.48 (Q1), 0.47

(Q2), 0.53 (Q3), 0.50 (Q4).

The second measure of interest was an implicit index of advisor’s influence on participant’s opinions, quantifying participants’ confidence change from pre- to post-advice:

$$\delta_C = C_{post} - C_{pre}. \quad (3)$$

where C_{pre} is, for Experiment 1 and 2, an integer value between +1 and +5 and C_{post} is an integer value between -5 and +5 (negative C_{post} representing changes of mind). Positive δ_C values mean increases in confidence from pre- to post-advice and negative values represent decreases in confidence. Notice that δ_C values have a negative skew, ranging from -10 (moving from highest confidence in one judgment to highest confidence in the opposite judgment) to +4 (moving from lowest to highest confidence rating for a single judgment). In Experiment 3, given the difference scale used, C_{post} can assume values between -50 and 50, while δ_C can range from -100 to 49. Agreement and disagreement trials typically have opposite effects on confidence change: agreement usually leads to increases in confidence while disagreement to confidence decreases. The absolute magnitude of confidence shifts in both agreement and disagreement trials can be expected to grow larger as the participant makes more use of the advice received. Thus a unitary measure of influence was obtained by subtracting average δ_C in disagreement from average δ_C in agreement:

$$I = \bar{\delta}_C^a - \bar{\delta}_C^d \quad (4)$$

where I assumes greater values as participant’s confidence increases in agreement and confidence decreases in disagreement become larger.

2 Supplementary Results

2.1 Confidence change by agreement

Previous research has shown that confidence inversely predicts advice taking and that people discount disagreeing advice. The crucial aim of our investigation is to show that confidence and agreement not only influence how each piece of agreeing/disagreeing advice is weighted, but are also aggregated across interactions to discern something about the advisors themselves (i.e., their overall reliability). This aim was reflected in the definition of our influence measure, which assesses pre- to post-advice changes in confidence separately for trials with agreeing vs. disagreeing advice from each advisor, and calculates influence as the difference between these changes. Thus, via our influence measure, we can show that advisors who more regularly disagree with a participants’ choices are less influential, even on those trials where their advice happens to agree with the participants’ view. Conversely, we can show that advisors who more regularly agree with the participant are more influential on those trials where their opinion diverges from the participants’ initial choice. To further explore these effects, and provide further evidence that the observed influence differences reflect gradually-learned evaluations of advisors’ relative reliability, we repeated the analysis and plots reported in the main text, but now assessing confidence change as a function of agreement, separately for each advisor. If people are simply discounting disagreeing advice (even if proportionally to their initial confidence) we would not expect to see differences across advisors, because all agreement

(/disagreement) trials should count the same, independently of who is agreeing. Instead if people are forming a durable representation of others' competence, we would expect to observe differences in confidence change when people agree (disagree) with different advisors. In particular, we would expect a greater influence (i.e., confidence change) of advisors that are believed to be more more competent.

2.2 Experiment 1

We replicate the results observed for the aggregate influence measure (Equation 3) also for agreement and disagreement trials separately. Figure S2 shows the confidence change observed for each advisor, broken down by agreeing and disagreeing trials. As expected, accurate and calibrated advisors produced larger confidence changes in both agreement and disagreement trials.

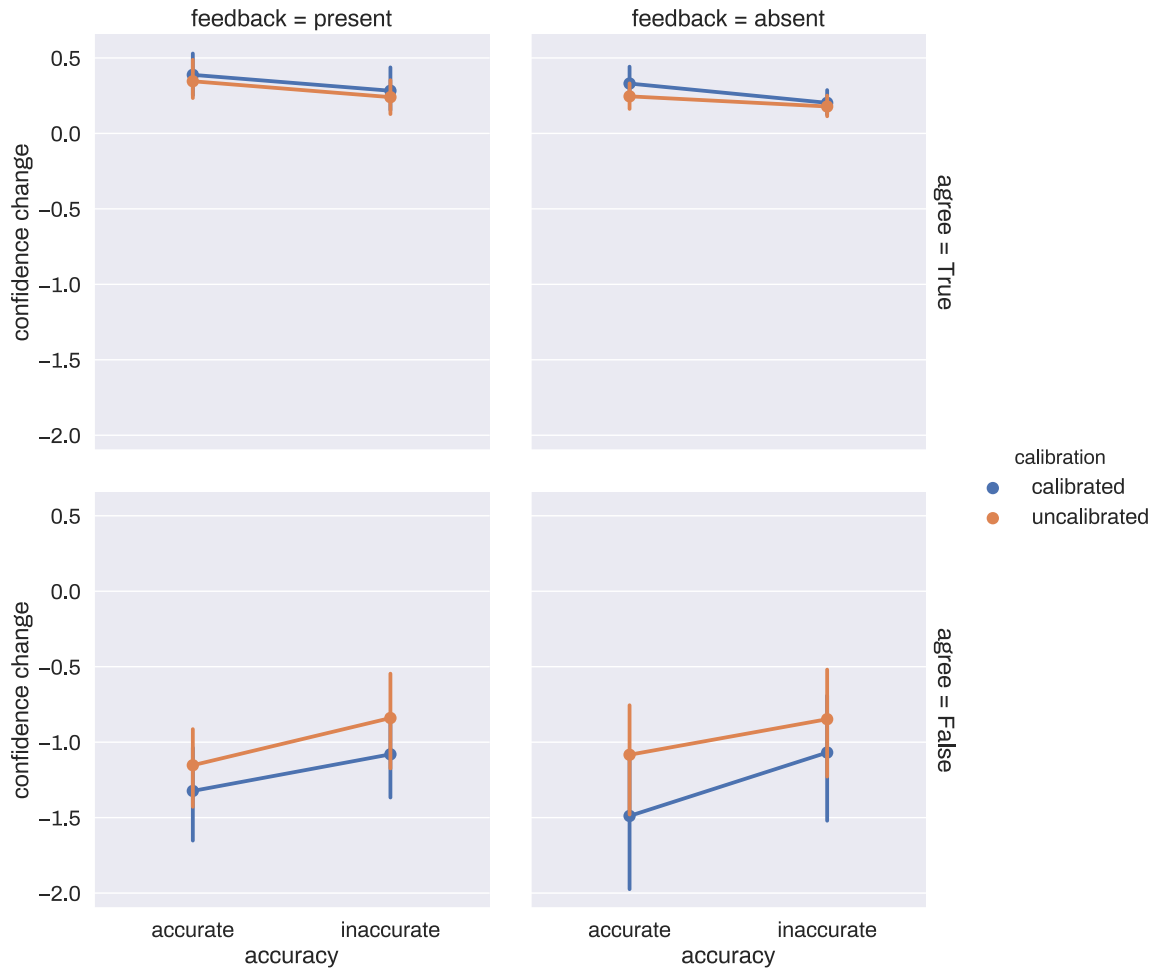


Figure S2: Experiment 1 - Confidence change in agreement and disagreement trials for different advisors. We observe an effect of both advisor's accuracy and calibration (see table S2) in both agreement and disagreement trials. This suggests that people accumulated a representation of others' competence irrespective, instead of simply discounting disagreeing evidence.

Effect	F(1,44)	p	ges
Agreement			
fb	0.75731824	3.888921e-01	1.082965e-02
acc	9.22407505	4.006146e-03**	1.969010e-02
cal	4.00170962	5.164876e-02.	4.576533e-03
conf	69.29001445	1.388161e-10***	1.311394e-01
fb:acc	0.01367809	9.074292e-01	2.978340e-05
fb:cal	0.06240406	8.038984e-01	7.169104e-05
fb:oc	4.65711595	3.642327e-02 *	1.004259e-02
acc:cal	0.41823365	5.211792e-01	4.502113e-04
acc:conf	0.43401602	5.134579e-01	3.094162e-04
cal:conf	3.37609059	7.290780e-02 .	2.128078e-03
fb:acc:cal	0.41563349	5.224712e-01	4.474136e-04
fb:acc:conf	0.17938831	6.739627e-01	1.279117e-04
fb:cal:conf	2.98407362	9.110176e-02 .	1.881440e-03
acc:cal:conf	0.74141098	3.938773e-01	2.551690e-04
fb:acc:cal:conf	1.83713762	1.822040e-01	6.320432e-04
Disagreement			
fb	0.006121536	9.379917e-01	1.011952e-04
acc	15.304270849	3.133341e-04 **	1.725678e-02
cal	18.823806847	8.256714e-05 ***	1.270240e-02
conf	38.697718603	1.598035e-07 ***	8.451080e-02
fb:acc	0.105297640	7.471003e-01	1.208017e-04
fb:cal	0.812641129	3.722487e-01	5.551212e-04
fb:conf	0.377914265	5.418879e-01	9.006905e-04
acc:cal	0.164358350	6.871398e-01	1.563503e-04
acc:conf	2.652761665	1.105102e-01	1.045018e-03
cal:conf	7.753482060	7.874907e-03 *	3.359669e-03
fb:acc:cal	0.814506124	3.717056e-01	7.743419e-04
fb:acc:conf	0.572707598	4.532184e-01	2.257951e-04
fb:cal:conf	0.189169288	6.657348e-01	8.223868e-05
acc:cal:conf	4.049781660	5.032305e-02 .	8.000137e-04
fb:acc:cal:conf	2.531817122	1.187317e-01	5.002976e-04

Table S2: Experiment 1 - confidence change broken down by agreement. We find similar results reported for influence reported in the main text, namely main effects for advisor's accuracy and calibration. Columns from left to right: Effect, F statistic (numerator's degrees of freedom, denominator's degrees of freedom), p-value, generalized η_G^2 measure of effect size. Effect abbreviations: fb (feedback), acc (advisor accuracy), cal (advisor calibration), conf (advisor confidence). Significance values: . (< .1), * (< .05), ** (< .01), *** (< .001).

2.3 Experiment 2

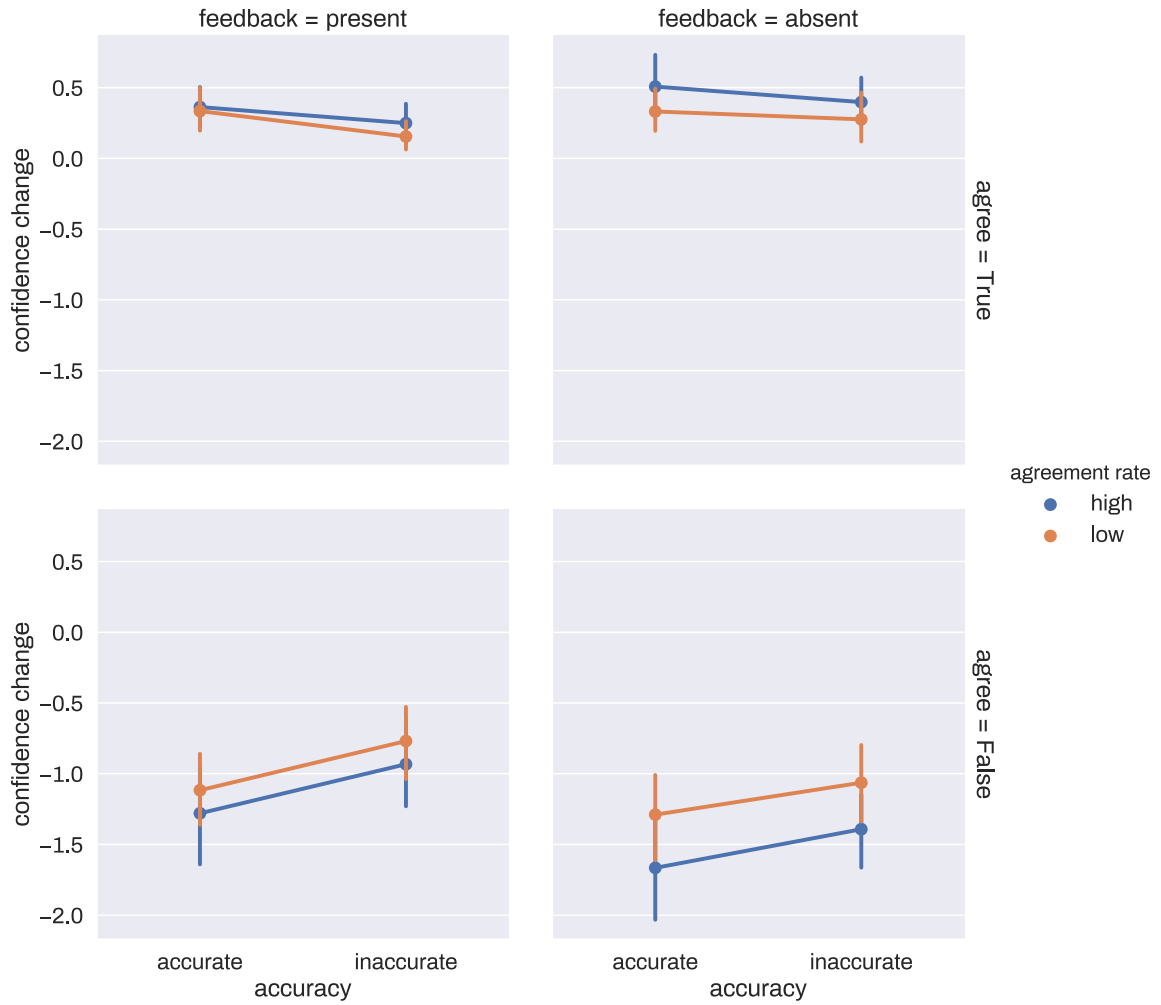


Figure S3: Confidence change in agreement and disagreement trials for different advisors. We observe an effect of both accuracy and agreement rates (see table S3) in both agreement and disagreement trials. This suggests that people accumulated a representation of others' competence irrespective, instead of simply discounting disagreeing evidence.

Effect	F(1,44)	p	ges
Agreement			
fb	0.9773923	0.328249136	1.718636e-02
acc	11.0912363	0.001762675 **	2.126074e-02
agr	11.6915961	0.001364394 **	1.796254e-02
fb:acc	0.8654221	0.357301932	1.692091e-03
fb:agr	2.0170753	0.162585940	3.145717e-03
acc:agr	0.0122558	0.912353006	1.609034e-05
fb:acc:agr	1.1344127	0.292646089	1.487152e-03
Disagreement			
fb	3.50257108	0.067931311 .	4.962705e-02
acc	13.31128416	0.000695565 ***	4.141194e-02
agr	12.24407362	0.001081278 **	3.128264e-02
fb:acc	0.35848098	0.552421231	1.162077e-03
fb:agr	1.64668779	0.206128216	4.324237e-03
acc:agr	0.03461100	0.853268206	6.699501e-05
fb:acc:agr	0.03595038	0.850489994	6.958741e-05

Table S3: We find the same main effects reported for influence reported in the main text. Columns from left to right: Effect, F statistic (numerator's degrees of freedom, denominator's degrees of freedom), p-value, generalized η_G^2 measure of effect size. Effect abbreviations: fb (feedback), acc (advisor accuracy rate), agr (advisor agreement rate). Significance values: . (< .1), * (< .05), ** (< .01), *** (< .001).

2.4 Experiment 3

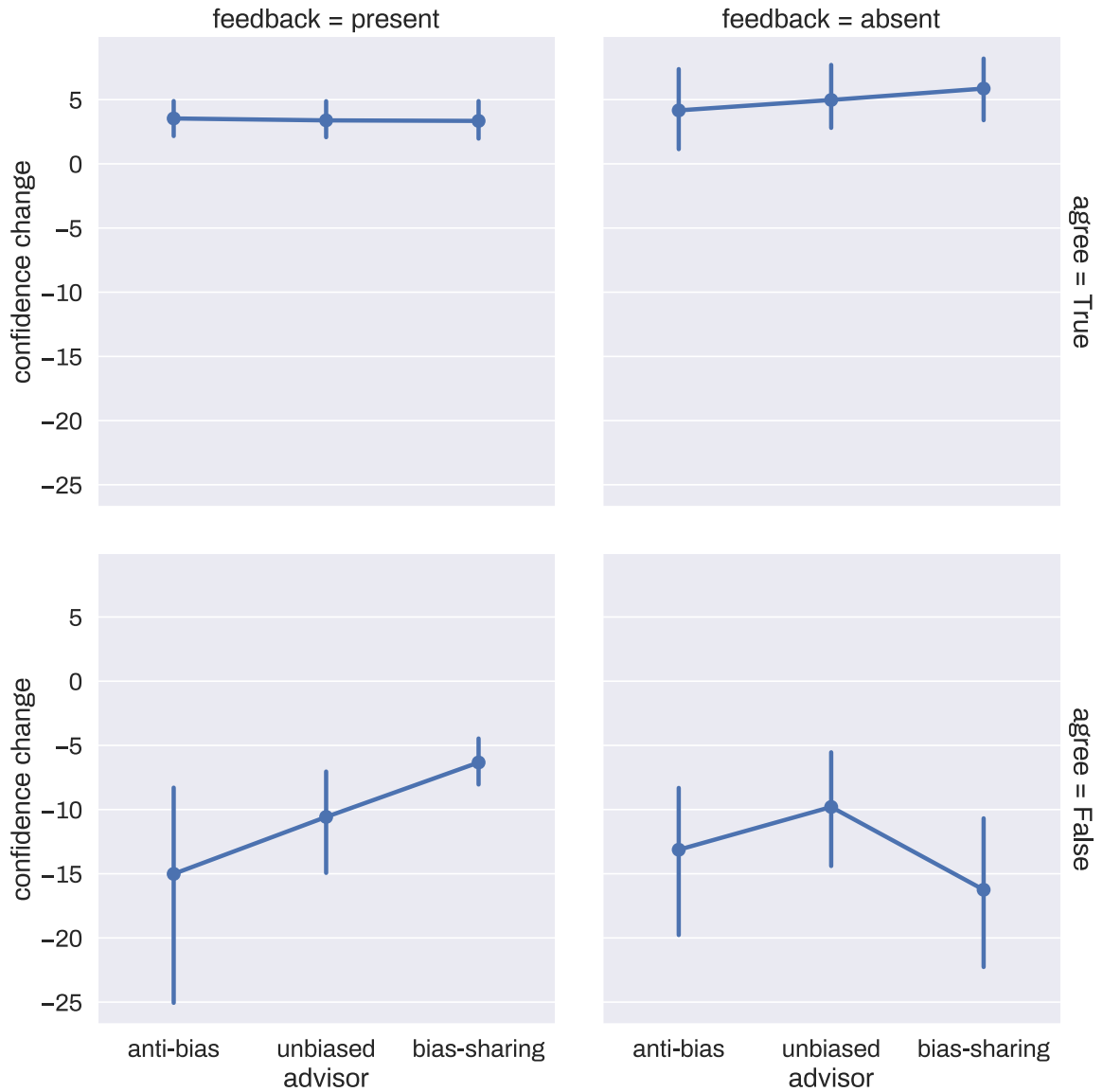


Figure S4: Confidence change in agreement and disagreement trials for different advisors. We observe an interaction between feedback and advisor type in both agreement (n.s.) and disagreement trials (Table S4). This suggests that people accumulated a representation of others' competence irrespective, instead of simply discounting disagreeing evidence.

Effect	DFn	DFd	F	p	ges
Agreement					
fb	1	46	1.3045963	0.2592860	0.022652313
adv	2	92	0.8914123	0.4135879	0.003529282
fb:adv	2	92	1.3948125	0.2530714	0.005511361
Disagreement					
fb	1	46	0.5992415	0.44282813	0.007701738
adv	2	92	1.6127341	0.20492582	0.013972922
fb:adv	2	92	4.2738736	0.01679238 *	0.036194829

Table S4: Experiment 3: confidence change broken down by agreement. In agreement trials, we observe the same numerical trend as reported in the main text, although these results are not significant. In disagreement trials on the contrary, the same interaction between feedback condition and advisor type is observed as the one reported in the main text, suggesting that disagreement trials might have been driving the main results. One possibility for this difference is the ceiling effect observed in agreement trials, which left more room for confidence change in the disagreement than in the agreement part. Columns from left to right: Effect, numerator’s degrees of freedom, denominator’s degrees of freedom, F statistic, p-value, generalized η_G^2 measure of effect size. Effect abbreviations: fb (feedback), adv (advisor type). Significance values: . (< .1), * (< .05), ** (< .01), *** (< .001).

3 Three heuristics for advisor competence estimation

We used the models to explore how people would estimate advisor reliability if they were using three simple heuristic algorithms that make use of readily available information in the decision process: objective feedback, if available; the degree to which advisors agree with the participants’ own initial judgments; and advisors’ agreement with the participants’ own judgments, scaled by the confidence with which those initial judgments are made. Notice that these models do not intend to be a faithful representation of the cognitive underpinnings of our human participants but a proof of concept showing that even in the absence of external feedback, confidence and agreement signals can be accumulated over time to form stable impressions of advisor’s accuracy and reliability. Importantly, once these impressions are formed, they can inform a more flexible use of advice. For example, instead of simply down-weighting advice by confidence (as prior studies have shown), a stable representation of the advisor’s underlying accuracy can be used to down-weight their advice also when their advice agrees with one’s own current opinion. However, as thoroughly investigated in the main text, this strategy relies on the independence of one’s own and one’s advisor’s judgments. If this independence is broken (e.g. if the self and the advisor are more likely to agree on incorrect choices) then this adaptive strategy backfires because it systematically overestimates the accuracy of highly agreeing individuals.

3.1 Model Description

Experiment 1 showed that people are sensitive to similar dimensions in the advice they receive both when objective feedback is available and when it was not provided. It is unclear what cues people are following to estimate their partners' reliability when feedback is taken away. Two simple explanations can be offered. The first one is that different advisors agreed differently often with participants and this in turn was taken as an indicator of good performance. In a binary choice task, if we assume that two people's judgments are independent, agreement rate between the two will scale linearly with the accuracy of each individual as long as performance is above chance. Thus, accumulating the number of agreement events over time for each individual advisor separately allows a subject to form a stable opinion about the other person's underlying accuracy.

A related, but subtler and potentially more powerful, strategy participants could have used is to accumulate over time the estimated probability of the advice being correct on a given trial. This quantity can be generated based on internal metacognitive signals of confidence, which provide a representation (albeit imperfect) of the uncertainty associated with a given perceptual judgment. In other words, given that confidence in a decision is a probabilistic representation of the correctness of that decision, it can also be used to estimate the likelihood that the advice received is correct or incorrect. Accumulating such evidence over time can help a decision maker to estimate the reliability underlying advice whenever more secure signals are not available.

To formalise such hypotheses we implemented a simple model that uses different pieces of information depending on different experimental conditions to estimate advisors' reliability. This simple model can then be compared with human observers to provide insight into the strategies they are using to evaluate advice reliability. Three different model variants are described below that account for the Feedback condition, the No-Feedback condition without metacognitive insight and the No-Feedback condition with metacognitive insight respectively.

3.2 Accuracy Model

When objective feedback is given to participants by the experimenter, the model can use it to infer the accuracy rate of its advisors. The accuracy of the advisor ($Acc = \{0, 1\}$) is the same as the accuracy of the subject in agreement trials, while is opposite in disagreement. By counting correct and error rates for each advisor separately, the model obtains a trial-by-trial estimation of the advisor's accuracy rate, θ , as the ratio between the number of advisor's correct trials and the total encounters with that advisor:

$$\theta^i = \frac{\alpha^i}{\alpha^i + \beta^i} \quad (5)$$

where α^i and β^i are the correct and error counts respectively, during the past trials with advisor i :

$$\alpha^i = \sum_{t=1}^n Acc_t \quad (6)$$

$$\beta^i = \sum_{t=1}^n 1 - Acc_t \quad (7)$$

$t = 1$ here represents the first encounter with advisor i while $t = n$ represents the last one. A slight complication in Experiment 1, however, is that advisors also provided a binary confidence judgment associated with the advice. A simple way for the model to make use of advisor’s confidence is by treating it as a linear scaler of the advice received. We applied a set of arbitrary weights to the four possible advice scenarios, namely the advisor is (1) correct and confident, (2) correct but unsure, (3) incorrect and unsure and (4) incorrect but confident (Table S5). Although arbitrary, any set of weights that preserves the order of such events would result in similar final advisor preferences.

	Event Observed			
	Inaccurate Confident	Inaccurate Unsure	Accurate Unsure	Accurate Confident
Feedback	-1	-0.5	+0.5	+1
No-Feedback	Disagree Confident	Disagree Unsure	Agree Unsure	Agree Confident
	-1	-0.5	+0.5	+1

Table S5: Model weights (w) applied to different advice events observed in the Feedback and No-Feedback scenario.

Thus instead of simple accuracies, α and β in equations 6 and 7 can now be reformulated as:

$$\alpha^i = \sum_{t=1}^n .5 + .5 * w_t \quad (8)$$

$$\beta^i = \sum_{t=1}^n .5 - .5 * w_t \quad (9)$$

This set of equations results in values of 1, 0.75, 0.25 and 0 for the four events listed above respectively. Although these values could be simply summed to obtain α and β values, the unusual formulation of the equations 8 and 9 was preferred to be coherent with the equations describing the following models. They show how a simple model can take into account feedback, advice received and advisors’ expressed confidence to track over time the objective reliability of its advisors.

3.3 Consensus Model

When feedback is removed from the participants, as in the No-Feedback conditions of our experiments, the model does not have access to the advisors’ objective accuracy. It must then rely on different proxies for objective accuracy and integrate those instead over time. The first cue to underlying accuracy rate we considered is agreement rate. When two independent agents express judgments on a binary task, the agreement rate between the two linearly scales with the accuracy of each whenever the accuracy rate is higher than chance: $Agr = Acc_1 * Acc_2 + (1 - Acc_1) * (1 - Acc_2)$. We thus adapted the equations of the *Accuracy* model above to exploit this covariation. Instead of tracking the accuracy rates of its advisors, the *Consensus* model tracks their agreement rates with subjective judgments. Thus equations 8 and 9 can be used to estimate a θ value by now using as

w_t the scaled agreement observed on encounter t as described in Table S5. To take into account the fact that in Experiment 1 advisors expressed a binary confidence judgment themselves associated with the advice, we used the same linear weights applied to the *Accuracy* model also to scale agreement (Table S5). In other words this model perfectly conflates accuracy with agreement, assuming that whenever an advisor agrees with the subjective original judgment, the advisor must be correct. Although this clearly is a simplifying assumption, the model offers a useful proof of concept to understand what inferences an agent lacking metacognitive insight can make simply by using heuristics. It can thus provide a benchmark to quantify the information that is present in the advice received.

3.4 Confidence Model

A more nuanced strategy that could be employed to estimate advisors' reliability when feedback is not directly available is through use of internal metacognitive signals. Trial-level variability in subjective confidence is known to covary with objective accuracy in a perceptual task (Henmon, 1911) and it theoretically represents the estimated likelihood of having made a correct judgment and/or selected the correct response (Pouget, Drugowitsch, & Kepecs, 2016). Thus, instead of simply using agreement rates as a cue for accuracy rate, a model endowed with metacognitive insight could accumulate over time the subjective probability that an advisor expressed a correct judgment. A *Confidence* model was created under the assumption that the trial-by-trial subjective reports of confidence are directly related to the true underlying estimated probabilities of having chosen the correct answer. On agreement trials the model estimates the probability of the advice being correct as the subjective probability of a correct answer. Conversely on disagreement trials the model estimates the probability of the advice being correct as the probability of having itself made an error. In other words trial-level agreement ($Agr = \{0, 1\}$) is scaled by trial-confidence expressed as a probability over outcomes (correct vs. incorrect response). Thus equations 8 and 9 above become according to this model:

$$\alpha^i = \sum_{t=1}^n .5 + (p_t(corr) - .5) * w_t \quad (10)$$

$$\beta^i = \sum_{t=1}^n .5 - (p_t(corr) - .5) * w_t \quad (11)$$

where w_t represents the scaled trial-level agreement as described in Table S5 and $p(corr)$ represents pre-advice confidence. As described below, rather than taking participants' confidence as a pure index of subjective $p(corr)$, we transformed the value to (1) reduce inter-subjects variability and (2) increase scale sensitivity. Regardless, the crucial point is that this model capitalises on the fact that being in agreement or disagreement with an advisor is more informative when the model is itself confident that it gave a correct answer than when it is more likely to have made a mistake.

Experiments 2-3. In Experiments 2 and 3, advisors did not express a level of confidence with their judgments. This allowed to simplify the above equations describing the three models. The *Accuracy* model could be simplified using equations 6 and 7 to

compute α and β for each advisor instead of equations 8 and 9. Similarly, the *Consensus* model now computes α and β values for each advisor i separately as:

$$\alpha_i = \sum_{t=1}^n .5 + .5 * Agr_t \quad (12)$$

$$\beta_i = \sum_{t=1}^n .5 - .5 * Agr_t \quad (13)$$

where Agr_t is the partner's consensus ($Agr = \{-1, 1\}$) observed on encounter t . Finally, the simplified *Confidence* model computes α and β values as:

$$\alpha = \sum_{t=1}^n .5 + (p(corr) - .5) * Agr_t \quad (14)$$

$$\beta = \sum_{t=1}^n .5 - (p(corr) - .5) * Agr_t \quad (15)$$

where $p(corr)$ is the pre-advice confidence expressed in probability scale as described in equation 18.

3.5 Bayesian update

All model variants can use the current estimated advisor's reliability θ to appropriately update the pre-advice probability of having selected the correct answer $p(corr)$ into a normative posterior, based on the binary advice A received (agree vs. disagree):

$$p(corr|A^i) = \frac{p(corr)p(A^i|corr)}{p(corr)p(A^i|corr) + p(err)p(A^i|err)} \quad (16)$$

where $p(err)$ is the subjective probability of making a mistake on the current trial and $p(A^i|corr)$ is the probability that advisor i agrees or disagrees given that the participant's choice is correct. Prior probability $p(corr)$ is estimated from a simple linear transformation of the pre-advice trial-level confidence data obtained from the participants after appropriate pre-processing. Pre-processing consisted in a parameter-free transformation that (a) brings all subjective confidence distributions on to a similar scale thus reducing the inter-subject variability and (b) expands the centre of the original subjective confidence distributions so to increase the informativeness of the average trial. This operation was inspired by recent models of adaptive information gain control (Cheadle et al., 2014). According to these proposals, the brain adapts the gain of neuronal firing to the range of information available over different time scales and cognitive domains (Carandini & Heeger, 2011; Cheadle et al., 2014). Here it serves the purpose of increasing the discriminability or information gain of different trials so that trials that are close together on confidence scale gets pulled apart on to a probability scale. The transformation uses parameters obtained from the data:

$$\hat{C}_{pre} = N * normcdf(C_{pre}) \quad (17)$$

where $normcdf(C)$ is the normal cumulative density function of the pre-advice confidence C ratings distribution, and N is the number of confidence ratings available on each interval of the scale (in Experiments 1,2: $N = 5$; in Experiment 3: $N = 50$). This

simple transformation has the property of translating a normal distribution into a uniform distribution in the range $[0, N]$. Notice that this transformation does not affect the ranking of confidence judgments but only their spacing along a probability scale. After pre-processing, confidence ratings were translated into a probability scale with the linear transformation:

$$p(\text{corr}) = 0.5 + (0.1 - \epsilon) * \hat{C}_{pre} \quad (18)$$

where ϵ is a small jitter ($\epsilon = .002$) introduced to avoid maximum confidence ratings being turned into probability of one and zero, which would in turn cause inconsistencies within the Bayesian formula (e.g., no confidence change regardless of advice reliability). Thus $p(\text{corr})$ represents trial-level confidence on a probability scale, which can be interpreted as the probability that the participant assigns to having given a correct answer on a given trial. From $p(\text{corr})$ we can also derive the subjective probability that a given trial will end up in an error: $p(\text{err}) = 1 - p(\text{corr})$.

To estimate the likelihood term $p(A^i | \text{corr})$ in equation 16 we applied a simple heuristic that uses the reliability θ of a given advisor:

$$p(A^i | \text{corr}) = \theta^A * (1 - \theta)^{1-A} \quad (19)$$

The equation above simply states that the probability of observing advisor i 's agreement ($A^i = 1$) when the participant is correct is equal to the accuracy rate of the advisor itself, assuming advisor's and participant's judgments are independent. Conversely, the probability of observing disagreement ($A^i = 0$) on the same trials is the advisor's error rate. In other words, the probability of agreement in trials when the participant is correct is the probability that the advisor too is correct. Similarly, the probability of disagreement in trials when the participant is correct is equal to the probability that the advisor is wrong.

3.6 Model results

Trial-by-trial agreement, reported confidence and objective feedback from the experimental data of the 46 participants in Experiment 1 were used to estimate rated competence and influence that the our three heuristic model variants would show with each advisor if they had experienced the corresponding advice profiles of the four virtual advisors. Separate model runs simulated the evolution of the accuracy estimate (theta parameter) according to the three learning rules described above: the *Accuracy* model that learns based on trial-by-trial feedback, which by hypothesis should capture patterns of rated competence and influence from participants the Feedback condition, and the *Consensus* and *Confidence* models, which provide distinct computational accounts of the evolution of rated competence and influence in the No Feedback condition—whether it depends purely on rates of agreement, or whether agreement is weighted according to participants' confidence in their own initial judgments.

Experiment 1 The three model variants were applied to the actual series of each participant's decisions and (where appropriate) their associated confidence, and the advice they received and (where appropriate) its accuracy, in Experiment 1. The aim of the following analyses was to verify how the pattern of final model's trust (Θ) in each advisor differed when different pieces of information were used to compute it. The models

are not intended as a mechanistic description of participants’ behaviour, but rather aim to explore how simple accuracy estimation rules lead to differentiated patterns of trust across advisors according to the type of information used to update estimates of advisor accuracy. For this analysis, data from the Feedback and No-Feedback groups were pooled together to increase statistical power, as the presence of feedback did not affect the variables that model’s variants were based on, namely advisors’ accuracy, agreement rates and participant’s pre-advice confidence ratings respectively. Our analysis focuses on the resulting values for each model variant across simulated participants, as a direct measure of the model’s belief about advisor accuracy.

The three model variants’ final Θ values were analyzed using a 2x2 repeated measures ANOVA with factors of advisor Accuracy (high vs low) and Calibration (high vs. low). Scaling factors were applied to agreement to take into account the fact that in this experiment advisors provided a confidence judgment with their advice. The Accuracy variant is, in this experiment, fully pre-determined by the advisors’ set accuracy rates and thus no statistical analysis was run due to the absence of variability across participants. We plot however its trust value for visual reference as it shows that when the model is provided with information about the objective performance of the participant (and thus of the advisors), it is able to distinguish advisors both in terms of their Accuracy rate and their confidence Calibration. The difference between calibrated and uncalibrated advisors is larger for inaccurate than accurate advisors, due to the weights used to convey advice confidence (Figure S5).

In the absence of objective feedback, both the Consensus variant—which estimates advisors’ reliability by assuming that advisors are correct whenever they agree with the participant’s own first decision, and wrong otherwise—and the Confidence variant—which uses agreement as a proxy for feedback like the Consensus variant, but scales them by pre-advice confidence—show greater trust for Accurate ($F(1, 45) > 165.84, p < .001, \eta_G^2 = .44$) and Calibrated ($F(1, 45) > 32.70, p < .001, \eta_G^2 = .13$) advisors compared to inaccurate or uncalibrated ones. Neither variant showed a significant interaction between the two factors ($F(1, 45) < 2.65, p > .11, \eta_G^2 = .01$). Notice that advisors were not constrained to agree with the participant a pre-determined number of times, thus explaining the variability observed across participants according to the specific sequence of decisions they made and advice they received. Taken together, these modeling results show that simple computations of advisor reliability perform well at this task even when trial-level feedback is absent, effectively capturing key patterns of trust observed in the human data across feedback conditions. For the task used in Experiment 1, the three variants do not make contradictory predictions on which advisors should be trusted. In particular the two No-Feedback variants (that base trust on simple agreement, or agreement weighted by confidence) cannot be disentangled using the data collected from this experiment, with both showing sensitivity to both accuracy and calibration of an advisor.

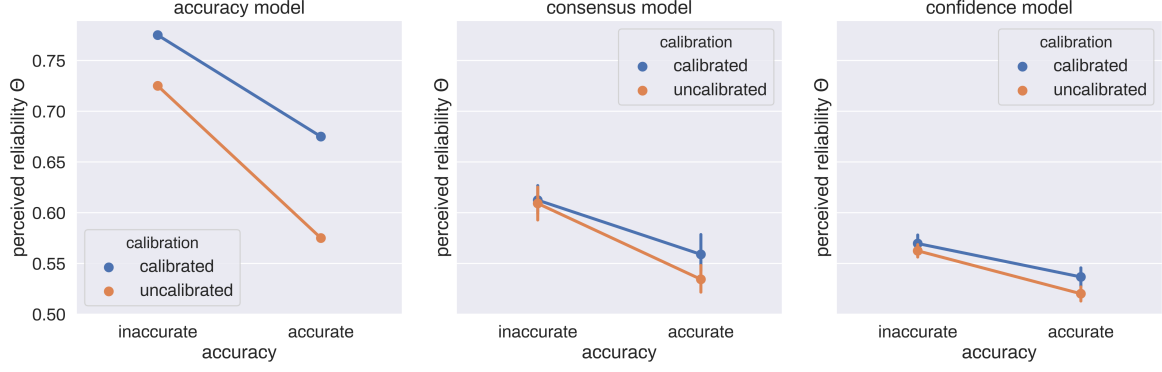


Figure S5: Experiment 1 - Heuristic models

Experiment 2 The models described above were applied to data from Experiment 2, to understand whether the Consensus and Confidence model variants behaved differently in scenarios where advice accuracy and advice agreement rate are dissociated. In this experiment, advisors did not express a confidence judgment about their opinions. Thus, all model variants could be simplified by not taking into account advice confidence. Trial-by-trial pre-advice confidence and advice were input to each of the three model variants and resulting Θ -values for each advisor were compared (Figure S6). Both the Accuracy and the Confidence models' Θ values showed a significant effect of Accuracy ($F(1, 45) > 8.85, p < .005, \eta_G^2 > .05$), while the Consensus model only showed a non significant marginal effect ($F(1, 45) = 3.22, p = .07, \eta_G^2 = .02$). Both the Confidence and Consensus models show a significant effect of Agreement ($F(1, 45) > 434.7, p < .001, \eta_G^2 > .71$), but no reliable interaction between the two factors ($F(1, 45) < 2.04, p > .15, \eta_G^2 < .007$).

Not surprisingly, when provided with objective feedback on trial-by-trial performance, a simple model of reliability estimation (Accuracy variant) distinguished advisors based on their accuracy but not their agreement profile. More surprisingly, a model without access to feedback but endowed with metacognitive insight (Confidence variant) was also able to discriminate between equally agreeing but differently accurate partners. As shown in Table 2 (main text), the accurate agreeing advisor tends to agree more often than the inaccurate agreeing advisor when the participant is objectively correct (6.5 times out of 7 against 5.5 times out of 7) and less often when the participant is objectively wrong (1.5 times out of 3 against 2.5 times out of 3). Trials when participants' initial judgment is correct are usually associated with greater confidence ratings (Fleming et al., 2014; Henmon, 1911; Koriati, 2012), thus a strategy of reliability estimation relying on confidence can exploit this covariation to detect differences in accuracy, notwithstanding equal agreement rates.

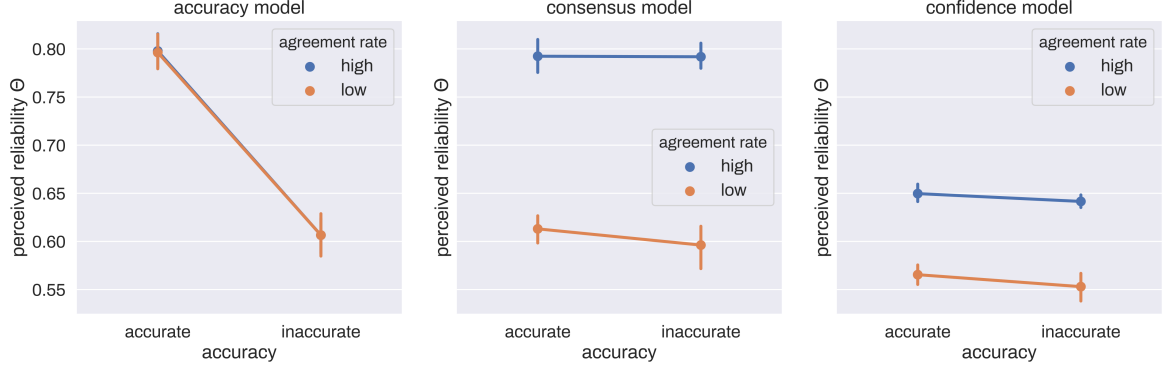


Figure S6: Experiment 2 - Heuristic models

Experiment 3 The following simulations and analyses explored the differing patterns of trust across advisors predicted by simple models of estimating advisor reliability. Simulations are also useful in this experiment as a check that our designs were controlled as intended (e.g., for agreement rates across advisors) even though we had less precise control over conditions because counterbalancing depended on an evolving estimate of participants' confidence distributions. Figure S7 shows the pattern of results (modeled $-$ values) that the three model variants produce.

Both when the model has access to trial-by-trial feedback (Accuracy variant), and when it only has access to past agreement (Consensus variant), no significant effect of Advisor is observed ($F(2, 94) < 1.70, p > .18, \eta_G^2 < .02$), nor is there a difference between the bias-sharing and the anti-bias advisor. These patterns are expected because the three advisors were matched for accuracy and agreement rates by design in this experiment. On the contrary, a Confidence variant which uses metacognitive information and past agreement (but lacked access to trial-level feedback) showed a significant effect of Advisor ($F(2, 94) = 7.95, p < .001, \eta_G^2 = .10$). Specifically, simulated trust was higher for the bias-sharing advisor than the anti-bias advisor ($t(47) = 3.54, p = .001, d = .74$), and higher for the unbiased advisor than the anti-bias advisor ($t(47) = 2.99, p = .004, d = .57$). Simulated trust was higher for the unbiased than the bias-sharing advisor, but this difference was not reliable ($t(47) = 1.21, p = .22, d = .25$). These findings indicate that by accessing metacognitive signals (as provided in the model by participants' confidence ratings) the model was able to discriminate among different advisors. This model correctly predicts greatest levels of trust in a bias-sharing advisor, but also predicts lowest levels of trust in anti-bias advisors, whereas our experimental participants expressed (numerically) lowest levels of trust in unbiased advisors.

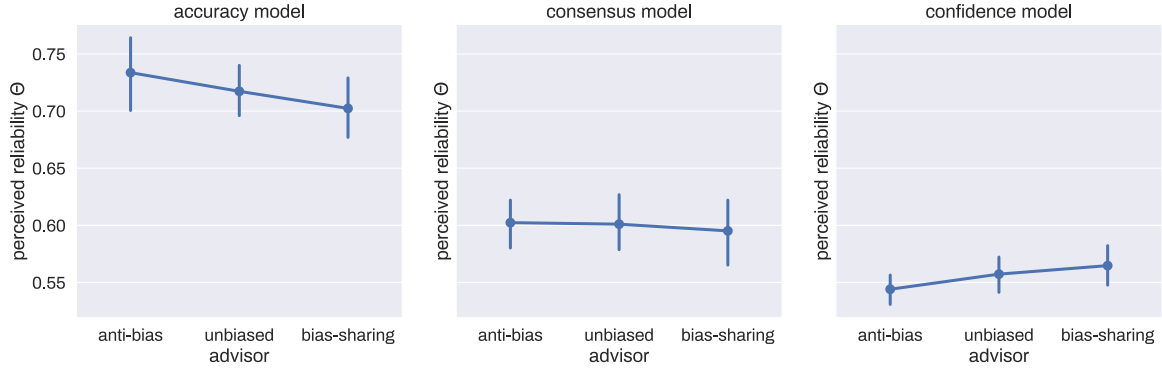


Figure S7: Experiment 3 - Heuristic models

4 Post-advice confidence correlations.

Our main analyses for each experiment focused on qualitative predictions arising from different strategies for inferring advisor reliability. Collectively, the results are consistent with the hypothesis that people use their internal sense of confidence in making these inferences—showing sensitivity to advisor accuracy in the absence of objective feedback, even when advisors are matched for agreement rate, and developing differing patterns of trust when advisor agreement rates vary with their own expressed decision confidence. Our final analysis of the empirical data attempted a more quantitative comparison of model predictions, specifically focusing on whether the *Consensus* or *Confidence* models better predicted post-advice confidence ratings across trials for the participants in the No Feedback conditions of Experiments 1-3.

For this analysis we used Bayes rule to infer the trial-by-trial post-advice confidence ratings that each variant would express given a participant’s expressed pre-advice confidence and advisor agreement (as defined above). The within-participant correlation between participants’ post-advice confidence and model’s post-advice confidence was computed for each experiment, for No-Feedback groups only. Second-order statistics were performed to test, across experiments, which variant was more strongly correlated with the human data. A 2x2 ANOVA on correlation coefficients with Model (*Consensus* vs. *Confidence*) and Experiment as factors showed that the *Confidence* variant was significantly more correlated with participants’ responses than the *Consensus* variant ($F(1, 22) = 8.18, p = 0.009, \eta_G^2 = 0.0049$). No significant effect of experiment nor interaction between the two were found ($F(2, 44) < 1.3, p > .25$), suggesting that, across experiments, the *Confidence* model’s post-advice confidence more strongly covaried with participants’ true post-advice responses. As a check for the soundness of this model comparison method, the same 2x2 ANOVA was run on the correlation coefficients between participants’ post-advice confidence and the model’s post-advice confidence predictions, after randomly shuffling trials within each participant. This operation should ensure that any advantage of the *Confidence* variant over the *Consensus* variant is not due to unspecific factors (like being overall more conservative in updating confidence), but rather to trial-level variability. After reshuffling, the *Consensus* and *Confidence* variants were not significantly different from each other ($F(1, 22) = 1.95, p = 0.17, \eta_G^2 = 9.54e - 04$), corroborating our conclusions.

5 Agent-based simulation

5.1 Model description

An agent-based model was programmed using NetLogo (Wilensky, 1999) and is available at https://github.com/chri4354/trust_formation_without_feedback. The model was initialised as a fully connected directed network of N agents. A directed edge from agent i to agent j represents the trust $\theta_{i,j}$ that i has in j 's opinions. We simulate agents on a lattice network performing repeated binary A/B decisions, receiving advice from other agents, inferring their reliability and updating their own initial decisions. We let the simulation run for a 1000 steps. A signal s with strength S was drawn from a uniform distribution between $-\frac{S}{2}$ and $+\frac{S}{2}$. This represents the decision quantity to estimate (e.g., difference in dots or true state of the world). The task of each agent was to determine if s was positive (event A) or negative (event B). Each agent estimated the posterior probability of A given the perceptual information generated by s as follows:

$$p'(A) = p(A|E_p) = \frac{p(A)E_p}{p(A)E_p + p(\bar{A})\bar{E}_p} \quad (20)$$

$$E_p = L(s + \mathcal{N}(0, \sigma)) \quad (21)$$

where $p(A)$ is the prior probability of observing A s before seeing any stimulus, L is a logistic sigmoid mapping from sensory evidence to probability; E_p is the perceptual evidence resulting from such mapping and \mathcal{N} is independent individual perceptual Gaussian noise with mean 0 and standard deviation σ . Bars represent complement probability. Each agent's perceptual noise (and thus accuracy) was manipulated by varying the noise parameter σ . Each agent's bias was manipulated by varying the initial value $p(A)$. Agents' confidence was represented as the distance from the uncertainty point 0.50:

$$C = .50 + |p'(A) - .50| \quad (22)$$

Trust, represented by the network's edges, was initialized to 0.50 for every agent and updated after social interaction. After making a judgment, agents selected one other agent to interact with either at random (random sampling) or proportionally to their trust (biased sampling). Agents then updated their initial judgment $p'(A)$ as follows:

$$\hat{p}(A) = p(A|E_s) = \frac{p'(A)E_s}{p'(A)E_s + p'(\bar{A})\bar{E}_s} \quad (23)$$

where E_s represents social evidence and is obtained from the advisor's judgment either by taking the advisor's raw judgment $p'(A)$ (without advice discounting) or by discounting the advisor's judgment proportionally to the agent's trust in the advisor (with advice discounting). Advice discounting consisted in a linear regression toward the uncertainty point 0.50 using the following equation:

$$E'_s = 0.5 + (\theta * (p'(A) - 0.5)) \quad (24)$$

The above equation regresses any confidence judgment $p'(A)$ toward the uncertainty point 0.50 proportionally to trust. A trust level of 1 would leave the advisor's judgment $p'(A)$ untouched, while a trust level of 0 would make any advisor's judgment equal to 0.50 and thus entirely uninformative. After updating their judgments, each agent i updated

its trust judgments (i.e., outward edges θ_i) based on the available information about other agents. If feedback is available, the agent updates its current trust in agent j by virtue of a delta rule in the form:

$$\theta_{i,j}^{t+1} = \theta_{i,j}^t + \alpha(F_j - \theta_{i,j}^t) \quad (25)$$

where F_j is the accuracy of agent j and α is a learning rate set to 0.1. If feedback is not available on the contrary, the agent replaces F with \hat{F} , or the *estimated* partner’s accuracy. \hat{F} was calculated using the agreement or agreement-in-confidence heuristics described above. In our simulations, we assessed the effect of feedback availability as it varied parametrically, from being available after every decision (i.e., p-feedback = 1.0), available after only some decisions (i.e., $0 < \text{p-feedback} < 1$), or never available (i.e., p-feedback = 0), rather than the simpler case of feedback presence/absence that we studied experimentally above.

The emergence of trust patterns when using agreement-based heuristics can be expected to track true accuracy when judgments are independent but generate clustering of populations when judgment correlations emerge within such populations. We defined a network’s clustering coefficient as the ratio between average trust toward agents who initially share the same bias (in-group trust) and total average trust: $\bar{\theta}_{in-group}/(\bar{\theta}_{in-group} + \bar{\theta}_{out-group})$. A ratio of 0.5 represents no preference (i.e., no difference in trust) toward agents sharing the same initial bias, while a ratio greater than 0.5 represents a preference toward agents sharing the same initial biases. We test how network clustering is shaped by the presence or absence of objective feedback and show that bias-specific segregation arises only when feedback is rarely available.

Finally, once bias-specific segregation is established, we ask whether such clustering remains stable. In particular, after 500 iterations we allow agents to dynamically change their original bias as a function of experience. For example, it is known that the bias observed in people performing binary judgments is influenced by their recent history of decisions (Akaishi, Umeda, Nagase, & Sakai, 2014; Zylberberg, Wolpert, & Shadlen, 2018). In the present context, if an agent systematically reports “A” but receives negative feedback, they should reduce their bias by decreasing their prior probability $p(A)$. Similarly, when feedback is absent, an agent who systematically reports “A” but finds themselves, after interacting with other agents, believing that B s are more frequent than expected, should reduce their bias towards A s. Conversely, bias should get stronger if the social contexts reinforces it (although see (Bail et al., 2018)). We modelled bias update with a delta rule:

$$p(A)^{t+1} = p(A)^t + \alpha(I^t - p(A)^t) \quad (26)$$

where I^t is an indicator variable that represents the final belief in the event A . When objective feedback is available, I takes the value of 1 if an event A occurred and 0 otherwise. When feedback is not available, I is set to the discrete or continuous final subjective belief in the event A . In the following section, we show the results obtained when a discrete final belief is used in the absence of feedback:

$$I = \begin{cases} 1, & \text{if } \hat{p}(A) \geq 0.5 \\ 0, & \text{otherwise} \end{cases}$$

Similar results were obtained setting I to the continuous belief $\hat{p}(A)$.

The following figures supplement figures in the main text.

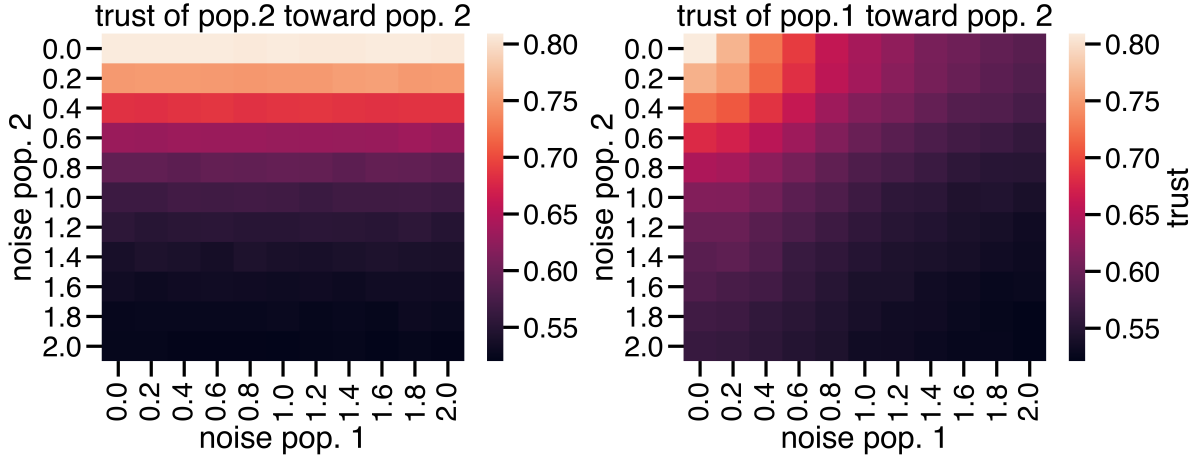


Figure S8: Trust of each subpopulation toward Population 2. Left panel: Trust of Population 2 towards Population 2 is inversely proportional to the noise of Population 2 agents (y-axis), but are (unsurprisingly) unaffected by the noise of Population 1 (x-axis). Right panel: Trust of Population 1 toward Population 2 is affected by both the noise of Population 2 and the noise of Population 1. Although the former correctly tracks Population 2's true underlying reliability, the latter reflects an interaction between the judge's characteristics and the advisor's characteristics, cf Kruger and Dunning (1999).

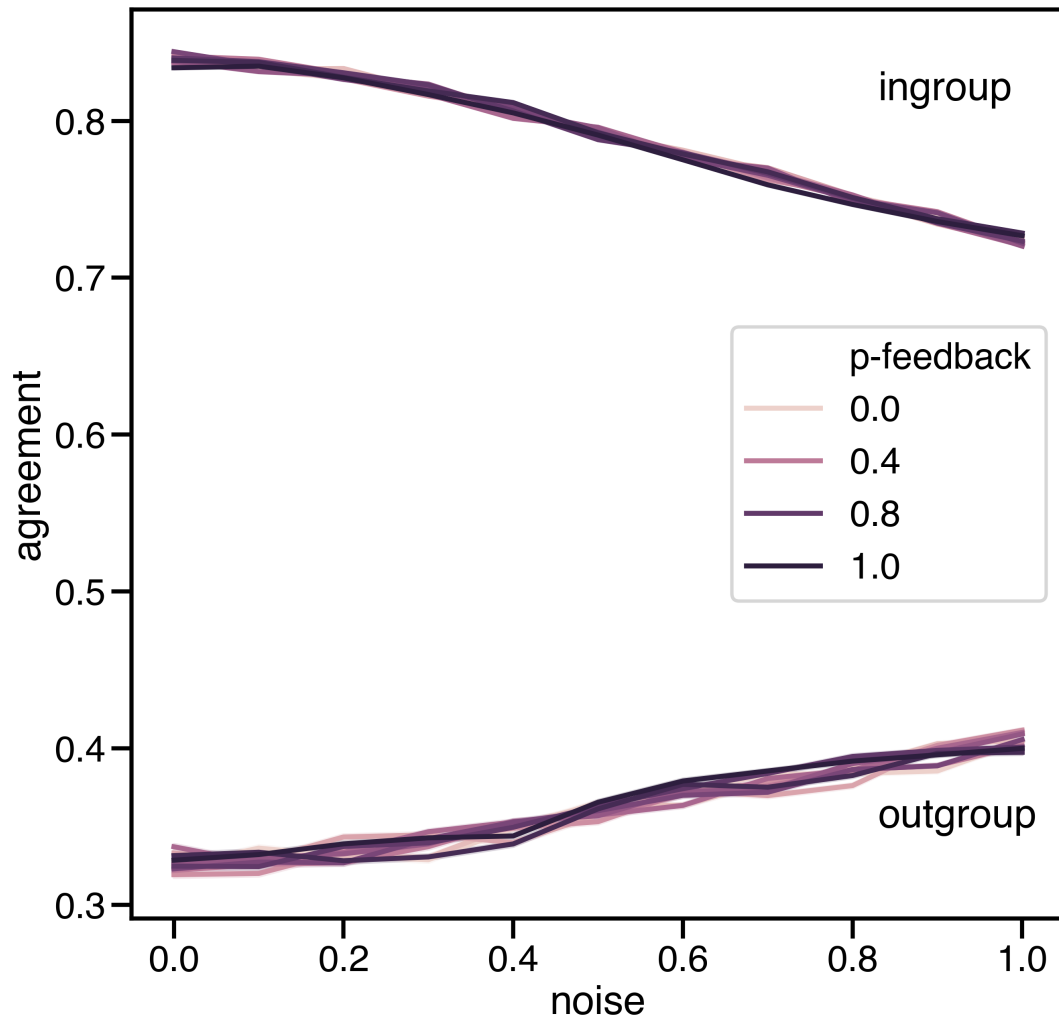


Figure S9: Average agreement rate as a function of probability of feedback and noise. Agreement with ingroup appears to decrease as a function of increasing noise, while agreement with outgroup tends to increase as a function of noise. On the contrary, feedback availability does not affect agreement rates.

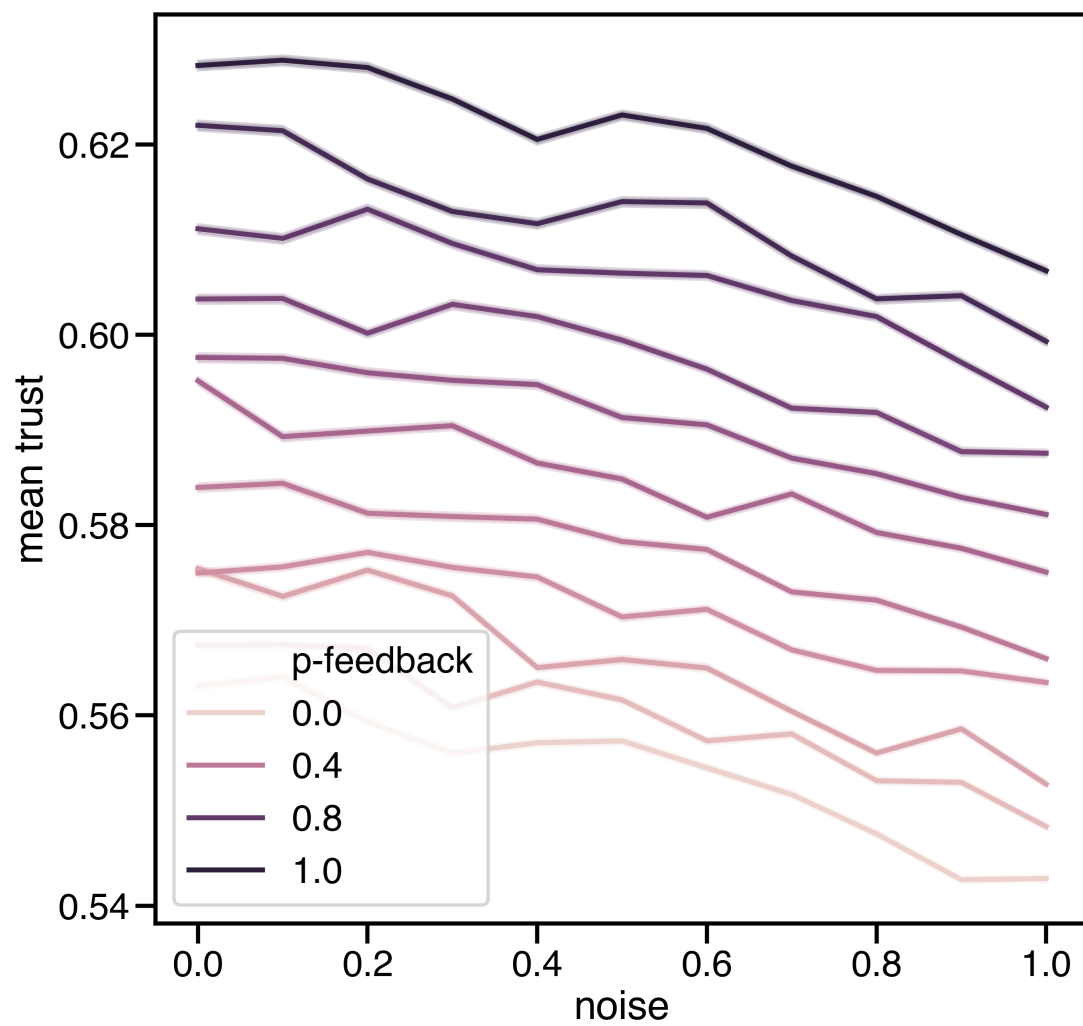


Figure S10: Average trust as a function of probability of feedback and noise. Trust appears to decrease as noise increases and feedback availability decreases.

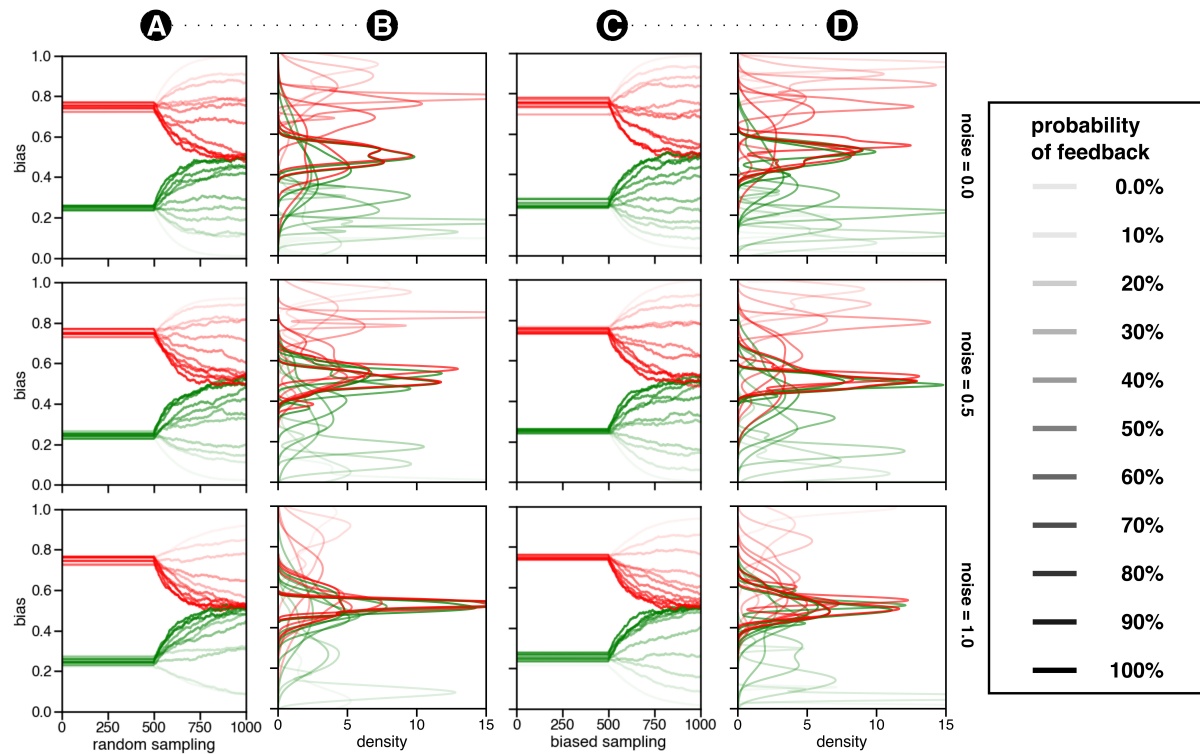


Figure S11: Bias distribution and evolution in simulations where agents discount advice proportionally to trust in the advisor.

References

- Akaishi, R., Umeda, K., Nagase, A., & Sakai, K. (2014, 1). Autonomous mechanism of internal choice estimate underlies decision inertia. *Neuron*, 81(1), 195–206. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/24333055> doi: 10.1016/j.neuron.2013.10.018
- Bail, C. A., Argyle, L. P., Brown, T. W., Bumpus, J. P., Chen, H., Hunzaker, M. B. F., ... Volfovsky, A. (2018, 9). Exposure to opposing views on social media can increase political polarization. *Proceedings of the National Academy of Sciences*, 115(37), 9216–9221. Retrieved from <http://www.pnas.org/lookup/doi/10.1073/pnas.1804840115> doi: 10.1073/pnas.1804840115
- Carandini, M., & Heeger, D. J. (2011, 11). Normalization as a canonical neural computation. *Nature Reviews Neuroscience*, 13(1), 51–62. Retrieved from <http://www.nature.com/doifinder/10.1038/nrn3136> doi: 10.1038/nrn3136
- Cheadle, S., Wyart, V., Tsetsos, K., Myers, N., de Gardelle, V., Hecce Castañón, S., & Summerfield, C. (2014, 3). Adaptive gain control during human perceptual choice. *Neuron*, 81(6), 1429–41. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/24656259> doi: 10.1016/j.neuron.2014.01.020
- Fleming, S. M., & Lau, H. C. (2014, 7). How to measure metacognition. *Frontiers in Human Neuroscience*, 8. Retrieved from http://www.frontiersin.org/Human_Neuroscience/10.3389/fnhum.2014.00443/abstract doi: 10.3389/fnhum.2014.00443
- Fleming, S. M., Maniscalco, B., Ko, Y., Amendi, N., Ro, T., & Lau, H. (2014). Action-Specific Disruption of Perceptual Confidence. *Psychological science*. doi: 10.1177/

0956797614557697

- Henmon, V. A. C. (1911). The relation of the time of a judgment to its accuracy. *Psychological Review*, 18(3), 186–201. Retrieved from <http://doi.apa.org/getdoi.cfm?doi=10.1037/h0074579> doi: 10.1037/h0074579
- Koriat, A. (2012, 4). When are two heads better than one and why? *Science (New York, N.Y.)*, 336(6079), 360–2. Retrieved from <http://www.sciencemag.org/cgi/doi/10.1126/science.1216549><http://www.ncbi.nlm.nih.gov/pubmed/22517862> doi: 10.1126/science.1216549
- Kruger, J., & Dunning, D. (1999). Unskilled and Unaware of It: How Difficulties in Recognizing One’s Own Incompetence Lead to Inflated Self-Assessments. *Journal of personality and social psychology*, 77(6), 1121–1134.
- Pouget, A., Drugowitsch, J., & Kepecs, A. (2016, 2). Confidence and certainty: distinct probabilistic quantities for different goals. *Nature Neuroscience*, 19(3), 366–374. Retrieved from <http://www.nature.com/doi/10.1038/nn.4240> doi: 10.1038/nn.4240
- Tottenham, N., Tanaka, J. W., Leon, A. C., McCarry, T., Nurse, M., Hare, T. A., ... Nelson, C. (2009, 8). The NimStim set of facial expressions: judgments from untrained research participants. *Psychiatry research*, 168(3), 242–9. Retrieved from <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3474329&tool=pmcentrez&rendertype=abstract> doi: 10.1016/j.psychres.2008.05.006
- Wilensky, U. (1999). *NetLogo*. Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL.
- Zylberberg, A., Wolpert, D. M., & Shadlen, M. N. (2018, 9). Counterfactual Reasoning Underlies the Learning of Priors in Decision Making. *Neuron*, 99(5), 1083–1097. Retrieved from <https://linkinghub.elsevier.com/retrieve/pii/S0896627318306330> doi: 10.1016/j.neuron.2018.07.035