

RESEARCH ARTICLE

A low-dimensional approximation of optimal confidence

Pierre Le Denmat^{1*}, Tom Verguts², Kobe Desender¹¹ Brain and Cognition, KU Leuven, Leuven, Belgium, ² Department of Experimental Psychology, Ghent University, Ghent Belgium* pierre.ledenmat@kuleuven.be, plendenmat1@gmail.com

Abstract

Human decision making is accompanied by a sense of confidence. According to Bayesian decision theory, confidence reflects the learned probability of making a correct response, given available data (e.g., accumulated stimulus evidence and response time). Although optimal, independently learning these probabilities for all possible data combinations is computationally intractable. Here, we describe a novel model of confidence implementing a low-dimensional approximation of this optimal yet intractable solution. This model allows efficient estimation of confidence, while at the same time accounting for idiosyncrasies, different kinds of biases and deviation from the optimal probability correct. Our model dissociates confidence biases resulting from the estimate of the reliability of evidence by individuals (captured by parameter α), from confidence biases resulting from general stimulus independent under and overconfidence (captured by parameter β). We provide empirical evidence that this model accurately fits both choice data (accuracy, response time) and trial-by-trial confidence ratings simultaneously. Finally, we test and empirically validate two novel predictions of the model, namely that 1) changes in confidence can be independent of performance and 2) selectively manipulating each parameter of our model leads to distinct patterns of confidence judgments. As a tractable and flexible account of the computation of confidence, our model offers a clear framework to interpret and further resolve different forms of confidence biases.

OPEN ACCESS

Citation: Le Denmat P, Verguts T, Desender K (2024) A low-dimensional approximation of optimal confidence. *PLoS Comput Biol* 20(7): e1012273. <https://doi.org/10.1371/journal.pcbi.1012273>

Editor: Lusha Zhu, Peking University, CHINA

Received: June 14, 2023

Accepted: June 24, 2024

Published: July 24, 2024

Copyright: © 2024 Le Denmat et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: All raw data and analysis code can be freely accessed at "https://github.com/plendenmat/ldc_paper".

Funding: The research was supported by a starting grant from the KU Leuven (STG/20/006, <https://www.kuleuven.be/kuleuven/>, awarded to KD), a Francqui Start-Up Grant from the Francqui Foundation (PXF-D8830, <https://www.francquifoundation.be/>, awarded to KD) and two grants from the Research Foundation Flanders, Belgium (FWO-Vlaanderen, <https://www.fwo.be/>) (G0B0521N awarded to KD and G010419N awarded to both KD and TV). The funders played

Author summary

Mathematical and computational work has shown that in order to optimize decision making, humans and other adaptive agents must compute confidence in their perception and actions. Currently, it remains unknown how this confidence is computed. We demonstrate how humans can approximate confidence in a tractable manner. Our computational model makes novel predictions about when confidence will be biased (e.g., over- or underconfidence due to selective environmental feedback). We empirically tested these predictions in a novel experimental paradigm, by providing continuous model-based feedback. We observed that different feedback manipulations elicited distinct patterns of confidence judgments, in ways predicted by the model. Overall, we offer a framework to

no role in the study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

both interpret optimal confidence and resolve confidence biases that characterize several psychiatric disorders.

Introduction

Decision confidence refers to a subjective feeling reflecting how confident agents feel about the accuracy of their decisions. This feeling of confidence often closely tracks the objective accuracy [1]: people usually report high confidence for correct trials and low confidence for errors. This observation is in line with the theoretical proposal that confidence reflects the Bayesian posterior probability that a decision is correct given available data [1–3]. As such, confidence represents valuable information that is taken into account to guide adaptive behavior, including learning [4–6]; speed-accuracy tradeoff adjustments [7,8]; and information seeking [9]. Therefore, having an accurate sense of confidence that best matches one's accuracy is of utmost importance to maintain adaptive behavior. However, estimating the Bayesian probability with limited data is computationally intractable. Additionally, empirical dissociations between confidence and accuracy are widespread, most prominently in cases of blindsight [10], change blindness [11] and anterior prefrontal lesions [12]. Such dissociations pose a serious challenge for the Bayesian interpretation of confidence. In this work, we reconcile these findings by proposing and empirically validating a low-dimensional approximation to the Bayesian probability, offering both a computationally tractable and flexible model for the computation of decision confidence.

Most attempts at modeling decision confidence have done so within the context of existing models of decision making. One highly influential account is based on the idea that decision making reflects a process of noisy accumulation of evidence until a decision boundary is reached [13]. For example, the drift-diffusion model (DDM) describes the decision-making process as the noisy accumulation of evidence in favor of one of two options. Here, evidence accumulates with a certain drift rate (representing the efficiency of evidence accumulation) until reaching a decision threshold, at which point a response is issued. Several approaches have been put forward to account for confidence within the DDM framework [14–16]. The most prominent approach relies on the Bayesian interpretation of confidence, modeling it as the probability of a choice being correct given the available data. Within the drift diffusion model, the available data to participants is the amount of accumulated evidence and the time spent accumulating, which are then combined into a probability that the decision was correct [2,15]. Such formalization of decision confidence is sometimes referred to as the “Bayesian readout” [17]. This Bayesian readout can be represented as a heatmap on the two-dimensional (data) space formed by both evidence and time. In Fig 1A, it can be seen that the Bayesian readout hypothesis predicts that confidence will be higher for trials with more accumulated evidence (reflected on the y-axis) and lower for trials with a longer decision duration (reflected on the x-axis). Consistent with these predictions, confidence indeed depends on evidence strength [1,2] and on elapsed decision time [14]. More generally, this modeling approach has been successful in explaining a wealth of data [17–19].

To compute confidence by reading out the probability correct given evidence and time, humans must have an accurate representation of the entire space created by crossing these two variables (i.e. the heatmap shown in Fig 1A). Previous accounts propose that individuals learn this mapping via experience [2]. However, computing the exact probability would in principle require estimating it at each point of the infinite-size (evidence, time) plane. Independently learning all positions on this heatmap in this way would either take a lot of time or yield very

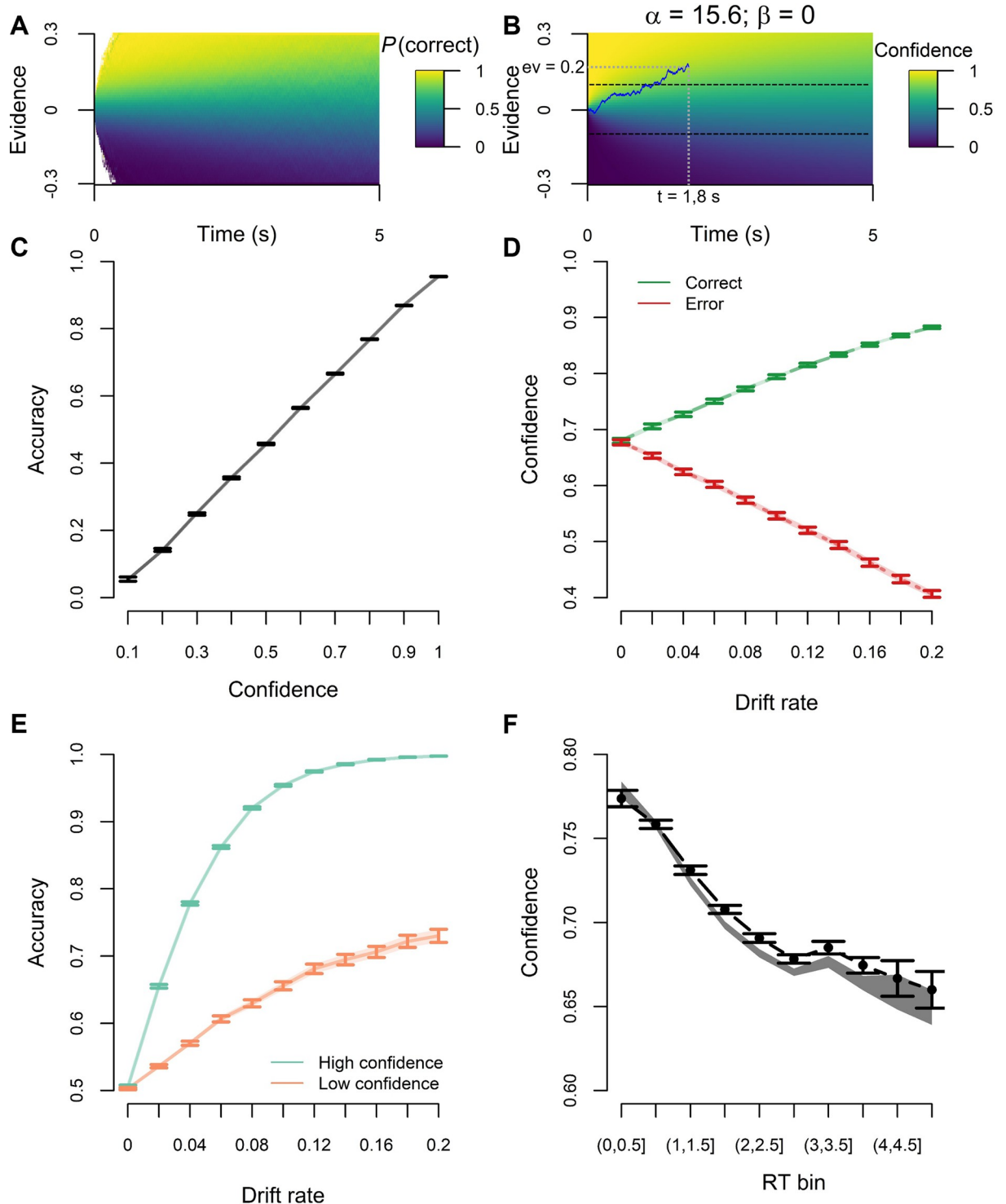


Fig 1. A. Confidence is thought to represent the Bayesian probability of a choice being correct conditional on evidence, time and choice. Within this theory, confidence is quantified as this probability, represented by the color on the heatmap. B. Because this optimal solution is intractable, the LDC model proposes a low-dimensional parametrization of this framework, which allows efficient estimation of confidence, while accounting for idiosyncrasies and confidence biases. The LDC model can generate a heatmap representing confidence which closely approximates the optimal Bayesian probability. Values of α and β were obtained by fitting the LDC model to the Bayesian probability of being correct over 1 000 000 simulated trials. Confidence for the trial plotted on top of the heatmap is given by Eq (3). Here, confidence = .85. C-F. To show the effectiveness of the LDC model we generated statistical signatures of confidence [1] based on the Bayesian read-out of confidence (error bars reflecting SEM, simulated $N = 100$) and based on the LDC model fits (shaded lines reflecting SEM). High and low confidence trials in panel E were obtained by performing a median split.

<https://doi.org/10.1371/journal.pcbi.1012273.g001>

noisy estimates. Thus, tractability is a key issue that needs to be addressed in order to understand how humans learn the probability correct given evidence and time. In typical Bayesian modeling approaches, the traditional computational solution to this high-dimensional problem is some kind of function approximation [20], whereby the probability map is approximated by a (much) lower dimensional function. Following that logic, in the current work we propose the Low-Dimensional Confidence (LDC) model, a simple yet efficient cognitive algorithm that approximates the optimal yet intractable Bayesian readout. Instead of learning the probability correct independently for all points on the (evidence, time) plane, the LDC model parameterizes the mapping in a way that only two values need to be learned by a cognitive agent to approximate the Bayesian readout. In the following sections, we describe how LDC allows to tractably compute the mapping from evidence and time to confidence. Using simulated data, we show that LDC provides a close approximation of Bayesian confidence. We then proceed to test and validate our model with human data.

Results

The low-dimensional approximation of confidence model (LDC model)

Constructing an accurate representation of confidence based on a limited number of samples is infeasible. However, under standard DDM assumptions, the probability of a correct choice given accumulated evidence and elapsed time can be expressed as the probability of drift rate ν being positive in case of upper boundary hit (and conversely $p(\nu < 0)$ in the lower boundary hit case). Such probability is characterized as [15]:

$$p(\nu > 0|e, t) = \Phi\left(\frac{e}{\sigma\sqrt{t}}\right) \quad (1)$$

where e is the accumulated evidence, t is the elapsed time, Φ is the cumulative distribution function of the standard normal distribution and σ is the within-trial noise of the DDM accumulator. Given that Φ is an integral without closed-form solution, it requires an infinite number of standard operations to be computed. We propose to approximate Φ with a more tractable logistic function (21 see [S1 Text](#) for a detailed derivation of our model):

$$\Phi\left(\frac{e}{\sigma\sqrt{t}}\right) \approx \frac{1}{1 + \exp\left(-\lambda\frac{e}{\sigma\sqrt{t}}\right)} \quad (2)$$

where $\lambda \approx 1.7$ is a constant that optimizes the approximation [21]. In its current form, the formalization of confidence proposed in Eq (2) cannot account for idiosyncrasies [22], diverse types of confidence biases and deviations from the optimal probability of a correct choice typically observed in empirical work [23–25]. We identify two distinct forms of deviation from the Bayesian probability correct: sensitivity in how evidence is treated and mapped onto confidence judgments, and a general increase or decrease in confidence ratings. In order to make the formulation of confidence flexible to these distinct forms of deviation, we thus further parameterize confidence in the following way:

$$\text{Confidence} = \frac{1}{1 + \exp\left(-\frac{1}{\sqrt{t}}(xze + \beta)\right)} \quad (3)$$

where $x \in \{-1; 1\}$ is the choice. The two free parameters of this equation capture how strongly individuals weigh evidence in their computation of confidence (α); and a stimulus-

independent confidence bias (β). Notably, the approximation of the optimal model in Eq (2) is retrieved from Eq (3) when $\alpha = \frac{z}{\sigma}$ and $\beta = 0$.

As a weighting parameter on evidence, α can be interpreted as individuals' estimate of the reliability of evidence. Intuitively, a very low α implies that during the computation of confidence participants consider the evidence as unreliable (e.g. individuals think that the stimulus is very noisy); a very high α implies that for the computation of confidence participants consider the evidence as overly reliable (e.g. individuals underestimate the amount of noise in the stimulus). In the extreme case where $\alpha = 0$, the model completely ignores evidence and the computation of confidence is entirely driven by β and time. If additionally $\beta = 0$, then confidence will always be .5. At the other end of the spectrum, if α tends to infinity, then the smallest amount of evidence will lead to extreme confidence judgments (i.e. either confidence = 1 if $ev > 0$ or confidence = 0 if $ev < 0$).

A positive confidence bias ($\beta > 0$) implies that the model has a general tendency to be overconfident. If $\beta = 0$, the model is unbiased and bases its confidence purely on the evidence accumulated and the time spent accumulating. A negative confidence bias ($\beta < 0$) indicates overall underconfidence.

It is important to note that α and β are parameters pertaining only to how the readout of evidence and time is mapped onto a confidence estimate. As such, a change in these parameters would only influence confidence, and leave the decision process unaffected. In contrast, a change in DDM parameters would influence both decision and confidence. For readers more familiar with SDT models of confidence, we note that similarly changes in μ and σ will influence both decision and confidence.

Simulations: The LDC model closely resembles Bayesian confidence

The aim of the current work is to provide a tractable and flexible approximation of the Bayesian readout of confidence. A first test of the LDC model is whether it can effectively approximate the Bayesian readout of confidence. For this sake, we generated data from 100 simulated participants from a range of typically observed DDM parameters. Our model was then fit to the true Bayesian posterior probability correct conditional on evidence, time and choice. LDC-predicted confidence almost perfectly correlated with the true probability of being correct (Spearman $r(999998) = .99$, $p < .001$). This close resemblance can be appreciated visually by comparing the model-based heatmap (created based on the estimated parameters; Fig 1B) to the heatmap based on the simulations (Fig 1A).

To further show that our model closely tracks the Bayesian readout of confidence, we tested its ability to reproduce statistical signatures that confidence should adhere to if it does reflect a Bayesian probability [1]. [1] identified three qualitative signatures, namely: confidence predicts choice accuracy, confidence increases with evidence strength for correct trials, but decreases with evidence strength for error trials (commonly called the folded X-pattern; [26,27]), and for any level of evidence strength above 0, high confidence trials should be linked with higher accuracy than low confidence trials [1].

Additionally, it is well known that confidence is negatively associated with the speed of responding [14,16,28,29]. This relationship is accounted for by the Bayesian readout of confidence under DDM assumptions (Eq (1)). The intuition behind this is that accumulation time informs on the difficulty of the decision and thus, on the probability of making a correct choice (if a lot of evidence is accumulated within a small amount of time, then one is likely to be in an easy trial, indicating high probability of making a correct choice). To test how well our model accounts for how time influences Bayesian confidence, we thus introduce a fourth signature,

namely average confidence for successive reaction time (RT) bins. As can be assessed on Fig 1C-1F, the simulated data showed an excellent fit to the signatures.

Empirically testing predictions of the LDC model

Having demonstrated that the LDC model can closely approximate the Bayesian readout of confidence on synthetic data, we next turned to empirical data from human participants. We tested two key predictions of the LDC model. First, the LDC model predicts that changes in confidence can be independent of performance. The two free parameters only describe how evidence and time are combined into a confidence judgment, but they do not affect the process that leads to specific levels of accumulated evidence and elapsed time. Any manipulation that selectively targets confidence while leaving performance unaffected should thus be captured by changes in α and/or β . A second novel prediction is that selective changes in each parameter of our model should lead to distinct modulations of confidence judgments. Thus, a manipulation targeting reliability (α) should lead to qualitatively distinct changes in confidence ratings compared to a manipulation targeting confidence bias (β).

Experiment 1: The LDC model accounts for performance-independent changes in confidence. We first tested a crucial prediction of our model, namely that changes in confidence can occur independent of changes in performance [9,30–32]. Although such dissociations have been observed since several decades (e.g., blindsight; 10), they pose a serious challenge for most current models of confidence. The LDC model naturally accounts for such dissociations. One particularly strong dissociation was observed in our recent work [19], in which a manipulation of participants' prior belief about their ability to perform a task selectively influenced their reported levels of confidence. In Experiment 1 of that paper, participants performed three perceptual tasks consecutively, each divided into a training and a testing phase (Fig 2). During the training phase, participants received feedback about their performance every 24 trials. Although participants were told that the feedback indicated how well they performed the task compared to a reference group, in reality the feedback was made up. Within each task, feedback indicated that performance was worse than of most other participants (*negative condition*); that it was on average (*average condition*); or that it was better than of most other participants (*positive condition*). During the testing phase, participants no longer received feedback; instead, they rated their confidence at the end of each trial. We observed a direct influence of the feedback manipulation on confidence, with more positive feedback leading to higher confidence, $F(2,47) = 16.65, p < .001, \eta^2 = .415, 95\% \text{ CI } [.193, .574]$. Importantly, this effect of feedback on confidence was not explained by objective performance, as RT and accuracy did not change as a function of feedback (accuracy: $X^2(2, N = 30987) = .30, p = .863, V = .003, 95\% \text{ CI } [0, .011]$; RT: $F(2,48) = 2.06, p = .14, \eta^2 = .079, 95\% \text{ CI } [0, .235]$).

We fitted the LDC model to the performance (accuracy and RT) and confidence reports in the test phase of this experiment, separately for each participant. LDC model predictions were then generated using the best fitting parameters for each individual. As can be seen in Fig 3, the LDC model provided an excellent fit to the data (see also S1 Fig for the observed and model predicted relationship between RT and confidence). Only confidence for incorrect trials at easier difficulty levels seemed to be overestimated by the model. This can be explained by the relatively small number of trials underlying these data points: 2.7% of all trials on average for the easy trial difficulty (i.e. less than 6 trials per feedback condition), and 5.3% for average trial difficulty (i.e. about 11 trials on average per feedback condition). Similar to the empirical data, feedback significantly influenced model-generated confidence ratings ($F(2,48) = 9.79, p < .001, \eta^2 = .290, 95\% \text{ CI } [.083, .465]$), but did not influence the performance data (RT: $F(2,48) = 1.19, p = .31, \eta^2 = .047, 95\% \text{ CI } [0, .185]$; Accuracy: $X^2(2, N = 30987) = .75, p = .69, V = .005,$

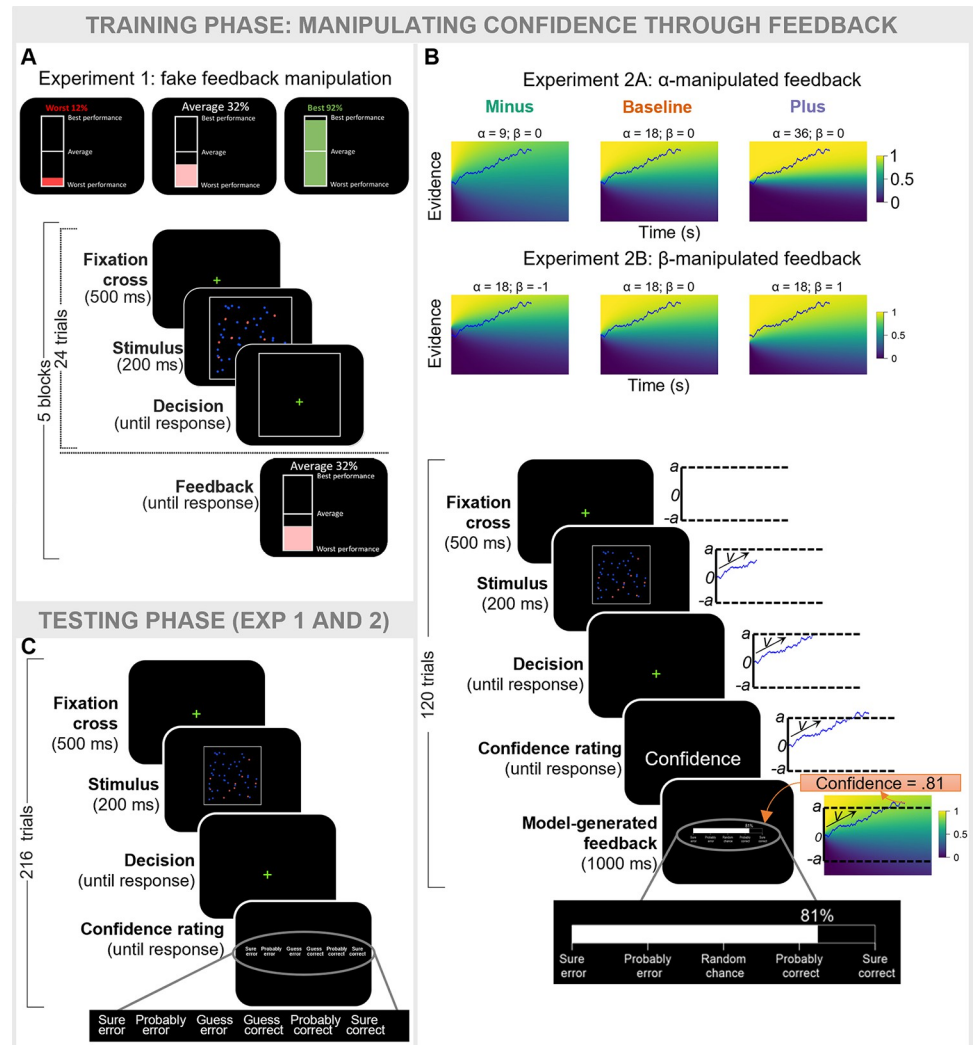


Fig 2. Experimental design. In both experiments, participants performed three different perceptual decision-making tasks (only one shown here). Each task started with a training phase during which a different feedback manipulation was induced. **A.** In Experiment 1, participants received fake feedback after each training block, framed as a comparison between their performance and the performance of a reference group. **B.** In Experiment 2, participants additionally rated their confidence before receiving trial-by-trial feedback reflecting their probability of making the correct choice. Unknown to participants, the feedback was actually generated by the LDC model behind the curtain. To do so, the evidence accumulation process for each trial was estimated using the mean drift rate and boundary from a previous pilot session (see [Methods](#) for full details). Feedback conditions differed in the α (resp. β) value used to generate feedback in Experiment 2A (resp. Experiment 2B). **C.** In both experiments, after each training phase participants completed a test phase during which they no longer received feedback but rated their decision confidence after each decision.

<https://doi.org/10.1371/journal.pcbi.1012273.g002>

95% CI [0, .014]). Thus, our model was able to capture the data pattern, namely that confidence reports can be influenced independently from behavioral performance.

We next investigated the estimated parameters of the model ([Fig 4](#)). Given that feedback selectively influenced confidence ratings, we expected a significant change in the confidence-specific parameters (i.e., α or β), but no variation in the DDM parameters (non-decision time, drift rate, decision threshold). Indeed, feedback had an influence on estimated α ($F(2,382) = 6.56, p = .002, \eta^2 = .033, 95\% \text{ CI } [.005, .073]$; [Fig 4A](#)) and β ($F(2,382) = 8.32, p < .001, \eta^2 = .042, 95\% \text{ CI } [.010, .085]$; [Fig 4B](#)). Tukey’s test for multiple comparisons found that estimated

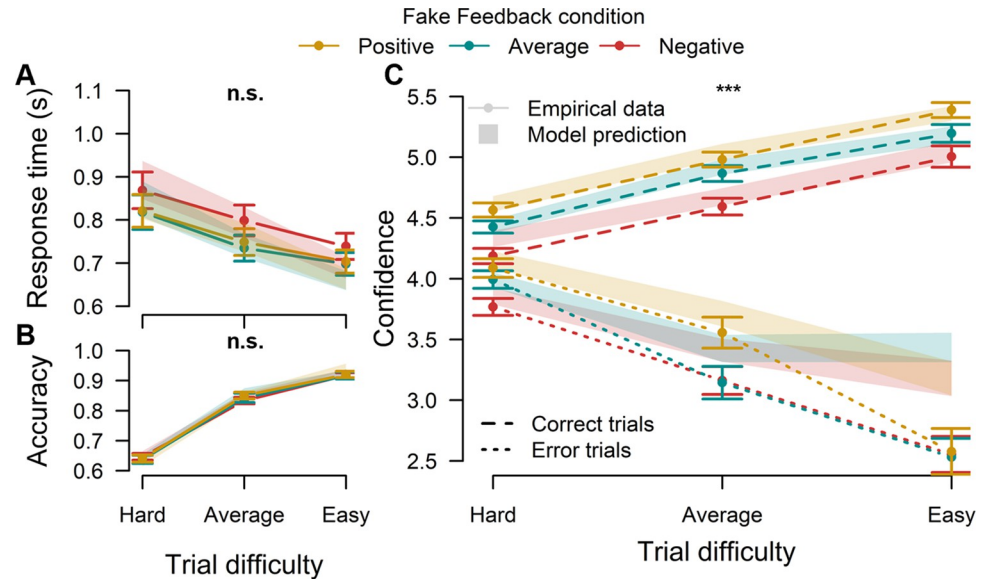


Fig 3. A key prediction of the LDC model is that confidence can vary independent from task performance. A-B. In Experiment 1, providing participants with fake feedback telling them their performance was better, equal or worse than a reference group indeed left RT (A) and accuracy (B) unaffected. C. On the other hand, fake feedback selectively influenced the reported level of confidence on correct trials. These results were closely captured by fitting the LDC model to these data. Note: Solid lines represent empirical data. Error bars represent standard error of the mean. Shades represent standard error of the mean of predictions of the LDC model. n.s. = $p > .05$; *** = $p < .001$. Significance indicators are about the effect of feedback conditions.

<https://doi.org/10.1371/journal.pcbi.1012273.g003>

α was lower in the negative condition than in the other two conditions (negative vs average: $p = .01$; negative vs positive: $p = .002$), whereas there was no difference in α between the average and positive conditions ($p = .88$). In a similar vein, β was higher in the positive condition compared to the other two (positive vs average: $p = .004$; positive vs negative: $p < .001$), whereas there was no difference between the negative and average conditions ($p = .83$). Finally, as expected estimated DDM parameters did not vary with feedback condition (drift rate: $F(2,48) = .18$, $p = .84$, $\eta^2 = .007$, 95% CI [0, .079], Fig 4E; non-decision time: $F(2,382) = .99$, $p = .37$, $\eta^2 = .005$, 95% CI [0, .002], Fig 4C) except for a minor effect on decision threshold ($F(2,382) = 3.30$, $p = .038$, $\eta^2 = .017$, 95% CI [0.00, .048]; Fig 4D). Post-hoc tests for the decision threshold revealed a slightly higher threshold in the negative condition compared to the positive condition and no difference with the other contrasts (negative—average: $p = .14$; negative—positive: $p = .04$; average—positive: $p = .85$).

Experiment 2: Dissociating parameter-specific effects on confidence ratings. Our next aim was to demonstrate that humans are sensitive to the specific parameterization of decision confidence proposed by the LDC framework. If confidence is computed using a low-dimensional solution, it should be possible to independently manipulate its parameters. Therefore, in a new set of two experiments, we aimed to induce selective changes in each parameter (reliability (α) or bias (β)) of the model.

The general design of both experiments was similar to Experiment 1: we manipulated the feedback during a training phase and investigated the impact of that manipulation on confidence ratings reported in a subsequent testing phase. Rather than presenting fake feedback every 24 trials, we adopted a novel approach where feedback during the training phase was presented after each trial in the form of a continuous value (Fig 2). Participants were told that this value reflected the probability that their response was correct (e.g., .8 vs .4 indicating that

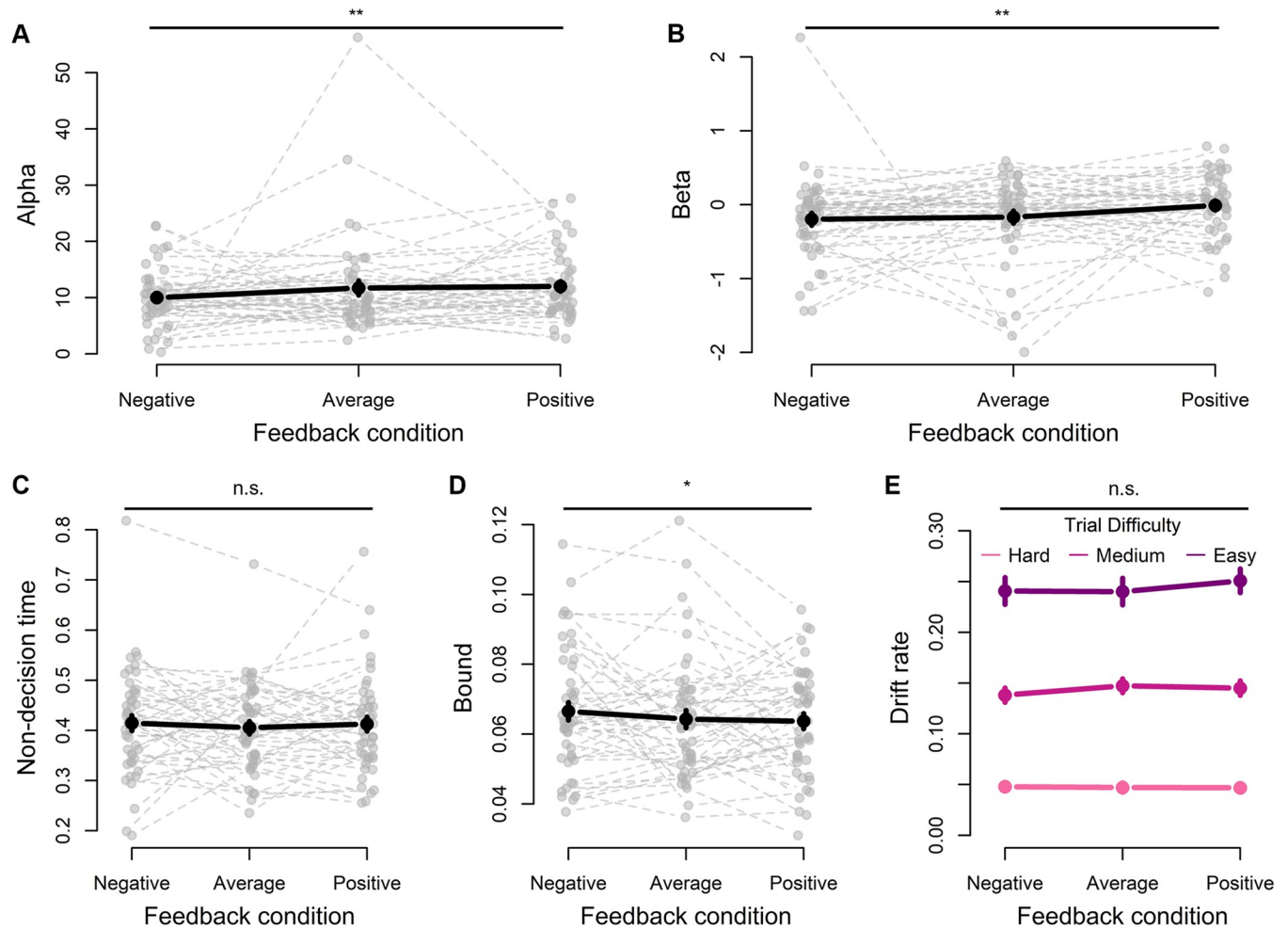


Fig 4. Best fitting parameters from model fits of Experiment 1. A-B The feedback manipulation influenced both α and β , with more positive feedback eliciting higher α and β . C-E The DDM parameters were not influenced by the feedback manipulation, except for a small significant effect on the bound. Grey dotted lines refer to individual fits. Note: Black lines are mean values. Error bars are SEM.

<https://doi.org/10.1371/journal.pcbi.1012273.g004>

there was a high vs low probability that they just made a correct choice). Unknown to participants the exact feedback value was generated by LDC behind the curtain (see [Methods](#) for full details). Both experiments comprise a baseline condition ($\alpha = 18$; $\beta = 0$) in which the feedback presented to participants reflected the model-approximated probability of a choice being correct. In Experiment 2A, the value of α that was used to generate the feedback was selectively manipulated between conditions. In addition to the baseline condition there was a minus condition where α was decreased ($\alpha = 9$), and a plus condition where α was increased ($\alpha = 36$). In Experiment 2B, the same procedure was used except that now the value of β was selectively manipulated between conditions ($\beta = -1$ in the minus condition and $\beta = 1$ in the plus condition).

A dissociable effect of manipulated feedback on confidence according to the parameter manipulated. Similar to Experiment 1, we expected participants in Experiment 2A and 2B to adapt how they compute confidence depending on the feedback received during the training phase. Given that the feedback was generated in the training phase by manipulating specific parameters in the LDC model, we expected that participants would learn to compute confidence using the same parameters settings in training and subsequent test phase. As previously

described, the reliability parameter α reflects how strongly individuals weigh evidence in their computation of confidence. Given that accuracy is closely related to the amount of available evidence, correct trials tend to have considerable supporting evidence when reporting confidence, whereas error trials usually have little to no supporting evidence. Given that α weighs evidence, a decrease (in the α -minus condition) or an increase (in the α -plus condition) of α is therefore expected to differently impact confidence for correct trials (strong influence) than for error trials (little to no influence). In contrast, the parameter β reflects a stimulus-independent confidence bias, so providing participants with β -manipulated feedback is expected to lead to changes in confidence irrespective of choice accuracy. The reasoning for this prediction is that β is not concerned with the evidence provided by the stimulus (nor by the response), as it simply adds (in the β -plus condition) or subtracts (in the β -minus condition) a constant to the (logit of the) confidence judgment regardless of what happens during the trial.

These intuitions are further illustrated in Fig 5A and 5B, which show the actual (i.e., manipulated) feedback that was presented to participants during the training phase of our experiments. Confirming the above intuition, there was an interaction in Experiment 2A between accuracy and the value of α ($F(2,12623) = 76.73, p < .001, \eta^2 = .012, 95\% \text{ CI } [.008, .016]$): feedback was more positive when generated by a higher α when considering correct trials only ($F(2,9911) = 723.77, p < .001, \eta^2 = .127, 95\% \text{ CI } [.116, .139]$), but did not change when considering error trials only, $F(2,44) = 1.41, p = .25, \eta^2 = .060, 95\% \text{ CI } [0, .212]$. For Experiment 2B, although there was a significant interaction between β value and accuracy ($F(2,10589) = 43.11, p < .001, \eta^2 = .008, 95\% \text{ CI } [.005, .012]$), feedback was more positive when generated by a higher β both when taking corrects ($F(2,8420) = 2260.84, p < .001, \eta^2 = .349, 95\% \text{ CI } [.334, .364]$) and errors ($F(2,35) = 183.56, p < .001, \eta^2 = .913, 95\% \text{ CI } [.852, .943]$) separately.

Behavioral results. We now turn to the effects of the feedback manipulations on the testing phase data. The results concerning task performance were as expected: we found no effect of feedback condition on performance (RT and accuracy) in the testing phase of Experiment 2A (RT: $F(2,40) = .15, p = .86, \eta^2 = .008, 95\% \text{ CI } [0, .086]$; Accuracy: $X^2(2, N = 25739) = 3.54, p = .17, V = .012, 95\% \text{ CI } [0, .023]$) and Experiment 2B (RT: $F(2,33) = .507, p = .61, \eta^2 = .030, 95\% \text{ CI } [0, .176]$; Accuracy: $X^2(2, N = 21939) = 1.85, p = .40, V = .009, 95\% \text{ CI } [0, .021]$). There was, however, the expected effect of trial difficulty on performance both in Experiment 2A (accuracy: $X^2(2, N = 25739) = 1023.00, p < .001, V = .199, 95\% \text{ CI } [.162, .214]$; RT: $F(2,25619) = 164.34, p < .001, \eta^2 = .013, 95\% \text{ CI } [.010, .015]$) and Experiment 2B (accuracy: $X^2(2, N = 21939) = 756.53, V = .186, 95\% \text{ CI } [.151, .200], p < .001$; RT: $F(2,21831) = 168.69, p < .001, \eta^2 = .015, 95\% \text{ CI } [.012, .019]$), with lower accuracy and higher RT when trial difficulty was higher (all post-hoc comparisons: $p_s < .02$). There was no interaction between feedback condition and trial difficulty on RT and accuracy in either Experiment 2A or 2B (all $p_s > .05$).

After demonstrating that the feedback did not influence task performance itself, we next turn towards confidence ratings (see also S1 Fig for an analysis of the relationship between RT and confidence). In line with the feedback presented during the training phase (Fig 5A and 5B), the data of the testing phase revealed that α -manipulated feedback in Experiment 2A had an effect on confidence ratings within correct trials ($F(2,39) = 4.86, p = .01, \eta^2 = .199, 95\% \text{ CI } [.011, .396]$), but not within error trials ($F(2,45) = .87, p = .43, \eta^2 = .037, 95\% \text{ CI } [0, .169]$; Fig 5C). Note that this finding should be interpreted with caution, since the interaction between accuracy and feedback was not significant ($F(2,44) = .62, p = .54, \eta^2 = .028, 95\% \text{ CI } [0, .151]$). Turning to Experiment 2B, in line with the predictions there was an effect of feedback condition on confidence ratings in both correct trials ($F(2,33) = 8.86, p < .001, \eta^2 = .348, 95\% \text{ CI } [.088, .544]$) and in error trials ($F(2,36) = 4.28, p = .02, \eta^2 = .190, 95\% \text{ CI } [.003, .392]$; Fig 4D). Here again, no interaction between accuracy and feedback condition was found ($F(2,35) = .29, p = .75, \eta^2 = .016, 95\% \text{ CI } [0, .132]$). Lastly, trial difficulty had an effect on confidence ratings

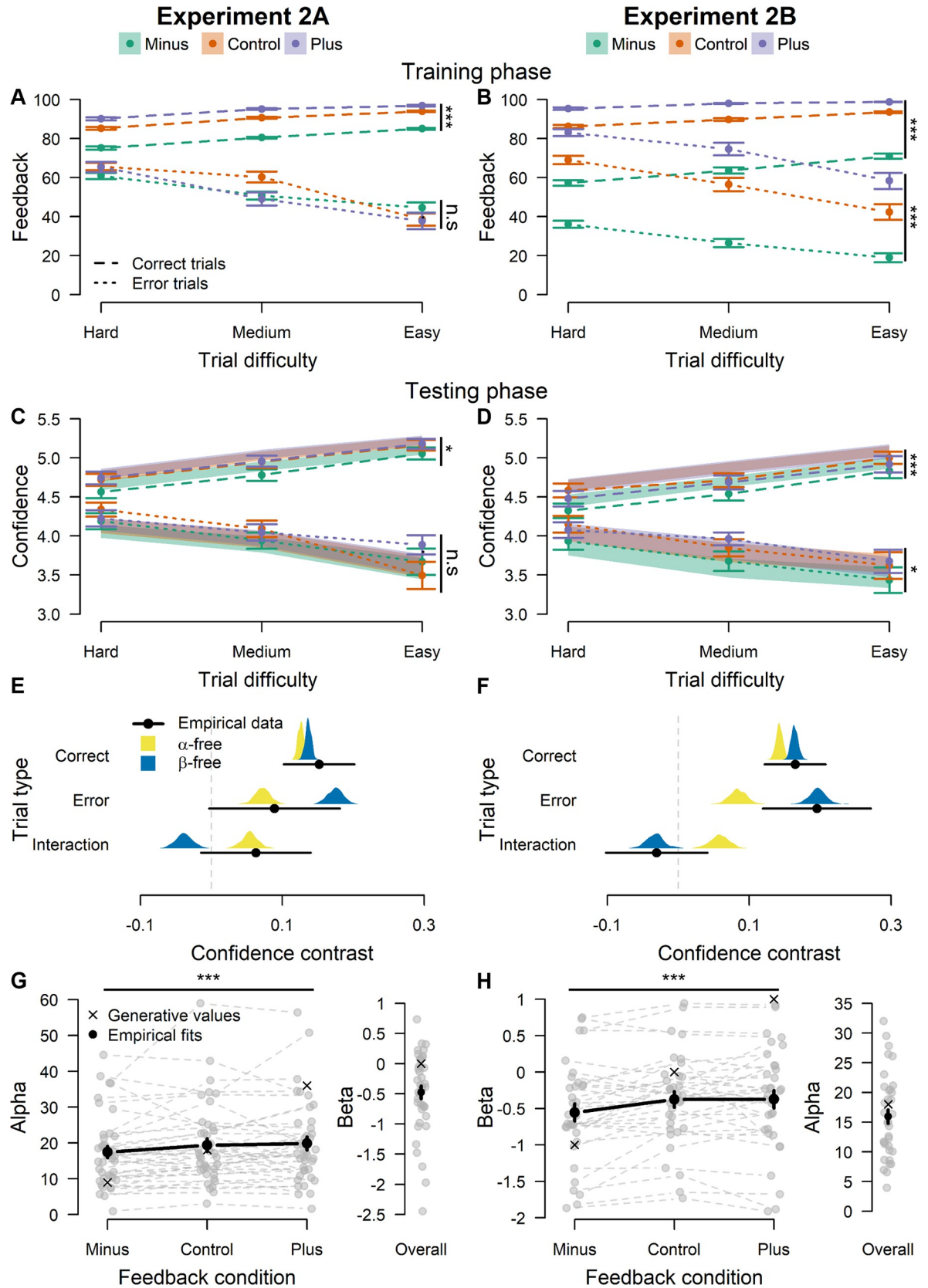


Fig 5. A key prediction of the LDC model is that participants should be sensitive to the specific parametrization of confidence proposed by the model. To test this, Experiment 2 provided participants with probabilistic feedback generated by the LDC model. Critically, LDC based feedback was generated using different levels of α or different levels of β . **A.** Changing α influences the confidence for correct trials but not for errors. **B.** Changing β influences feedback for both corrects and errors. The pattern that we saw in the feedback (which effectively are our predictions) was also seen in the behavioral data. **C.** α -manipulated feedback influenced confidence reports for correct but not error trials. **D.** β -manipulated feedback influenced confidence reports on both correct and error trials. **E.** Fitting the LDC model to the empirical data of Experiment 2 revealed that data in the α -manipulated feedback was best explained by a model in which α was allowed to vary. **F.** Data from the β -manipulated feedback was best explained by a model in which β was allowed to vary. To visualize this, we computed confidence contrasts for the empirical data (black lines), as well as for the α -free (yellow distribution) and β -free (blue distribution) model fit, separately for corrects and errors. “Interaction” refers to the difference between the confidence contrast in corrects and errors. **G.** Parameter estimates from the α -free model. Estimated α was higher when feedback was generated by higher α . **H.** Parameter estimates from the β -free model. Estimated β was higher when feedback was generated by higher β . As a reference, the values of α and β used to generate feedback are added as crosses in panels G and H. Note: grey dots and lines refer to individual estimates. Black dots correspond to sample means, distributions correspond to the bootstrapped mean predicted confidence contrasts. Error bars and shaded areas represent empirical and model-simulated SEM, respectively. n.s. = $p > .05$; * = $p < .05$; *** = $p < .001$

<https://doi.org/10.1371/journal.pcbi.1012273.g005>

in both Experiment 2A ($F(2,25633) = 75.21, p < .001, \eta^2 = .006, 95\% \text{ CI } [.004, .005]$) and 2B ($F(2,253) = 33.49, p < .001, \eta^2 = .209, 95\% \text{ CI } [.0125, .292]$). We found no interaction between trial difficulty and feedback condition (Experiment 2A: $F(4,25625) = 2.37, p = .05, \eta^2 < .001, 95\% \text{ CI } [0, .001]$; Experiment 2B: $F(4,20385) = 1.70, p = .15, \eta^2 < .001, 95\% \text{ CI } [0, .001]$). Overall, these results corroborate the predicted pattern and show a clearly dissociable effect of feedback on confidence ratings according to the parameter manipulated in the feedback generation.

LDC model fits. We next performed model comparison to explore whether the different patterns of confidence ratings observed in Experiment 2A and 2B would be best explained by a change in the targeted parameter (i.e. a change in α in Experiment 2A and a change in β in Experiment 2B). Two candidate LDC models were fit to the accuracy, RT and confidence data of both experiments. Each model differed in whether α or β was fixed between feedback conditions: in the α -free model, only α was allowed to vary between feedback conditions, whereas in the β -free model, only β was allowed to vary between feedback conditions. As recommended in [33], we investigated how well simulations from the best-fitting parameters from both the α -free and the β -free models were able to reproduce the observed behavioral effects. Specifically, we defined a confidence contrast that captured the qualitative signatures seen in the feedback presented. Since the difference in feedback between the baseline and the plus condition was negligible relative to how both conditions differed from the minus condition in both experiments, we computed our confidence contrast as average confidence in the minus condition subtracted from average confidence in the baseline and the plus condition. Fig 5E and 5F show the empirical confidence contrast as well as the distribution of the mean predicted confidence contrast for both the α -free and the β -free model obtained via bootstrapping. In Experiment 2A, the confidence contrasts predicted by both the α -free and the β -free model was highly similar for correct trials, and both matched well to the empirical data. However, while the α -free model closely captured the confidence contrast in errors and hence the interaction, the β -free model overestimated the effect in errors, which led it to underestimate the interaction. Similarly, in Experiment 2B, both models accurately captured the empirical confidence contrast in correct trials. Additionally, the β -free model precisely reproduced both the empirical confidence contrast in error trials and the interaction, whereas the α -free model clearly underestimated the confidence contrast in error trials, which led to predicting an interaction that was not present in the empirical data. We conclude that the α -free (resp. β -free) model fits best to the experiment where α (resp. β) was manipulated.

To further confirm that the α -free (resp. β -free) model is the most likely to explain the results of Experiment 2A (resp. 2B), we additionally quantified the goodness-of-fit of each model using Bayesian information criterion (BIC). Four additional candidate models were

included in that analysis. To ensure that the feedback effects were best captured by a change in one parameter only, we included a null model where neither α nor β varied between feedback conditions and a full model where both α and β were free to vary between conditions. To investigate whether a pure Bayesian readout could capture such data, we also included two models computing the exact probability of making a correct choice given accumulated evidence and time, based on the estimated drift rates. In the first Bayesian readout model, the drift rates were fixed between conditions (just like with the other LDC candidate models). Importantly, the estimated DDM parameters for the LDC models and this Bayesian readout model were identical. These models therefore only differ in how inputs from post-decisional DDM (i.e. post-decision evidence and post-decision RT) are converted into confidence, and can only be distinguished on the goodness of their fit to confidence ratings. Since fixing drift rates between feedback condition effectively prevents the Bayesian readout model from returning any difference in confidence between conditions, we included a second Bayesian readout model where drift rates were allowed to vary between feedback conditions. [Table 1](#) reports the difference in mean BIC across participants of each candidate model compared to the best model, separately for both experiments. A first conclusion that can be drawn, is that all LDC models considerably outperformed the Bayesian readout models. Even the Null LDC model, effectively blind to the feedback effects on confidence, showed better BIC than the Bayesian readout model with varying drift rates between feedback conditions (in theory able to account for feedback effects). Second, both the α -free and β -free model outperformed the null model (i.e., providing strong evidence for a change in the parameters) as well as the full model (i.e., providing strong evidence for a *selective* change in the parameters). Third, as expected the α -free model showed the lowest BIC for the data of Experiment 2A. Surprisingly though, the α -free model also slightly outperformed the β -free model in Experiment 2B. Overall, the difference in BIC between the α -free and the β -free models appears marginal compared to how strongly they each outperformed the null and full models. Additionally, the difference in BIC between the α -free and the β -free models was bigger in Experiment 2A ($\Delta_{BIC} = 4.15$), where the α -free model was expected to be the best performing model, compared to the difference observed in Experiment 2B ($\Delta_{BIC} = 2.54$). Applying categorical cutoffs to describe the magnitude of the evidence in favor of the α -free model in both experiments, such as the rule of thumb proposed by [34], leads to conclude that the α -free model has considerably more support than the β -free model in Experiment 2A, but only weak support in Experiment 2B.

We also did group-level Bayesian model selection to determine which model better accounts for the data [35,36]. Using BIC weights as model evidence, we performed two different analyses: in the first one, we included all four LDC candidate models. Second, since the α -free and β -free models did not differ much in mean BIC and both clearly outperformed the other candidate models, we also looked at the results of group-level Bayesian model selection with only these two models included in the analysis. In this analysis framework, candidate

Table 1. Model comparison expressed in distance in mean BIC across participants from the best-fitting model.

Model	Δ_{BIC}	
	Experiment 2A	Experiment 2B
Null	21.02	23.81
α -free	0	0
β -free	4.15	2.54
Full	19.58	19.42
Bayesian readout (fixed drifts)	113.56	108.43
Bayesian readout (free drifts)	134.63	137.14

<https://doi.org/10.1371/journal.pcbi.1012273.t001>

models are treated as random effects with fixed (unknown) distribution in the population (i.e. the data of all participants is no longer assumed to originate from the same model). Therefore in the following results, we reason in terms of model frequency within the population. The Bayesian Omnibus Risk (BOR) was consistently low when all four LDC candidate models were included (Experiment 2A: BOR = .005, Experiment 2B: BOR = .027), suggesting that one model was more frequent than the others. We next looked at the maximum Exceedance Probability (maxEP) to further investigate which model was found to be more frequent than the others. Consistent with the results from mean BIC, the α -free model had the highest EP in both experiments (Experiment 2A: maxEP = .996, Experiment 2B, maxEP = .987). Interestingly, when only the α -free and β -free models were included in the analysis, the BORs were much higher in both experiments (Experiment 2A: BOR = .365, Experiment 2B: BOR = .253), indicating that the models have a significant chance of being equally frequent. This last result further highlights the difficulty of drawing a clear conclusion from such small differences in mean BIC between both models. In light of this, and following the recommendations of [33], it appears more relevant to select the winning model according to its ability to reproduce the observed behavioral effect of interest (namely, the confidence contrasts). Still, it remains an open question for future work to determine why exactly the α -free model obtained a relatively better fit overall.

Lastly, to verify that the model-validated effects observed in Experiment 2A and 2B are caused by increased α (resp. β) values in the conditions where feedback was generated from higher α (resp. β) values, we looked at the estimated α and β from the α -free model fits in Experiment 2A and the β -free model fits in Experiment 2B (Fig 5G and 5H). As expected, there was an effect of feedback condition on estimated α in Experiment 2A, $F(2,318) = 16.56$, $p < .001$, $\eta^2 = .09$, 95% CI [.04, .16], with lower estimated α in the minus condition compared to the control and plus conditions ($ps < .001$), but no significant difference between the control and the conditions ($p = .55$). Similarly in Experiment 2B, estimated β was found to be different across feedback conditions, $F(2,270) = 26.44$, $p < .001$, $\eta^2 = .16$, 95% CI [.09, .24], with again lower estimated β in the minus condition compared to the control and plus conditions ($ps < .001$) but no difference between the control and the plus conditions ($p = .99$). Taken together, these results suggest that theoretically motivated confidence manipulations can lead to specific and theoretically predicted changes in confidence.

Experiment 3: Feedback manipulations influence confidence in a transfer task. In Experiment 1 and 2, trials in the training and testing phase were identical (i.e. same task, same difficulty levels). Therefore, a final potential concern is that instead of changing their mapping of confidence (i.e. changed their α and β), participants might instead have learned to give a specific confidence rating for the trials they were trained on. We resolved this ambiguity by testing whether the effect of feedback can also be observed when trials in the testing phase are perceptually different from the trials in the training phase. In Experiment 3, participants were constantly trained on either the letter discrimination task or the dot color task where they received β -manipulated feedback (i.e. the β -minus and β -plus conditions of Experiment 2B) and then tested on either the same task (Same condition) or a different task (Transfer condition).

In the Same condition, there was a main effect of feedback condition on confidence ratings ($F(1, 74) = 7.96$, $p = .006$, $\eta^2 = .10$, 95% CI [.01, .24]), replicating the findings of Experiment 1 and 2. More importantly, a main effect of feedback condition on confidence was also found in the Transfer condition ($F(1, 50) = 4.64$, $p = .036$, $\eta^2 = .09$, 95% CI [0, .26]). Like in Experiments 1 and 2, manipulated feedback influenced neither RT nor accuracy, in both the Same and the Transfer conditions (Same RT: $F(1,34) = .82$, $p = .371$, $\eta^2 = .02$, 95% CI [0, .20], Accuracy: $F(1, N = 24352) = 3.10$, $p = .078$, $V = .011$, 95% CI [0, .024]; Transfer RT: $F(1,34) = .04$, $p = .849$, η^2

= .001, 95% CI [0, .10], Accuracy: $\chi^2(1, N = 24412) = 0.29$, $p = .587$, $V = .003$, 95% CI [0, .016]). These results invalidate the alternative explanation that individuals only learn to give different confidence ratings for the specific trials they were trained on. Instead, they generalize this learning even to different tasks.

Discussion

How to incorporate the sense of confidence in models of decision-making has been the focus of much recent work. An influential framework is based on the Bayesian interpretation of confidence [3,37–39], namely that confidence reflects the probability of being correct given both accumulated evidence and elapsed time [14,15,17]. In order to accurately compute this probability, it is necessary to know how to compute confidence based on the agent's available data (evidence and time). Currently, a computationally plausible account describing how individuals learn this mapping is lacking. In the current work, we introduced the LDC model, which provides a tractable and flexible account of decision confidence. Using simulations, we first showed that LDC provides a highly reliable approximation of the true probability correct. Fitting this model to empirical data revealed that LDC accounts very well for human confidence ratings. Critically, using a novel feedback manipulation, we validated two key predictions from the model, namely that changes in confidence can be independent of performance, and that independently manipulating the reliability (α) and bias (β) parameters elicit clearly dissociable and identifiable effects on confidence.

Introducing tractability and flexibility to decision confidence modelling

The LDC model belongs to the family of DDM-based models of decision confidence. Here, confidence is conceptualized as a (Bayesian) readout of the probability of a correct choice given evidence, time and choice. Existing models following that approach have been successful in explaining a wealth of data, including the link between confidence and RT [14,17], and deviations from accuracy through the contribution of priors [18,19]. Estimating the probability correct based on the available data, however, is computationally intractable. The LDC model therefore proposes to approximate the Bayesian readout with a logistic function, offering a tractable approach of how humans compute confidence. Importantly, even though it describes an approximation to an optimal solution, the LDC model follows the same principles as the optimal Bayesian readout and uses the same information to compute confidence. It thus fundamentally differs from heuristic models that describe the computation of confidence as based on performance-unrelated cues [40]. On another note, while the LDC model solves the problem of estimating confidence for each point on the (evidence, time) plane, it leaves open how agents know about the appropriate parameter settings (i.e. how they learn the correct mapping). Our behavioral results suggest that the parameters might be learned from external feedback, this perspective should be investigated further in future works by looking at the dynamics of this potential learning process.

To increase flexibility and account for deviations from optimality, the LDC model relies on two free parameters, which control the reliability of evidence (α) and bias (β) in the computation of confidence. A different class of confidence models that can account for biases and deviations between confidence and accuracy is based on SDT framework [41–46]. For that purpose, these models typically either assume the existence of metacognitive noise [43–45], and/or consider that confidence is not entirely derived from the same signal as the primary decision [41–44,46]. A recent study comparing the different SDT models of confidence on simple perceptual tasks proposed that confidence is simply computed as a noisy readout of the evidence used for the primary decision [47]. Although the LDC model is grounded within the

DDM tradition which conceptualizes confidence as the Bayesian probability correct, it does not critically hinge upon the specifics of the DDM. It would be straightforward to construct a simplified version of the LDC model which ignores the element of time. This would allow to directly compare the LDC approach to recent SDT models of confidence. Crucially, with its parameters, our model can flexibly account for the different types of idiosyncrasies, biases and deviation from the optimal Bayesian readout [22–25], which are all merged into a single meta-cognitive noise parameter in most SDT frameworks.

Confidence can vary independently from task performance

In all Experiments, we observed that decision confidence was influenced by the feedback manipulation, whereas objective performance was not. This finding rules out an interpretation whereby the feedback influenced task performance, and changes in confidence simply reflect this change in performance. Indeed, some previous work has shown that changes in confidence can be explained by subtle differences in RT [14,48]. This was not the case in the current experiments because both accuracy and RT were not influenced by the feedback manipulations. As such, it is unlikely that pure Bayesian read-out models can account for the effects observed in the current work, as they do not allow for confidence-specific parameters [14–16]. In contrast, LDC accurately captured the effect of feedback on confidence in the absence of changes in objective performance, thus attesting to the flexible nature of the LDC model. Previous studies have unraveled several other factors that influence the reported level of decision confidence, while leaving task performance unaffected, for example emotional states [30,49], working memory content [32] and age [50,51]. Besides, dissociations between performance and metacognition have long been reported in cases such as blindsight [10,52], where individuals with lesions in primary visual cortex show above chance level performance at visual tasks despite reporting no awareness of the stimuli. The opposite pattern of low performance linked with high confidence has also been observed. Change blindness (i.e. failure to detect major differences between two images while they flicker off and on) is a typical example of such meta-cognitive error where individuals believe they would be able to detect major changes, despite being unable to do so [11]. These examples highlight how ubiquitous dissociations between performance and metacognition are. By incorporating free parameters controlling for evidence reliability and bias into the computation of confidence, the LDC model is in principle flexible enough to account for such dissociations.

Humans can independently tune evidence reliability and bias in confidence

In Experiment 2A and 2B, we aimed to selectively manipulate confidence ratings according to each parameter of the LDC. By providing model-generated feedback from different α 's in Experiment 2A and different β 's in Experiment 2B, we revealed clearly distinct patterns of confidence ratings according to the parameter manipulated. Moreover, the empirically observed patterns were best captured by models where the manipulated parameter was set as a free parameter (e.g. α -free model when feedback was α -manipulated). One might argue that the magnitude of the observed effects of feedback on confidence was on a significantly smaller scale than the actual differences in feedback between conditions. Similarly, the across conditions difference between the best-fitting parameters was on a much smaller scale than the differences in feedback-generating values. This mainly shows that participants do not come to the experiments with a blank slate, and do not use external feedback only to adjust their estimation of confidence. Instead, participants have their own prior beliefs on their performance at the start of a new task, and partially use external feedback in combination with internally-gathered information about the task to update how they compute confidence (similar to how

individuals integrate both internal and external forms of feedback to adjust their performance, [53]). Given the number of trials in the training phase used to induce a change in confidence (120), and the subtlety of the feedback manipulations, it was actually non-trivial that any change in the fitted values of α and β was observed.

These results imply that individuals can change their computation of confidence consistently with our parameterization of confidence, providing strong validating evidence in favor of LDC. This observation raises the intriguing possibility that individuals might exert control over the parameters governing the computation of confidence in a way that maximizes utility. Intuitively, computing confidence in such a way that it closely matches the Bayesian readout seems like the rational strategy to optimize utility, as it would allow to optimize behavior based on the best possible internal evaluation of that behavior [5,7,9]. In some contexts, however, other factors than informativeness play a role in the utility of confidence. When competing for shared limited resources, expressing overconfidence plays a key role in convincing other agents not to compete for the resource (i.e. “bluffing”; [54,55]). Errors caused by overconfidence, though, bear a high cost in such strategy. In such a context, the optimal way to compute confidence seems to be an increase in the evidence reliability estimate (α), which will lead to higher confidence for scenarios with much evidence (i.e., overconfidence when you are likely to win the competition) but lower confidence for scenarios with little evidence (i.e., when you are likely to lose the competition). Increasing β in this scenario is likely suboptimal because this produces overall high confidence, also for scenarios with little evidence. The opposite scenario might be true in a social decision-making context. If confidence is used to assert influence rather than to convey accuracy [56], the optimal strategy might be an overall increase in β , resulting in general overconfidence (i.e. irrespective of accuracy) to push forward one’s choice. These examples show that what is traditionally treated as deviations from the optimal Bayesian readout can sometimes be considered as optimal through the lenses of utility maximization.

Beyond dichotomies with model-informed feedback

In contrast with the binary “correct/error” feedback typically provided in lab experiments, feedback received in daily life is not always clear-cut. Individuals must often make sense of noisy and probabilistic feedback cues (e.g. how should a street-artist interpret a subtle nod from a spectator?). Continuous feedback has been used in the past to communicate performance relative to other (hypothetical) participants [19,57] or to give average accuracy over several past trials [58,59]. However, in the current work we designed a novel feedback manipulation which provides continuous feedback about choice accuracy on a trial-by-trial basis. It is important to note that our instructions simply stated that feedback would reflect the probability of being correct on a single trial, without much more explanation as to how this proportion was calculated. A skeptical participant could reasonably doubt the trustworthiness of the feedback, since it might seem unlikely that we provide an “accurate” probability of being correct on a single trial basis (e.g. is a feedback of 80% vs 70% really informative, or are the values pure noise added by the experimenter). Despite these potential obstacles, our feedback manipulation did produce the confidence patterns we predicted, hence validating our model-generated feedback approach. This nuanced way of providing feedback goes beyond the mere distinction between dichotomous valid versus invalid feedback [60], and offers a promising framework to control the level of ambiguity and informativeness of trial-by-trial feedback, allowing to study in a more fine-grained manner how individuals process and are impacted by more realistic, ambiguous feedback [61,62].

Interpreting the LDC parameters

An appealing property of computational models is that their parameters often have clear interpretations, and can be selectively manipulated [13,63], although it is subject of recent debate [64]. Similarly, in LDC, evidence and time are mapped onto confidence by means of a reliability parameter, α , and a confidence bias parameter, β . Our reliability parameter, α , can be interpreted as an individual's estimate of the precision of evidence. This interpretation is similar to the recently proposed concept of "meta-uncertainty", which is described as "the subject's uncertainty about the uncertainty of the variable that informs their decision" [65]. In both the LDC model and [65]'s CASANDRE model, one's estimate of evidence reliability weighs how evidence is used to compute confidence. Note that an important difference is that in CASANDRE the estimate is assumed to be correct on average (i.e. individuals are assumed to have an uncertain, but on average correct estimate of evidence reliability), whereas one of the key points of the LDC model is that participants can have incorrect values of α .

The second parameter of LDC, β , globally increases or decreases confidence. Interestingly, β is scaled by $1/\sqrt{t}$. Another possibility for our model could be to have a time-independent confidence bias parameter. Such hypothetical parameter would become more important in the computation of confidence with time relative to evidence, which is scaled by $1/\sqrt{t}$. In contrast, β reflects the same constant shift in confidence for all accumulation time. It therefore straightforwardly relates to the metacognitive bias described in other models of confidence that ignore RTs [66]. It is interesting to note that β has no lower boundary, and therefore in theory it allows for cases where agents with a very low (negative) β would constantly judge that they made an error. While this scenario may seem irrational, some individuals may indeed be so insecure about their decisions that they would often rate their confidence below the guess threshold. It remains unclear whether such behavior occurs empirically. For now, we simply note that the LDC model can in principle account for it, potentially opening new doors to understand metacognitive impairments in future work [23–25,45].

In light of this interpretation of α and β , one can further interpret specific patterns in the data. For example, in Experiment 1, we observed a change in α in response to negative feedback (with a significantly lower estimated α compared to the other two conditions), indicating that participants judged evidence as less reliable after receiving negative feedback. On the contrary, we observed a change in β after positive feedback (with a significantly higher estimated β compared to the other two conditions), suggesting a general overconfidence bias after receiving positive feedback. This dissociation suggests that despite similar effects at the behavioral level, the LDC model allows to further tease apart the origins of confidence biases e.g. in response to positive vs negative feedback. It will be interesting to investigate in future work how experimental factors influencing confidence can be differentially explained in terms of evidence reliability (α) or general under- or overconfidence (β).

Finally, we note that in the current parameterization of confidence, identical to the Bayesian readout, confidence always depends on \sqrt{t} . However, the influence of time on confidence might vary according to the task or individual. To account for such hypothetical sources of variability, one could expand the LDC model by further parameterizing the influence of time with a third parameter, γ , and replace \sqrt{t} in Eq (3) with t^γ . The model then has an accurate calibration of how time influences confidence when $\gamma = 0.5$, and overweighs (resp. underweighs) time in the computation of confidence when $\gamma > 0.5$ (resp. $\gamma < 0.5$). By doing so, future work might investigate whether variability in the relation between confidence and decision time can be captured by the extended LDC model.

Conclusion

We introduced the LDC model, a new model of decision confidence that offers a tractable and flexible approximation of confidence as the Bayesian probability of making the correct decision. The model provides a low-dimensional parametrization of decision confidence which allows efficient estimation of confidence, while at the same time accounting for idiosyncrasies and different kinds of confidence biases. This parameterization of confidence was validated in two experiments showing a distinct pattern of confidence ratings after specifically manipulating the mapping according to each parameter of the model.

Methods

Experiment 1

Ethics statement. All procedures were approved by the Katholieke Universiteit Leuven (KU Leuven) ethics committee.

Participants. Fifty participants (eight men, one third gender, age: $M = 19$, $SD = 4.9$, range 17–52) took part in Experiment 1 (two excluded due to chance level performance). All participants participated in return for course credit and read and signed a written informed consent at the start of the experiment. Detailed methods and analyses for Experiment 1 have already been reported in [19]. We briefly report the general procedure here.

Procedure. Participants completed three decision-making tasks: a dot color task, a dot number task and a letter discrimination task. Each task started with 120 training trials. Feedback during training was presented at the end of blocks of 24 trials. Unknown to participants, feedback was predetermined to be either good, average or bad for a specific task, and feedback scores were randomly sampled according to the feedback condition. Each participant received good feedback on one task, average feedback on another task, and bad feedback on a third task (order and mapping with tasks counterbalanced between participants). After the training phase of a task, participants performed 216 test trials where feedback was no longer provided. Instead, confidence ratings were queried at the end of each trial. For each task, there were three levels of stimulus difficulty (easy, average, or hard).

Dot color task. On each trial, participants decided whether a box contained more (static) blue or red dots. The total number of dots was always 80, with differing proportions of red or blue dots depending on the difficulty condition. The position of dots was randomly generated on each trial.

Dot number task. On each trial, two boxes were presented, one of which contained 50 dots and the other more or less than 50 dots. Participants decided which of the two fields contained the largest number of dots. The exact number of dots in the variable field differed depending on the difficulty condition. The position of dots was randomly generated on each trial.

Letter discrimination task. On each trial, participants decided whether a field contained more X's or O's. The total number of X's and O's was always 80, with differing proportions of X's or O's depending on the difficulty condition. The position of the letters was randomly generated on each trial.

Experiment 2

Ethics statement. All procedures were approved by the KU Leuven ethics committee.

Participants. Forty-three participants (8 men, 35 women, age: $M = 18.49$, $SD = 1.03$, range 16–22) took part in Experiment 2A. Forty-two participants (9 men, age: $M = 18.83$, $SD = 2.05$, range 17–29) took part in Experiment 2B. Due to chance performance on at least one of the tasks, we removed 3 participants from Experiment 2A and 3 participants from

Experiment 2B. Five additional participants were removed from Experiment 2B due to (almost) no variability in their confidence reports (i.e. used the same report on more than 90% of the trials). Two participants passed our a priori inclusion criteria but showed behavior that could indicate a misunderstanding of the confidence scale (i.e. chance-level performance or higher for all confidence ratings). Removing these two participants from the analysis did not significantly change the behavioral results. All participants took part in return for course credit and signed informed consent at the start of the experiment.

Stimuli and apparatus. All experiments were conducted on 22-inch DELL monitors with a 60 Hz refresh rate, using PsychoPy3 [67]. All stimuli were presented on a black background centered on the middle of the screen (radius 2.49° visual arc). Stimuli for the dot number task were presented in two equally sized boxes (height 20°, width 18°) at an equal distance from the center of the screen. Stimuli for all other tasks were presented in one box (height 22°, width 22°).

Procedure. In both experiments, participants completed three decision-making tasks: a dot color task, a shape discrimination task and a letter discrimination task. Each task started with 108 training trials. After each choice, participants rated their confidence level and then received (continuous) feedback about their performance. After the training phase of a task, a test phase of 216 trials followed which was identical to the training phase, except that feedback was omitted. Every trial was assigned one of three possible difficulty levels. The difficulty levels were matched between the three tasks based on a pilot staircase session. For all tasks, a trial started with a fixation cross that was presented for 500 ms, after which the stimulus appeared for 200 ms or until a response was given. Participants indicated their choice using the C or N key using the thumbs of both hands. There was no time limit for responding, although participants were instructed to respond as fast and accurately as possible. After each choice, participants rated their confidence on a 6-point scale, labeled from left to right: 'sure error', 'probably error', 'guess error', 'guess correct', 'probably correct', and 'sure correct' (reversed order for half the participants). Confidence was indicated using the 1, 2, 3, 8, 9 and 0 keys at the top of the keyboard with the ring, middle and index fingers of both hands. There was no time limit for indicating confidence. During the training phase only, a trial ended with a visual presentation of feedback. An empty horizontal rectangle was filled in white starting from the left end of the rectangle (reversed order for half the participants, matched to the confidence counterbalancing). The proportion filled corresponded to the probability that the response was correct (e.g. halfway filled if feedback is 50%). Ticks at the 0, 25, 50, 75 and 100 percent marks were respectively labeled 'sure error', 'probably error', 'random chance', 'probably correct' and 'sure correct'.

On each trial, participants decided whether a box contained more elements from one out of two categories. In the letter discrimination task, elements were A's or B's, in the dot color task, blue or red dots and in the shape discrimination task, squares and circles. The total number of elements in a box was always 80, with the exact proportion of each element depending on the difficulty condition. The position of the elements was randomly generated on each trial.

Experiment 3

Participants. Thirty-four participants (12 men, age: $M = 19.5$, $SD = 3.28$, range 18–32) took part in both sessions of Experiment 3. All participants took part in return for course credit and signed informed consent at the start of the experiment. All procedures were approved by the KU Leuven ethics committee.

Procedure. Similar to the previous experiments, Experiment 3 was divided into a training phase (120 trials where participants received model-generated feedback) and a subsequent testing phase (180 trials where no feedback was provided). Crucially, participants were always

trained on the same task and then tested on the same task (Same condition) and on a different task (Transfer condition, randomized order between participants). The testing tasks were the letter discrimination task and the dot color task. Half of the participants were trained on the letter discrimination task while the other half was trained on the dot color task. Feedback was manipulated according to β in two within-participants conditions, identical to Experiment 2B's minus and plus conditions. The experiment was organized in two sessions of 1 hour, separated by 6 to 8 days. A different feedback condition was assigned to each session (randomized order across participants). There were two training and testing phases for each feedback and transfer condition, for a total of 240 training trials and 360 testing trials per feedback and transfer condition.

Model-generated feedback. Instead of binary feedback (correct/error), feedback during the training phase after each trial was provided in the form of a continuous value. Participants were told that this probability reflected the probability that their response was correct. In reality, the feedback was generated by our model of confidence. To do so, we estimated the single-trial evidence accumulation process online (i.e., during the experiment). To do so, we assumed that performance was equivalent to the average performance observed in piloting sessions. In other words, we assumed that the current decision threshold and drift rate were equal to the average decision threshold and drift rate from piloting sessions. At the moment a decision was made, the evidence accumulation process just reached the decision threshold. We thus inferred that the amount of accumulated evidence at the time of decision was equal to the average decision threshold estimated from the pilot sessions. Then, to estimate the total amount of accumulated evidence at the time of the confidence report, we added the post-decisional evidence estimated by running a random-walk for a duration fixed to the observed confidence RT and with a drift rate set to the average drift rate estimated from the pilot sessions (the sign of which varied whether the response was correct or not). Feedback was thus equal to model confidence computed according to a fixed (α , β) pairing (the value of which depended on the condition and experiment one is in) from that total evidence and the total time (decision RT + confidence RT).

Feedback conditions. In a baseline condition, the feedback presented to participants reflected the actual model-generated probability of a choice being correct. To get the value of α and β that best approximate the true probability of a choice being correct, we estimated both parameters based on the heatmap generated by the drift rates observed in the pilot sessions. In the baseline condition, α was thus set to 18 and β to 0. In Experiment 2A, for one task feedback was computed using a lower value of α (namely 9), and for another task feedback was computed using a higher value of α (namely 36; termed “ α -plus” condition). The association between the manipulation of α and the task was counterbalanced across participants. In Experiment 2B, feedback was provided according to the baseline condition in one task, using a lower value of β in another task (-1), and using a higher value of β in another task (1).

Statistical analyses. All data were analyzed using mixed effects models. We started from models including the fixed factors and their interaction(s), as well as a random intercept for each participant. These models were then extended by adding random slopes, only when this significantly improved model fit. Confidence ratings and RT were analyzed with linear mixed effects models, for which we report F statistics and the degrees of freedom as estimated by Satterthwaite's approximation. Accuracy was analyzed using a generalized linear mixed model, for which we report X^2 statistics. We computed Cramer's V as effect sizes for Chi-squared tests [68]. All model fit analyses were done using the lmerTest R package [69].

Bounded evidence accumulation. We modeled choice and RT data using the drift diffusion model (DDM), a popular variant of the wider class of accumulation-to-bound models. In the DDM, noisy evidence (representing the difference between the evidence for both options)

is accumulated, the strength of which is controlled by a drift rate ν , until one of two decision thresholds a or $-a$ is reached. Non-decision components are captured by a non-decision time parameter ter . To simulate data from the model, random walks were used as a discrete approximation of the continuous diffusion process of the drift diffusion model. Each simulated random walk process started at z^*a (here, z was an unbiased starting point fixed to 0). At each time step τ , accumulated evidence changed by Δ with Δ given in Eq (4):

$$\Delta = \nu\tau + \sigma\sqrt{\tau}N(0, 1) \quad (4)$$

Within-trial variability is given by σ . In all simulations, τ was set to 1 ms, and σ was fixed to .1.

Model fitting. Model predictions were obtained from the random-walk simulation described above. Evidence continued to accumulate after threshold crossing for a duration that was sampled from the confidence RT distribution of the trials being fitted. Note that this sampling was done without replacement, ensuring that the simulated confidence RT distribution exactly matched the empirically observed confidence RT distribution. The number of trials being simulated was equal to 20 times the number of empirical trials being fitted to ensure that every trial of the empirical confidence RT distribution is being simulated an equal amount of time. Given that the model-generated confidence comes on a continuous scale from 0 to 1, we binned the model output into 6 equally-spaced bins.

Accuracy and RT data of each task and participant was estimated using 5 DDM parameters: non-decision time, decision threshold and three drift rate parameters (one for each trial difficulty level). Additionally, α and β were fitted to the confidence judgments, separately for each feedback condition. We implemented quantile optimization, and computed the proportion of trials falling within each of six groups formed by quantiles .1, .3, .5, .7 and .9 of RT, separately for corrects and errors. Similarly with confidence ratings, we computed the proportion of trials resulting at each of the 6 levels of confidence judgment separately for corrects and errors. The resulting objective function consisted in minimizing the sum of squared errors described in Eq (5):

$$SSE = \sum_{k \in \{0,1\}} \sum_{i=1}^{N_q} (oRT_{i,k} - pRT_{i,k})^2 + \sum_{j=1}^{N_d} (oCJ_{j,k} - pCJ_{j,k})^2 \quad (5)$$

with $N_q = N_{cl} = 6$ the number of RT groups/possible confidence value, $oRT_{i,k}$ and $pRT_{i,k}$ respectively the proportions of observed and predicted trials falling within quantile i of RT, separately for corrects ($k = 1$) and errors ($k = 0$), and $oCJ_{j,k}$ and $pCJ_{j,k}$ reflecting their counterpart for confidence. Models were fitted using a differential evolution algorithm [70], as implemented in the DEoptim R package [71]. The population size was set to 10 times the number of parameters to estimate. The algorithm stopped once no improvement of the objective function was observed for the last 100 generations.

Model comparison. All candidate models for the model comparison were based on the same estimated DDM parameters fitted separately to accuracy and RT data (i.e. minimizing the first term of the SSE in Eq (5)). Each candidate model was then fit to confidence ratings (i.e. minimizing the second term of the SSE in Eq (5)). BIC values for model comparison were computed as follows:

$$BIC = k \ln(n) + n \ln\left(\frac{SSE}{n}\right) \quad (6)$$

with k the number of free parameters and n the number of data points. This formulation of BIC holds true assuming normally distributed model errors with zero mean [72]. BIC values

for each model represented in [Table 1](#) correspond to the mean BIC over participants. Bootstrapped 95% confidence intervals of confidence contrasts were obtained by simulating 500 datasets based on the fits of each participant and then computing the mean confidence contrasts of each repetition. The 95% confidence interval was computed as the .025 and .975 quantiles of the distribution formed by the bootstrapping.

Parameter recovery. To make sure that the estimated parameters from our model fits are interpretable, we performed a parameter recovery analysis that we report here. We simulated data from 200 simulated participants using parameters from the ranges observed in both Experiment 1 and 2. In order to reproduce the experimental settings, each simulated participant had 3 different trial difficulty levels (i.e. drift rates). We simulated 648 trials per participant (216 per drift rate). Recovery for all parameters was excellent (all $r_s > .93$). Additionally, there was no significant correlation between estimated α and β , $r(198) = .07$, $p = .30$, suggesting that the two parameters are not trading off one against the other.

Supporting information

S1 Text. Explicit derivation from DDM assumptions to our formulation of confidence.
(PDF)

S1 Fig. Confidence decreases with RT in both behavioral data and LDC model predictions.
(PDF)

Acknowledgments

We thank Boris Burle for useful discussions.

Author Contributions

Conceptualization: Pierre Le Denmat, Tom Verguts, Kobe Desender.

Data curation: Pierre Le Denmat.

Formal analysis: Pierre Le Denmat.

Funding acquisition: Kobe Desender.

Investigation: Pierre Le Denmat.

Methodology: Pierre Le Denmat, Tom Verguts, Kobe Desender.

Project administration: Tom Verguts, Kobe Desender.

Supervision: Tom Verguts, Kobe Desender.

Validation: Tom Verguts, Kobe Desender.

Visualization: Pierre Le Denmat.

Writing – original draft: Pierre Le Denmat.

Writing – review & editing: Pierre Le Denmat, Tom Verguts, Kobe Desender.

References

1. Sanders JI, Hangya B, Kepecs A. Signatures of a Statistical Computation in the Human Sense of Confidence. *Neuron*. 2016 May 4; 90(3):499–506. <https://doi.org/10.1016/j.neuron.2016.03.025> PMID: 27151640

2. Kiani R, Shadlen MN. Representation of Confidence Associated with a Decision by Neurons in the Parietal Cortex. *Science*. 2009 May 8; 324(5928):759–64. <https://doi.org/10.1126/science.1169405> PMID: 19423820
3. Meyniel F, Sigman M, Mainen ZF. Confidence as Bayesian Probability: From Neural Origins to Behavior. *Neuron*. 2015 Oct 7; 88(1):78–92. <https://doi.org/10.1016/j.neuron.2015.09.039> PMID: 26447574
4. Boldt A, Blundell C, De Martino B. Confidence modulates exploration and exploitation in value-based learning. *Neurosci Conscious*. 2019 Jan 1; 2019(1):niz004. <https://doi.org/10.1093/nc/niz004> PMID: 31086679
5. Drugowitsch J, Mendonça AG, Mainen ZF, Pouget A. Learning optimal decisions with confidence. *Proc Natl Acad Sci*. 2019 Dec 3; 116(49):24872–80. <https://doi.org/10.1073/pnas.1906787116> PMID: 31732671
6. Frömer R, Nassar MR, Bruckner R, Stürmer B, Sommer W, Yeung N. Response-based outcome predictions and confidence regulate feedback processing and learning. *eLife*. 2021; 10:e62825. <https://doi.org/10.7554/eLife.62825> PMID: 33929323
7. Balsdon T, Wyart V, Mamassian P. Confidence controls perceptual evidence accumulation. *Nat Commun*. 2020 Apr 9; 11(1):1753. <https://doi.org/10.1038/s41467-020-15561-w> PMID: 32273500
8. Desender K, Boldt A, Verguts T, Donner TH. Confidence predicts speed-accuracy tradeoff for subsequent decisions. *eLife*. 2019 Aug 20; 8:e43499. <https://doi.org/10.7554/eLife.43499> PMID: 31429827
9. Desender K, Boldt A, Yeung N. Subjective Confidence Predicts Information Seeking in Decision Making. *Psychol Sci*. 2018 May 1; 29(5):761–78. <https://doi.org/10.1177/0956797617744771> PMID: 29608411
10. Weiskrantz L, Warrington E, Sanders MD, Marshall J. Visual capacity in the hemianopic field following a restricted occipital ablation. *Brain J Neurol [Internet]*. 1974 [cited 2023 Mar 1]; 97(4). Available from: <https://ora.ox.ac.uk/objects/uuid:bae71800-5f6f-4e75-a4af-c36cd6a5719d> PMID: 4434190
11. Levin DT, Momen N, Drivdahl SB, Simons DJ. Change Blindness Blindness: The Metacognitive Error of Overestimating Change-detection Ability. *Vis Cogn*. 2000 Jan 1; 7(1–3):397–412.
12. Fleming SM, Ryu J, Golfinos JG, Blackmon KE. Domain-specific impairment in metacognitive accuracy following anterior prefrontal lesions. *Brain*. 2014 Oct 1; 137(10):2811–22. <https://doi.org/10.1093/brain/awu221> PMID: 25100039
13. Ratcliff R, McKoon G. The Diffusion Decision Model: Theory and Data for Two-Choice Decision Tasks. *Neural Comput*. 2008 Apr; 20(4):873–922. <https://doi.org/10.1162/neco.2008.12-06-420> PMID: 18085991
14. Kiani R, Corthell L, Shadlen MN. Choice Certainty Is Informed by Both Evidence and Decision Time. *Neuron*. 2014 Dec 17; 84(6):1329–42. <https://doi.org/10.1016/j.neuron.2014.12.015> PMID: 25521381
15. Moreno-Bote R. Decision Confidence and Uncertainty in Diffusion Models with Partially Correlated Neuronal Integrators. *Neural Comput*. 2010 Jul; 22(7):1786–811. <https://doi.org/10.1162/neco.2010.12-08-930> PMID: 20141474
16. Pleskac TJ, Busemeyer JR. Two-stage dynamic signal detection: A theory of choice, decision time, and confidence. *Psychol Rev*. 2010 Jul 26; 117(3):864. <https://doi.org/10.1037/a0019737> PMID: 20658856
17. Calder-Travis J, Charles L, Bogacz R, Yeung N. Bayesian confidence in optimal decisions. *Psychol Rev [Internet]*. 2024 [cited 2024 Jun 30]; Available from: <https://ora.ox.ac.uk/objects/uuid:f5696994-8d91-472e-a7a0-f1637c5c6573>
18. Khalvati K, Kiani R, Rao RPN. Bayesian inference with incomplete knowledge explains perceptual confidence and its deviations from accuracy. *Nat Commun*. 2021 Sep 29; 12(1):5704. <https://doi.org/10.1038/s41467-021-25419-4> PMID: 34588440
19. Van Marcke H, Denmat PL, Verguts T, Desender K. Manipulating Prior Beliefs Causally Induces Under- and Overconfidence. *Psychol Sci*. 2024 Apr 1; 35(4):358–75. <https://doi.org/10.1177/09567976241231572> PMID: 38427319
20. Sutton RS, Barto AG. *Reinforcement learning: An introduction*, 2nd ed. Cambridge, MA, US: The MIT Press; 2018. xxii, 526 p. (Reinforcement learning: An introduction, 2nd ed).
21. Bowling SR, Khasawneh MT, Kaewkuekool S, Cho BR. A logistic approximation to the cumulative normal distribution. *J Ind Eng Manag*. 2009 Jul 10; 2(1):114–27.
22. Ais J, Zylberberg A, Barttfeld P, Sigman M. Individual consistency in the accuracy and distribution of confidence judgments. *Cognition*. 2016 Jan 1; 146:377–86. <https://doi.org/10.1016/j.cognition.2015.10.006> PMID: 26513356
23. Hauser TU, Allen M, Rees G, Dolan RJ. Metacognitive impairments extend perceptual decision making weaknesses in compulsivity. *Sci Rep*. 2017 Jul 26; 7(1):6614. <https://doi.org/10.1038/s41598-017-06116-z> PMID: 28747627

24. Rollwage M, Dolan RJ, Fleming SM. Metacognitive Failure as a Feature of Those Holding Radical Beliefs. *Curr Biol*. 2018 Dec 17; 28(24):4014–4021.e8. <https://doi.org/10.1016/j.cub.2018.10.053> PMID: 30562522
25. Rouault M, Seow T, Gillan CM, Fleming SM. Psychiatric Symptom Dimensions Are Associated With Dissociable Shifts in Metacognition but Not Task Performance. *Biol Psychiatry*. 2018 Sep 15; 84(6):443–51. <https://doi.org/10.1016/j.biopsych.2017.12.017> PMID: 29458997
26. Kepecs A, Mainen ZF. A computational framework for the study of confidence in humans and animals. *Philos Trans R Soc B Biol Sci*. 2012 May 19; 367(1594):1322–37.
27. Rausch M, Zehetleitner M. The folded X-pattern is not necessarily a statistical signature of decision confidence. *PLoS Comput Biol*. 2019 Oct 21; 15(10):e1007456. <https://doi.org/10.1371/journal.pcbi.1007456> PMID: 31634359
28. Herregods S, Denmat PL, Desender K. Modelling Speed-Accuracy Tradeoffs in the Stopping Rule for Confidence Judgments [Internet]. *bioRxiv*; 2023 [cited 2023 Jun 13]. p. 2023.02.27.530208. Available from: <https://www.biorxiv.org/content/10.1101/2023.02.27.530208v1>
29. Moran R, Teodorescu AR, Usher M. Post choice information integration as a causal determinant of confidence: Novel data and a computational account. *Cognit Psychol*. 2015 May 1; 78:99–147. <https://doi.org/10.1016/j.cogpsych.2015.01.002> PMID: 25868113
30. Allen M, Frank D, Schwarzkopf DS, Fardo F, Winston JS, Hauser TU, et al. Unexpected arousal modulates the influence of sensory noise on confidence. Shan H, editor. *eLife*. 2016 Oct 25; 5:e18103. <https://doi.org/10.7554/eLife.18103> PMID: 27776633
31. Boldt A, de Gardelle V, Yeung N. The impact of evidence reliability on sensitivity and bias in decision confidence. *J Exp Psychol Hum Percept Perform*. 2017 Apr 6; 43(8):1520. <https://doi.org/10.1037/xhp0000404> PMID: 28383959
32. Maniscalco B, Lau H. Manipulation of working memory contents selectively impairs metacognitive sensitivity in a concurrent visual discrimination task. *Neurosci Conscious*. 2015 Jan 1; 2015(1):niv002.
33. Palminteri S, Wyart V, Koehlin E. The Importance of Falsification in Computational Cognitive Modeling. *Trends Cogn Sci*. 2017 Jun 1; 21(6):425–33. <https://doi.org/10.1016/j.tics.2017.03.011> PMID: 28476348
34. Burnham KP, Anderson DR. Multimodel Inference: Understanding AIC and BIC in Model Selection. *Sociol Methods Res*. 2004 Nov 1; 33(2):261–304.
35. Rigoux L, Stephan KE, Friston KJ, Daunizeau J. Bayesian model selection for group studies—Revisited. *NeuroImage*. 2014 Jan 1; 84:971–85. <https://doi.org/10.1016/j.neuroimage.2013.08.065> PMID: 24018303
36. Stephan KE, Penny WD, Daunizeau J, Moran RJ, Friston KJ. Bayesian model selection for group studies. *NeuroImage*. 2009 Jul 15; 46(4):1004–17. <https://doi.org/10.1016/j.neuroimage.2009.03.025> PMID: 19306932
37. Adler WT, Ma WJ. Comparing Bayesian and non-Bayesian accounts of human confidence reports. *PLoS Comput Biol*. 2018 Nov 13; 14(11):e1006572. <https://doi.org/10.1371/journal.pcbi.1006572> PMID: 30422974
38. Constant M, Pereira M, Faivre N, Filevich E. Prior information differentially affects discrimination decisions and subjective confidence reports. *Nat Commun*. 2023 Sep 6; 14(1):5473. <https://doi.org/10.1038/s41467-023-41112-0> PMID: 37673881
39. Geurts LS, Cooke JRH, van Bergen RS, Jehee JFM. Subjective confidence reflects representation of Bayesian probability in cortex. *Nat Hum Behav*. 2022 Feb; 6(2):294–305. <https://doi.org/10.1038/s41562-021-01247-w> PMID: 35058641
40. Ackerman R, Thompson VA. Meta-Reasoning: Monitoring and Control of Thinking and Reasoning. *Trends Cogn Sci*. 2017 Aug 1; 21(8):607–17. <https://doi.org/10.1016/j.tics.2017.05.004> PMID: 28625355
41. Fleming SM, Daw ND. Self-evaluation of decision-making: A general Bayesian framework for metacognitive computation. *Psychol Rev*. 2017; 124:91–114. <https://doi.org/10.1037/rev0000045> PMID: 28004960
42. Jang Y, Wallsten TS, Huber DE. A stochastic detection and retrieval model for the study of metacognition. *Psychol Rev*. 2012; 119:186–200. <https://doi.org/10.1037/a0025960> PMID: 22059901
43. Maniscalco B, Lau H. The signal processing architecture underlying subjective reports of sensory awareness. *Neurosci Conscious*. 2016 Jan 1; 2016(1):niw002. <https://doi.org/10.1093/nc/niw002> PMID: 27499929
44. Rausch M, Hellmann S, Zehetleitner M. Confidence in masked orientation judgments is informed by both evidence and visibility. *Atten Percept Psychophys*. 2018 Jan 1; 80(1):134–54. <https://doi.org/10.3758/s13414-017-1431-5> PMID: 29043657

45. Shekhar M, Rahnev D. The nature of metacognitive inefficiency in perceptual decision making. *Psychol Rev.* 2021; 128:45–70. <https://doi.org/10.1037/rev0000249> PMID: 32673034
46. Zylberberg A, Barttfeld P, Sigman M. The construction of confidence in a perceptual decision. *Front Integr Neurosci* [Internet]. 2012 [cited 2023 Feb 17];6. Available from: <https://www.frontiersin.org/articles/10.3389/fnint.2012.00079>
47. Shekhar M, Rahnev D. How do humans give confidence? A comprehensive comparison of process models of perceptual metacognition. *J Exp Psychol Gen.* 2024; 153(3):656–88. <https://doi.org/10.1037/xge0001524> PMID: 38095983
48. Zylberberg A, Fetsch CR, Shadlen MN. The influence of evidence volatility on choice, reaction time and confidence in a perceptual decision. *eLife.* 2016; 5:e17688. <https://doi.org/10.7554/eLife.17688> PMID: 27787198
49. Massoni S. Emotion as a boost to metacognition: How worry enhances the quality of confidence. *Conscious Cogn.* 2014 Oct 1; 29:189–98. <https://doi.org/10.1016/j.concog.2014.08.006> PMID: 25286128
50. Overhoff H, Ko YH, Feuerriegel D, Fink GR, Stahl J, Weiss PH, et al. Neural correlates of metacognition across the adult lifespan. *Neurobiol Aging.* 2021 Dec 1; 108:34–46. <https://doi.org/10.1016/j.neurobiolaging.2021.08.001> PMID: 34487950
51. Weil LG, Fleming SM, Dumontheil I, Kilford EJ, Weil RS, Rees G, et al. The development of metacognitive ability in adolescence. *Conscious Cogn.* 2013 Mar 1; 22(1):264–71. <https://doi.org/10.1016/j.concog.2013.01.004> PMID: 23376348
52. Ko Y, Lau H. A detection theoretic explanation of blindsight suggests a link between conscious perception and metacognition. *Philos Trans R Soc B Biol Sci.* 2012 May 19; 367(1594):1401–11. <https://doi.org/10.1098/rstb.2011.0380> PMID: 22492756
53. Balsdon T, Pisauro MA, Philiastides MG. Distinct basal ganglia contributions to learning from implicit and explicit value signals in perceptual decision-making. *Nat Commun.* 2024 Jun 22; 15(1):5317. <https://doi.org/10.1038/s41467-024-49538-w> PMID: 38909014
54. Johnson D, Fowler J. The Evolution of Overconfidence. *Nature.* 2011 Sep 15; 477:317–20. <https://doi.org/10.1038/nature10384> PMID: 21921915
55. Li K, Szolnoki A, Cong R, Wang L. The coevolution of overconfidence and bluffing in the resource competition game. *Sci Rep.* 2016 Feb 17; 6(1):21104. <https://doi.org/10.1038/srep21104> PMID: 26883799
56. Bang D, Ershadmanesh S, Nili H, Fleming SM. Private–public mappings in human prefrontal cortex. Frank MJ, Kahnt T, Chang SW, editors. *eLife.* 2020 Jul 23; 9:e56477. <https://doi.org/10.7554/eLife.56477> PMID: 32701449
57. Lewthwaite R, Wulf G. Social-comparative feedback affects motor skill learning. *Q J Exp Psychol.* 2010; 63:738–49. <https://doi.org/10.1080/17470210903111839> PMID: 19691002
58. Herzog MH, Fahle M. The role of feedback in learning a vernier discrimination task. *Vision Res.* 1997 Aug 1; 37(15):2133–41. [https://doi.org/10.1016/s0042-6989\(97\)00043-6](https://doi.org/10.1016/s0042-6989(97)00043-6) PMID: 9327060
59. Shiu LP, Pashler H. Improvement in line orientation discrimination is retinally local but dependent on cognitive set. *Percept Psychophys.* 1992 Sep 1; 52(5):582–8. <https://doi.org/10.3758/bf03206720> PMID: 1437491
60. Ernst B, Steinhauser M. Effects of invalid feedback on learning and feedback-related brain activity in decision-making. *Brain Cogn.* 2015 Oct 1; 99:78–86. <https://doi.org/10.1016/j.bandc.2015.07.006> PMID: 26263382
61. Becker MPI, Nitsch AM, Schlösser R, Koch K, Schachtzabel C, Wagner G, et al. Altered emotional and BOLD responses to negative, positive and ambiguous performance feedback in OCD. *Soc Cogn Affect Neurosci.* 2014 Aug 1; 9(8):1127–33. <https://doi.org/10.1093/scan/nst095> PMID: 23893850
62. Gu R, Ge Y, Jiang Y, Jia Luo Y. Anxiety and outcome evaluation: The good, the bad and the ambiguous. *Biol Psychol.* 2010 Oct 1; 85(2):200–6. <https://doi.org/10.1016/j.biopsycho.2010.07.001> PMID: 20619312
63. Van den Brink RL, Murphy PR, Desender K, Ru N de, Nieuwenhuis S. Temporal Expectation Hastens Decision Onset But Does Not Affect Evidence Quality. *J Neurosci.* 2021 Jan 6; 41(1):130–43.
64. Rafiei F, Rahnev D. Qualitative speed-accuracy tradeoff effects that cannot be explained by the diffusion model under the selective influence assumption. *Sci Rep.* 2021 Jan 8; 11(1):45. <https://doi.org/10.1038/s41598-020-79765-2> PMID: 33420181
65. Boundy-Singer ZM, Ziemba CM, Goris RLT. Confidence reflects a noisy decision reliability estimate. *Nat Hum Behav.* 2022 Nov 7; 1–13.
66. Guggenmos M. Reverse engineering of metacognition. Wyart V, Frank MJ, Fleming S, editors. *eLife.* 2022 Sep 15; 11:e75420. <https://doi.org/10.7554/eLife.75420> PMID: 36107147

67. Peirce J, Gray JR, Simpson S, MacAskill M, Höchenberger R, Sogo H, et al. PsychoPy2: Experiments in behavior made easy. *Behav Res Methods*. 2019 Feb 1; 51(1):195–203. <https://doi.org/10.3758/s13428-018-01193-y> PMID: 30734206
68. Kim HY. Statistical notes for clinical researchers: Chi-squared test and Fisher's exact test. *Restor Dent Endod*. 2017 May; 42(2):152–5. <https://doi.org/10.5395/rde.2017.42.2.152> PMID: 28503482
69. Kuznetsova A, Brockhoff PB, Christensen RHB. lmerTest Package: Tests in Linear Mixed Effects Models. *J Stat Softw*. 2017 Dec 6; 82:1–26.
70. Price K, Storn RM, Lampinen JA. *Differential Evolution: A Practical Approach to Global Optimization*. Springer Science & Business Media; 2006. 544 p.
71. Mullen K, Ardia D, Gil DL, Windover D, Cline J. DEoptim: An R Package for Global Optimization by Differential Evolution [Internet]. Rochester, NY; 2009 [cited 2022 Oct 14]. Available from: <https://papers.ssrn.com/abstract=1526466>
72. Solway A, Botvinick MM. Evidence integration in model-based tree search. *Proc Natl Acad Sci*. 2015 Sep 15; 112(37):11708–13. <https://doi.org/10.1073/pnas.1505483112> PMID: 26324932