



**HAL**  
open science

# Does nonlinear neural network dynamics explain human confidence in a sequence of perceptual decisions ?

Kevin Berlemont, Jean-Remy Martin, Jérôme Sackur, Jean-Pierre Nadal

## ► To cite this version:

Kevin Berlemont, Jean-Remy Martin, Jérôme Sackur, Jean-Pierre Nadal. Does nonlinear neural network dynamics explain human confidence in a sequence of perceptual decisions?. 2019. hal-02138028v2

**HAL Id: hal-02138028**

**<https://hal.science/hal-02138028v2>**

Preprint submitted on 24 Jun 2019 (v2), last revised 5 Mar 2020 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# DOES NONLINEAR NEURAL NETWORK DYNAMICS EXPLAIN HUMAN CONFIDENCE IN A SEQUENCE OF PERCEPTUAL DECISIONS ?

Kevin Berlemont<sup>1,2\*</sup>, Jean-Rémy Martin<sup>2</sup>, Jérôme Sackur<sup>3</sup>, Jean-Pierre Nadal<sup>1,4</sup>

<sup>1</sup> Laboratoire de Physique de École Normale Supérieure, PSL University, CNRS, Sorbonne Université, Université Paris-Diderot, Sorbonne Paris Cité, 75005 Paris, France.

<sup>2</sup> Centre for Research in Cognition & Neurosciences, Faculté des Sciences Psychologiques et de l'Éducation, Université Libre de Bruxelles (ULB), B-1050 Bruxelles, Belgium.

<sup>3</sup> Laboratoire de Sciences Cognitives et Psycholinguistique, École des Hautes Études en Sciences Sociales (EHESS), PSL University, Département d'études cognitives, (CNRS/ENS/EHESS), 75005 Paris, France

<sup>4</sup> Centre d'Analyse et de Mathématiques Sociales, École des Hautes Études en Sciences Sociales, PSL University, CNRS, 75006 Paris, France.

---

## ABSTRACT

Electrophysiological recordings during perceptual decision tasks in monkeys suggest that the degree of confidence in a decision is based on a simple neural signal produced by the neural decision process. Attractor neural networks provide an appropriate biophysical modeling framework, and account for the experimental results very well. However, it remains unclear whether attractor neural networks can account for confidence reports in humans. We present the results from an experiment in which participants are asked to perform an orientation discrimination task, followed by a confidence judgment. Here we show that an attractor neural network model quantitatively reproduces, for each participant, the relations between accuracy, response times and confidence, as well as sequential effects. Our results suggest that a metacognitive process such as confidence in one's decision is linked to the intrinsically nonlinear dynamics of the decision-making neural network.

---

## INTRODUCTION

A general understanding of the notion of confidence is that it quantifies the degree of belief in a decision [Meyniel et al., 2015b, Mamassian, 2015]. Many cognitive and psychology studies have tackled the problem of confidence estimation either by directly requiring participants to provide an estimation of their confidence [Peirce and Jastrow, 1884, Zylberberg et al., 2012, Adler and Ma, 2018], or by using postdecision wagering (subjects can choose a safe option, with low reward regardless of the correct choice) [Vickers, 1979 (reedited in 2014, Seth, 2008, Fleming et al., 2010, Kepecs and Mainen, 2012, Massoni, 2014)]. Postdecision wagering has been used in behaving animals in order to study the neural basis of confidence [Smith et al., 2003, Kepecs et al., 2008, Kiani and Shadlen, 2009, Komura et al., 2013, Lak et al., 2014].

In the context of two alternative forced choices (2AFC), researchers have studied the process of confidence readout with various models: using Bayesian decision theory and signal detection theory [Clarke et al., 1959, Galvin et al., 2003, Fleming et al., 2010, Kepecs and Mainen, 2012], as a measure of precision [Yeung and Summerfield, 2012, Meyniel et al., 2015a], or using integration of evidence over time [Smith and Vickers, 1988, Kepecs et al., 2008, Pleskac and Busemeyer, 2010, Drugowitsch et al., 2012]. Authors [Kepecs and Mainen, 2012, Meyniel et al., 2015b] have suggested that choice and confidence can be read out from the same neural representation. In the experiment by Kiani and Shadlen [2009], monkeys perform a two alternative forced choices task, but after a certain delay a sure target is presented for a certain but small reward. The probability of choosing this sure target reflects the monkey's degree of choice uncertainty, assuming that risk aversion strongly correlates with this uncertainty. In order to account for the experimental findings in this uncertain option task, Wei and Wang [2015] and Jaramillo et al. [2019] have proposed biophysical attractor neural network models. They show that these models capture both behavioral performance and the associated physiological recordings from Lateral Intraparietal (LIP) neurons.

---

\*Corresponding contact: kevin.berlemont@lps.ens.fr

In the present work, we address the issue of the ability of attractor networks to quantitatively account for confidence reports in human. For this, we first experimentally investigate confidence formation and its impact on sequential effects in human experiments. Participants perform an orientation discrimination task on Gabor patches that deviate clockwise or counter-clockwise with respect to the vertical. In some blocks, after reporting their decisions, participants perform a confidence judgment on a visual scale. Then, we fit an attractor neural network model [Wong and Wang, 2006, Berlemont and Nadal, 2019] on the behavioral data. More precisely, for each participant, we calibrate a network specifically on his/her behavioral data, the fit being only based on mean response times and accuracy. With the model so calibrated for each participant, and making simulations that replicate the experimental protocol, here, for the first time, we quantitatively confront an attractor neural network behavior with human behavior during full sequences of perceptual decisions. Following Wei and Wang [2015], we assume that confidence is an increasing function of the difference, measured at the time the decision is made, between the mean spike rates of the two neural pools specific to one or the other of the two possible choices. We show that in this way, behavioral effects of confidence can be accurately estimated for each participant. We find that the attractor neural network accurately reproduces an effect of confidence on serial dependence which is observed in the experiment: participants are faster (respectively slower) on trials following high (resp. low) confidence trials. We show that this effect is intrinsic to the non-linear nature of the network dynamics.

## RESULTS

### EXPERIMENT AND NEURAL MODEL

Participants completed a visual discrimination task between clockwise and counter-clockwise orientated stimuli, followed, or not, by a task in which they were asked to assess the confidence in their decision. We used three kinds of blocks, comprising either sequences of pure decision trials (*pure* blocks), trials with feedback (*feedback* blocks) or trials with confidence judgments (*confidence* blocks). In *feedback* blocks, on each trial, participants received auditory feedback on the correctness of their choice. In the *confidence* block, after each trial, they were asked to report their confidence on a discrete scale of ten levels, from 0 to 9. In *feedback* blocks, participants were not asked to report their confidence, and in *confidence* blocks they did not receive any feedback. We illustrate the experimental protocol in Figure 1, panels A - D.

For the modelling of the neural correlates, we consider a decision-making recurrent neural network governed by local excitation and feedback inhibition, based on the biophysical models of spiking neurons introduced and studied in Compte et al. [2000] and Wang [2002]. We work with the reduced version derived in Wong and Wang [2006], allowing for large-scale numerical simulations and for better analytic analysis. More precisely, we consider the model variant introduced in Berlemont and Nadal [2019], which takes into account a corollary discharge (see Figure 1.E) allowing the network to engage in a sequence of perceptual decisions. The model consists of two competing units, each one representing an excitatory neuronal pool, selective to one of the two available response options, here  $C$  (clockwise) or  $AC$  (anti-clockwise). Each population receives a task-related input signaling the perceived evidence for each option. The difference between these inputs varies inversely with the difficulty of the task, thus it varies with the absolute value of the Gabor orientation.

The decision, ' $C$ ' or ' $AC$ ', is made when one of the two units reaches a threshold  $z$ . Once a decision is made (threshold is reached), an inhibitory current (the corollary discharge) is injected into the two neural pools, causing a relaxation of the neural activities towards a low activity, neutral state, therefore allowing the network to deal with consecutive sequences of trials, as illustrated in Figure 1.F. For a biologically relevant range of parameters, relaxation is not complete at the onset of the next stimulus, hence the decision made in this new trial will depend on the one at the previous trial. In a previous work, we showed [Berlemont and Nadal, 2019] that the model accounts for the main sequential and post-error effects observed in perceptual decision making experiments in human and monkeys.

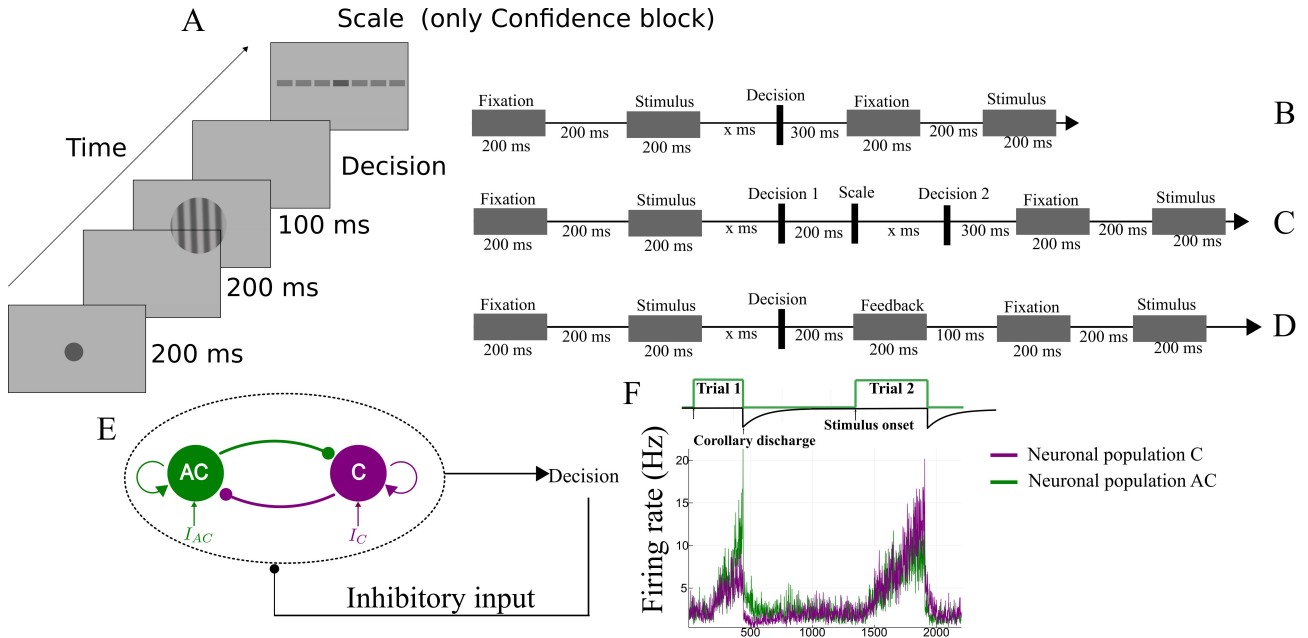
Full details about the experiment and the model can be found in the Material and Methods Section.

### CALIBRATION OF THE MODEL ONTO THE BEHAVIORAL RESULTS

In Figure 2, we show for each participant response times and accuracies with respect to stimulus orientation (absolute value of the orientation angle). All subjects exhibit improved accuracies and shorter response times for less difficult (larger orientation) stimuli, as classically reported in the literature [Baranski and Petrusic, 1994].

We fit the model to these behavioral data. As detailed in the Material and Methods Section, we perform model calibration in order to reproduce both the mean response times and the accuracy (success rates). For each participant, this is done separately for the three types of blocks.

First, we note that the model correctly reproduces the behavioral results of the different participants, as can be seen in



**Figure 1: Experimental protocol and model architecture.** Procedure of the discrimination task, for the three blocks. (A) Structure of a trial: Following a fixation period, the circular grating (Gabor patch, oriented clockwise, C, or counter-clockwise, AC) appears and participants make the decision (C or AC). In confidence blocks, after a delay, participants report their confidence with respect to their choice, on a discrete scale with 10 levels. (B) Time course of a pure block trial. (C) Time course of a confidence block trial. (D) Time course of a feedback block trial. (E) Decision-making network structure. The network consists of two neural pools (specific to clockwise (C) and anti-clockwise (AC) stimuli), endowed with self-excitation and mutual inhibition. After a decision is made (threshold crossed), a non specific inhibitory input (corollary discharge) is sent onto both units. (F) Time course of the neural activities of both pools during two consecutive trials.

Figure 2. Second, we compare the values of the parameters obtained for the pure and confidence blocks. We find that participants have higher decision threshold (Signed Rank test [Wilcoxon, 1945]  $p = 0.03$ ), higher stimulus strength level by angle (Signed Rank test,  $p = 0.031$ ) and higher mean non-decision times (Signed Rank test  $p = 0.03$ ). Two of the authors of the present paper (J-R. Martin, J. Sackur, personal communication, April, 2018) have obtained analogous results when analyzing similar data within the Drift Diffusion framework: non-decision time, drift rate and decision threshold are modified by the confidence context in the experimental setup.

**Non-decision times.** Our fitting procedure allows estimating the non-decision times. In Figure 3, we represent the histogram of the response times across participants for the pure and confidence blocks. The red curve shows the distribution of non-decision times in the model, and the black curve the response times distribution. We note that, with a fit only based on the *mean* response times and accuracies, the model also accurately account for the *distributions* of response times. We find that the minimum value of non-decision time is 75 ms for the pure block, and 100 ms for the confidence block, and the average non-decision times are within the order of magnitude of saccadic latency [Luce et al., 1986, Mazurek et al., 2003]. Finally, we observe that the non-decision times distributions clearly show a right skew for several participants, in agreement with Verdonck and Tuerlinckx [2016]. This justifies the modelling of non-decision times with an exponentially modified Gaussian distribution (EMG) [Grushka, 1972], instead of simply adding a constant non-decision time to every decision time.

## CONFIDENCE MODELING

Recent studies have reported a choice-independent representation of confidence signal in monkeys [Ding and Gold, 2011] and in rats [Kepecs et al., 2008], as well as evidence for a close link between decision variable and confidence – in monkeys from LIP recordings [Kiani and Shadlen, 2009] and in humans from fMRI experiments [Hebart et al., 2014]. In an experiment with monkeys, Kiani and Shadlen [2009] introduce a ‘sure target’ associated with a low reward, which can be chosen instead of the categorical targets. The probability of *not* choosing the sure target is then a proxy for the

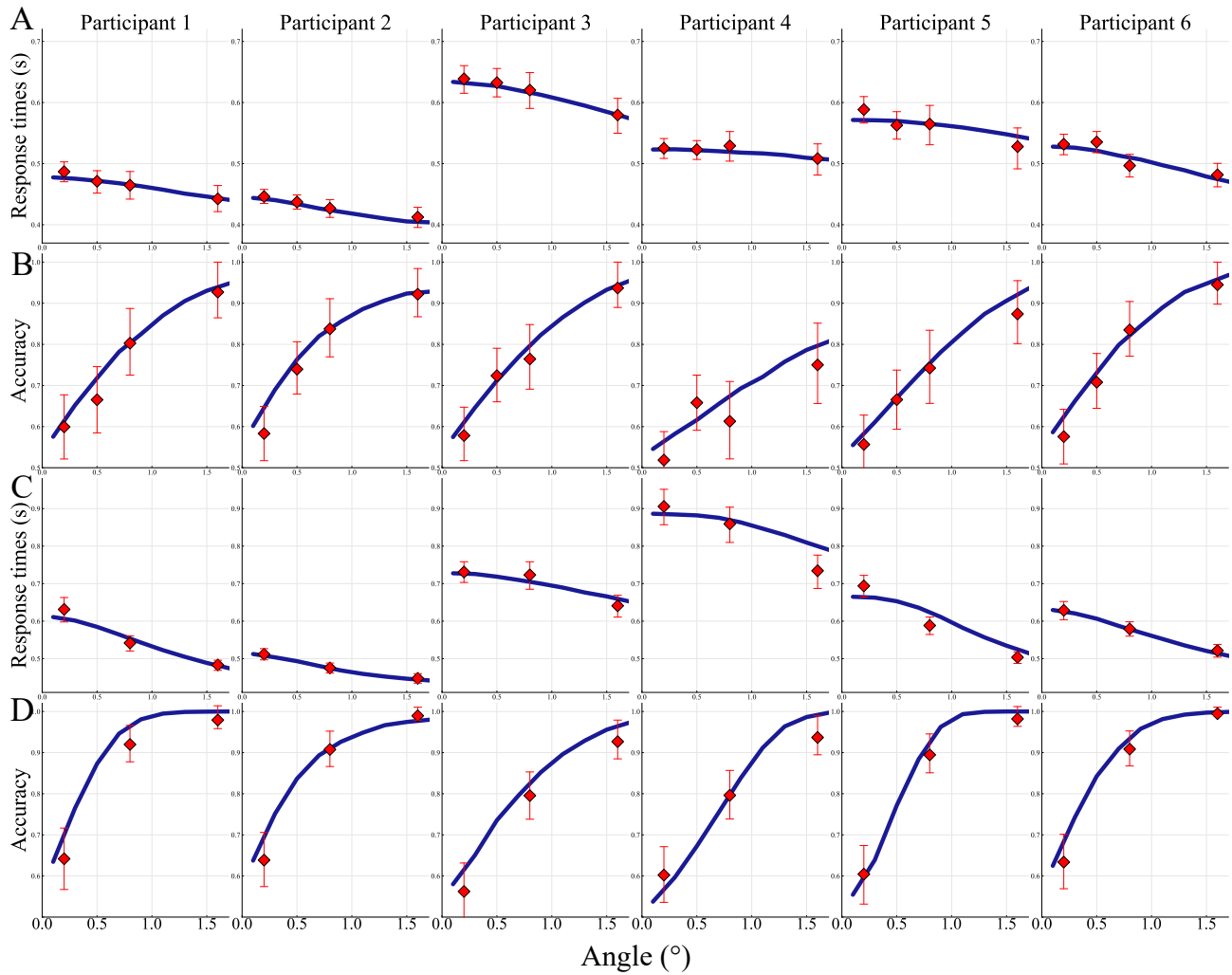


Figure 2: Mean response times (A,C) and accuracies (B,D) as a function of the absolute value of stimulus orientation, in the pure (A and B) and confidence (C and D) blocks. For each subject we represent the behavioral data (red dots) and the associated fitted model (blue line). Error bars are 95% confidence interval using the bootstrap method.

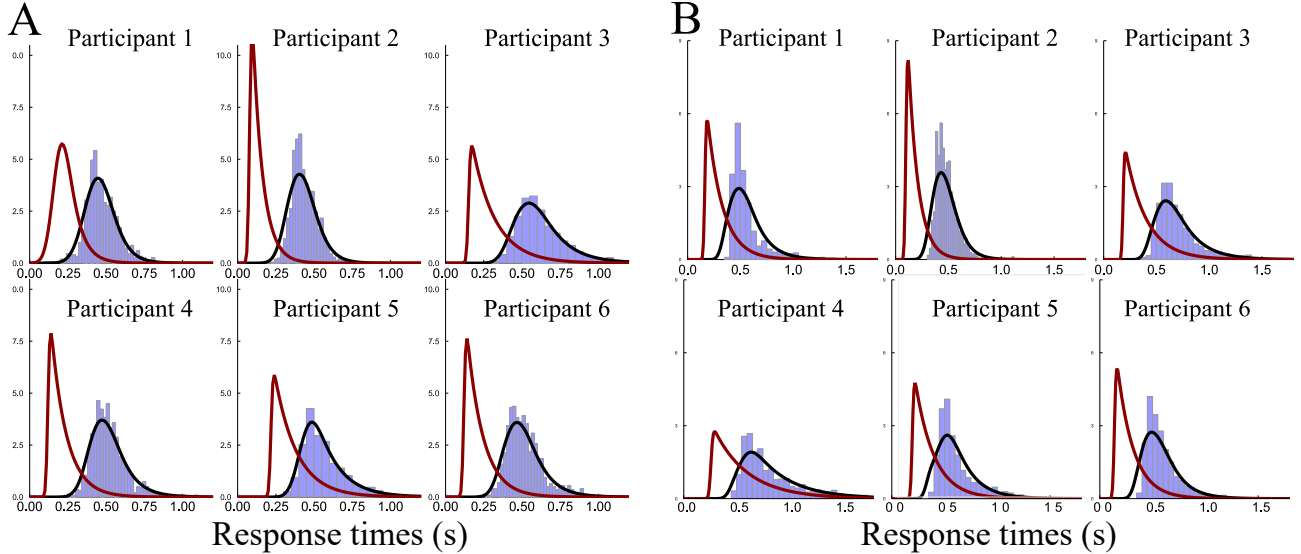


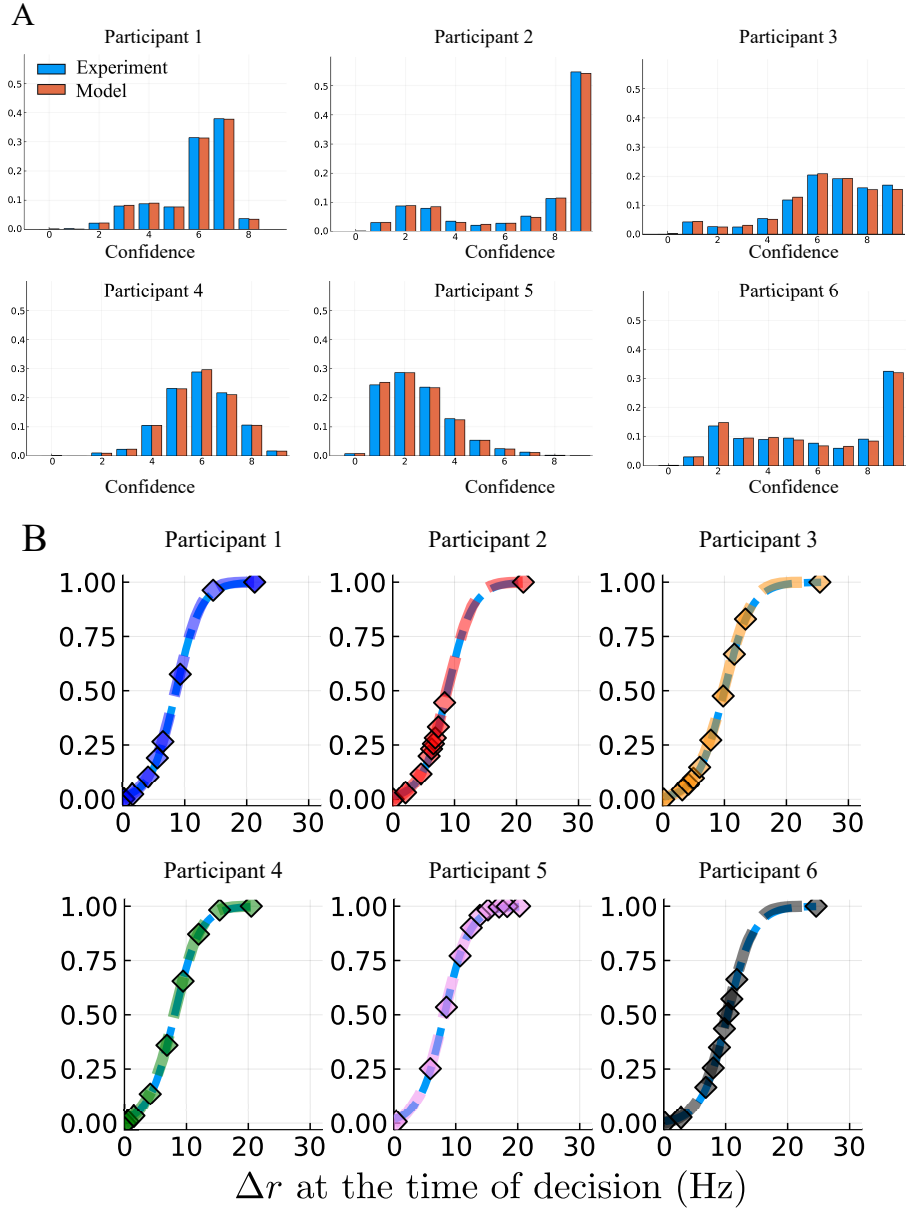
Figure 3: **Distributions of RTs for each subject (A) Pure block data, (B) confidence block data.** For both panels: In blue, participants’ histograms of the response times; Black curve: density of response times of the simulated network model; Red curve: the associated non decision response times distribution.

confidence level. [Wei and Wang \[2015\]](#) model the neural correlates of confidence within the framework of attractor neural networks. They assume that the confidence level (as given by the probability of not choosing the sure target) is a sigmoidal function of the difference, at the time of decision, between the activities of the winning and losing neural pools. This hypothesis is in line with similar hypothesis in the framework of DDMs and other decision-making models [[Vickers, 1979](#) (reedited in 2014, [Mamassian, 2015](#), [Drugowitsch et al., 2014](#))]. They then show that the empirical dependencies of response times and accuracies in the confidence level are qualitatively reproduced in the simulations of the neural model.

Following [Wei and Wang \[2015\]](#), we make here the hypothesis that the confidence in a decision is based on the difference  $\Delta r$  between the neural activities of the winning and losing neural pools, measured at the time of the decision: the larger the difference, the greater the confidence. In our experiment, the measure of confidence is the one reported by the subjects on a discrete scale, and it is this reported confidence level that we want to model. Within our framework, we quantitatively link this empirical confidence to the neural difference  $\Delta r$  by matching the distribution of the neural evidence balance with the empirical histogram of the confidence levels. In Figure 4, we show, for each participant, the matching between the histogram of confidence levels, as reported by the participant, and the distribution of  $\Delta r$ , as obtained in the model calibrated on the participant performance. We note that the main difference between the participants’ histograms lies in the percent of trials (and level on the confidence scale) for which a participant reports the highest confidence level. This last point is highly dependent on the participant, and can be at a very low value of  $\Delta r$  (see e.g. Participant 2 on Figure 4.B).

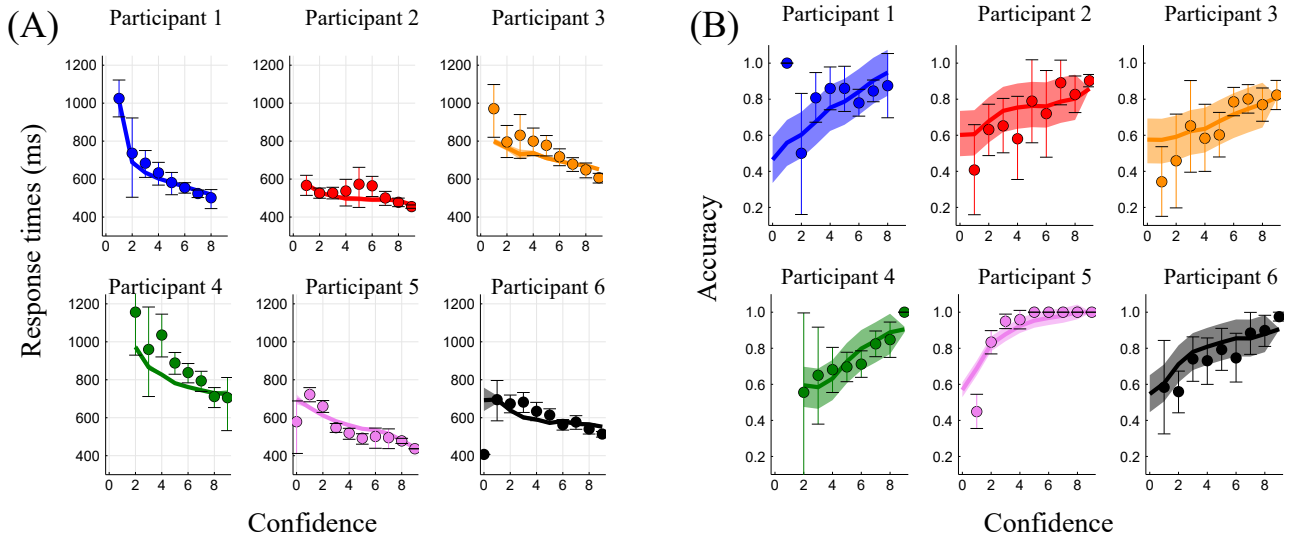
In our analysis, the shape of the mapping is not chosen a priori but non-parametrically inferred from the experimental data. This is in contrast with previous studies in which the sigmoidal shape is imposed [[Beck et al., 2008](#), [Kepecs et al., 2008](#), [Kepecs and Mainen, 2012](#), [Wei and Wang, 2015](#)]. We find however, that for each participant, the empirical mapping is very well approximated by a sigmoidal function of the type  $1/(1 + \exp(-\beta(\Delta r - \kappa)))$ , with participant-specific parameters  $\kappa$  and  $\beta$ . In [Wei and Wang \[2015\]](#), the authors exhibit a link between a probabilistic measure of confidence and  $\Delta r$ , under the form of a sigmoid function. The similarity of our findings thus suggests that the human reported confidence can be understood as a discretization of a probabilistic function.

**Response times and accuracies vs. Confidence** Studies have shown that confidence ratings are closely linked to response times [[Baranski and Petrusic, 1994](#), [Desender et al., 2018a](#)] and choice accuracy [[Peirce and Jastrow, 1884](#), [Baranski and Petrusic, 1994](#), [Sanders et al., 2016](#), [Urai et al., 2017](#)]. The behavioral confidence in our model is assumed to be based on a simple neural quantity measured at the time of the decision. In what follows, we study whether this hypothesis on the neural correlates of confidence can account for the links between the behavioral data: response times, accuracy and confidence. In Figure 5, we represent the response times (Figure 5.A) and choice accuracy (Figure 5.B) with respect to the reported confidence level for each participant. The data points show the experimental results (with the error bars as the bootstrapped 95% confidence interval), and the colored line the result of the simulation (with the light colored area the bootstrapped 95% confidence interval). Response times decrease [[Baranski and Petrusic, 1994](#), [Desender et al.,](#)



**Figure 4: Matching network confidence measure to empirical behavioral confidence.** (A) Confidence histograms. The x-axis gives the value of the confidence on a discrete scale from 0 to 9. Each sub-panel corresponds to a different participant with, in blue, the histogram of the reported confidence, and in orange, the one from the model. For clarity we plot the blue and orange bars side by side, but the bins of the histograms are, by construction, identical. (B) Transfer function  $F$  for each participant. The x-axis denotes the difference in neural pools activities  $\Delta r$  at the time of the decision, and the y-axis the cumulative distribution of  $\Delta r$ . Each point represents the levels of  $\Delta r$  delimiting the level of confidence (from left to right, confidence level 0 to confidence level 9). The dashed colored curve is the cumulative distribution function (CDF) and the light blue dashed curve is the fit of the CDF by a sigmoid.





**Figure 5: Response times and Accuracy as a function of confidence.** (A) Response times, (B) Accuracy. For both panels: each sub-panel represents a different participant. Dots are experimental data with 95% bootstrapped confidence interval as error bars. Lines are averages over 20 simulations of the attractor neural network model (calibrated as explained in the Material and Methods Section). The shaded area represents the 95% bootstrapped confidence interval on the mean.

2018a] and accuracies increase with confidence [Geller and Whitman, 1973, Vickers and Packer, 1982, Sanders et al., 2016, Desender et al., 2018a]. We find a monotonic dependency between response times and confidence, and between accuracy and confidence, but with specific shapes for each participant. Note that some values of confidence are only observed for a few trials, resulting then in large error bars especially for accuracy as we take the mean of a binary variable. For the numerical simulations, the relatively large size of the confidence interval is due to the limited number of trials, since we limit ourselves to the same protocol as the experimental one.

For comparison, we also fit on the same data another non-linear model with mutual inhibition, the Usher-McClelland model [Usher and McClelland, 2001] (more details in the Material and Methods section). We fit the Usher-McClelland model separately for each participant using the same optimization algorithm as done for the attractor network model. In Figure 6, we represent the response times and accuracy with respect to confidence. This figure has to be contrasted with Figure 5. We note that the model fits the response times with respect to confidence, but only at intermediate levels of confidence. For some participants, we observe a strong divergence at high confidence (Participant 1, 4 and 5). This can be understood by the fact that, in this model, 'firing rate' variables can take negative values (in fact, for the steady state in the absence of any input, the firing rate variables take negative values). This leads to extreme value of confidence for long trials. However, the trend in accuracy is not correct for some participants (Participants 1 et 4). Accuracy is an increasing function of confidence (except for participant 5), but the experimental data do not fall within the bootstrapped confidence interval of the simulations.

In contrast, we see that our more biophysical model correctly reproduces the psychometric and chronometric functions with respect to confidence for each participant, despite the important difference of response times between participants.

Previous studies found that, during a perceptual task, reported confidence increases with stimulus strength for correct trials, but decreases for error trials [Kepecs et al., 2008, Sanders et al., 2016, Desender et al., 2018a]. This effect of confidence has been correlated to patterns of firing rates in experiments with rats [Kepecs et al., 2008] and to the human feeling of confidence [Sanders et al., 2016]. This effect is in accordance with a prediction of statistical confidence, defined as the Bayesian posterior probability that the decision-maker is correct [Griffin and Tversky, 1992, Ernst and Banks, 2002, Sanders et al., 2016]. In Figure 7, we represent the mean confidence as a function of stimulus strength, for correct and error trials. We observe the same type of variations of confidence with respect to stimulus strength, both in the experimental results and in the model simulations. We thus see that the attractor network model reproduces a key feature of statistical confidence.



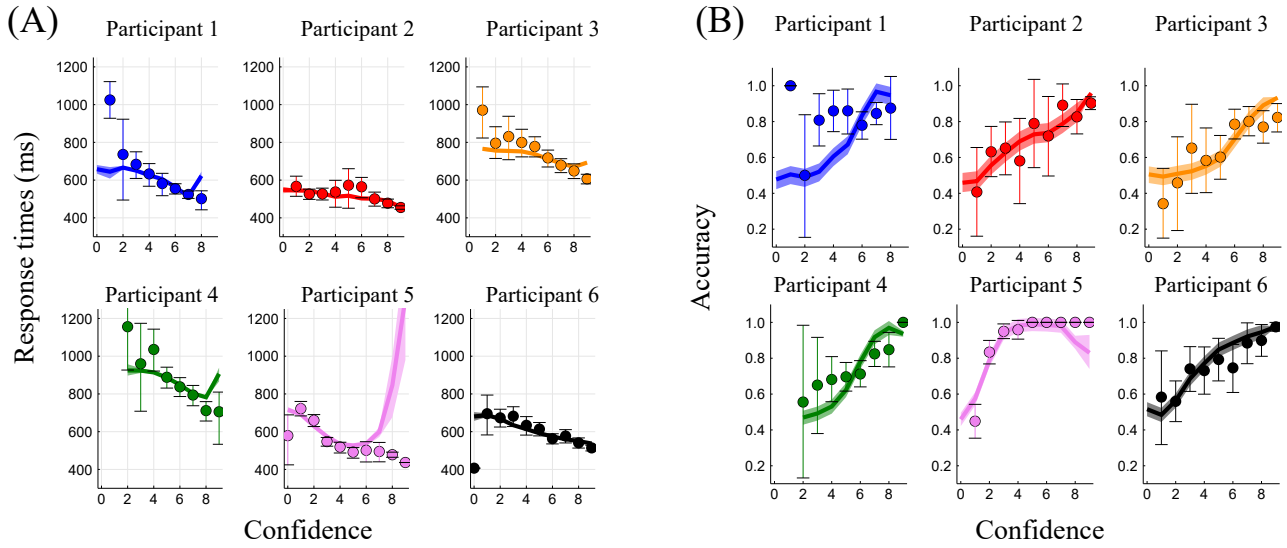


Figure 6: **Fit of the Usher-McClelland model.** (A) Response times, and (B), Accuracy, with respect to confidence. For both panels: Each sub-panel represents a different participant. The x-axis gives the reported confidence. The dots stand for the experimental data, and the error bar the 95% bootstrapped confidence interval. The line denotes the results of the simulations with the Usher McClelland model.

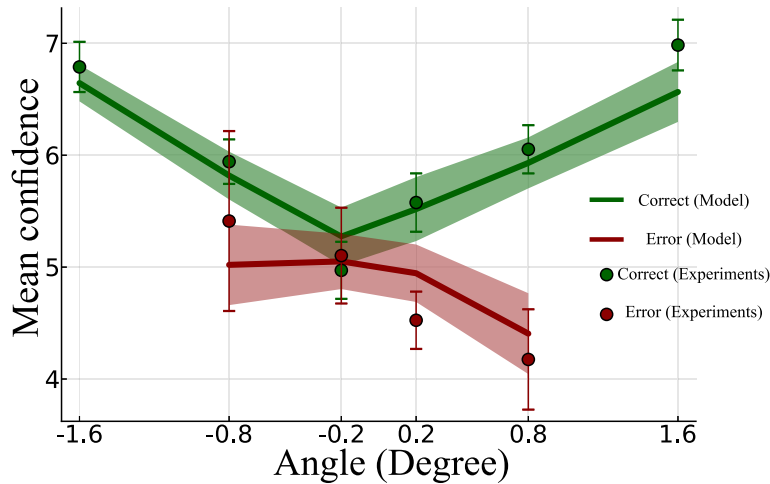


Figure 7: **Confidence as a function of stimulus strength.** We represent the mean confidence as a function of stimulus orientation in correct trials (green), and in error trials (red), for the experimental data (points) and the model simulations (lines). With parameters resulting from the fit on the confidence block, the numerical protocol mimic the experimental one (same number of trials, and same angle values). Due to the discrete levels of confidence, and the high performance in the task, to get enough statistics we combined the data of all subjects. The shaded areas (resp. error bars) denote the 95% bootstrapped confidence interval on the mean for the simulation (resp. data)

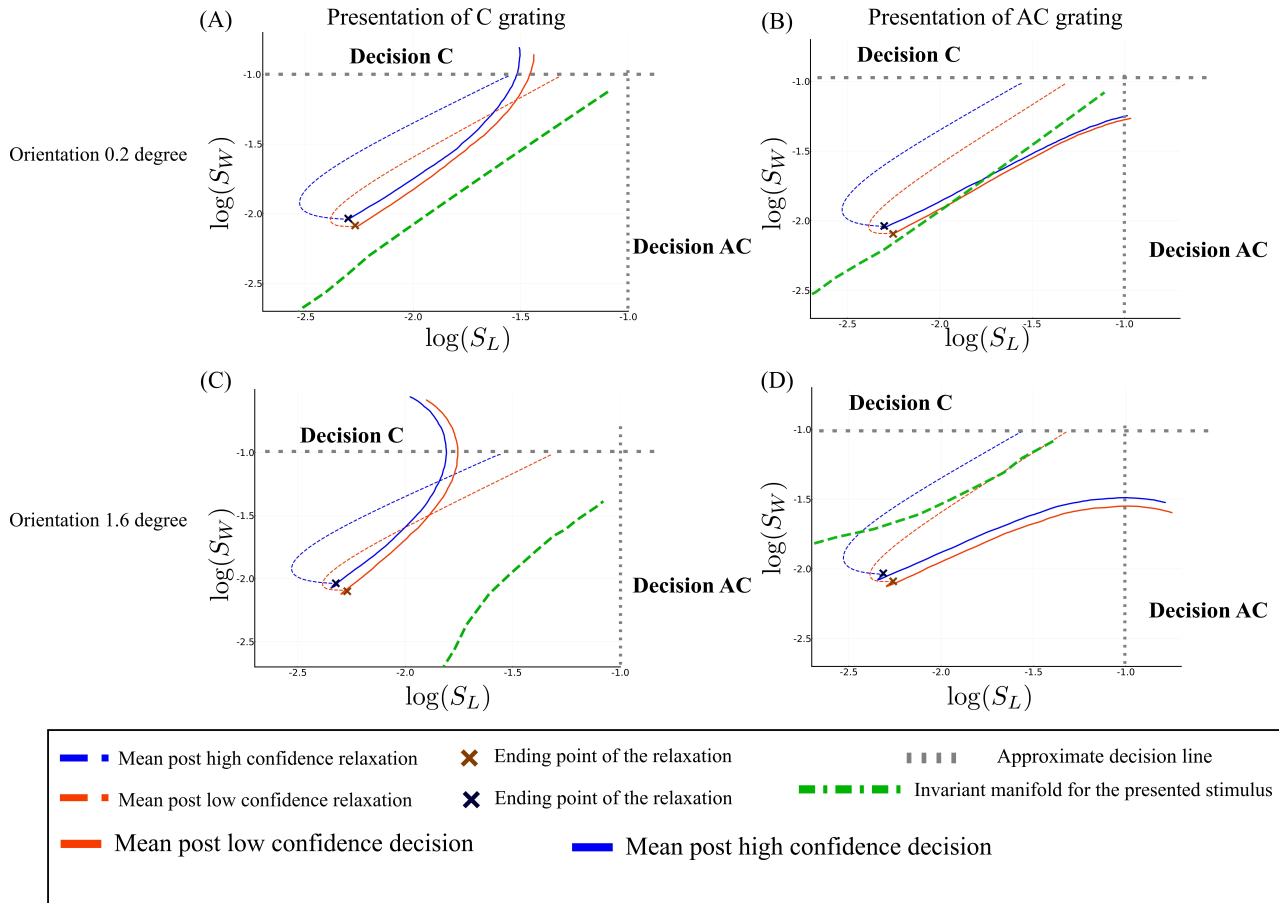
**Statistical analysis of sequential effects** Perceptual decisions made by humans in behavioral experiments depend not only on the current sensory input, but also on the choices made at previous trials. Various sequential effects have been reported [Fernberger, 1920, Laming, 1979, Leopold et al., 2002, Gold et al., 2008], and researchers have proposed different models to account for them [Cho et al., 2002, Angela and Cohen, 2009, Glaze et al., 2015, Bonaiuto et al., 2016, Berlemont and Nadal, 2019]. When the subject does not receive any feedback, confidence in his/her decision might be important for controlling future behaviors [Yeung and Summerfield, 2012, Meyniel et al., 2015b]. Recently, the effects of confidence on the history biases have been experimentally investigated [Braun et al., 2018, Samaha et al., 2018]. One main finding is that decisions with high confidence confer stronger biases upon the following trials. Here, we investigate the influence of confidence upon the next trial in the empirical data, and we show that the results are well reproduced by the behavior of the dynamical neural model.

First, we perform a statistical analysis of the effect of history biases on response times in the experimental data. For this, for each participant, we classify each trial into *low* and *high* confidence: a trial is considered as *low confidence* (resp. *high confidence*) if the reported confidence is below (resp. above) the participant’s median. We analyze the history biases making use of linear mixed effects models (LMM) [Gelman and Hill, 2007]. We find that higher orientations lead to faster response times (slope of  $-0.1$ ,  $p = 0.002$ ) and the repetition biases on response times [Cho et al., 2002] (coefficient of  $-0.034$ ,  $p < 2 \times 10^{-6}$ ). In line with previous works [Desender et al., 2018a], high confidence has the effect of speeding up the following trial (coefficient of  $-0.0206$ ,  $p = 0.0005$ ). Finally, we find that the previous response time has an effect on the subsequent one (coefficient of  $0.015$ ,  $p = 3.23 \times 10^{-5}$ ), meaning that the participants have the tendency to show sequences of fast (or slow) response times.

Next, following a numerical protocol replicating the experimental one, for each participant, we make numerical simulations with the model specifically calibrated on the participant’s data. Recall that the fit has been done on mean response times and accuracy, hence without taking into account serial dependencies. We then study the correlations between decisions made in successive trials by the neural attractor network, performing the same type of statistical analysis as done on the experimental data. We note that the attractor neural network captures the variation of response times with respect to angle orientation, as expected from Wong and Wang [2006] (coefficient of  $-0.017$ ,  $p = 9.47 \times 10^{-5}$ ). We find that the dependency in the choice history (through the repetition of responses), as observed in the experimental data, is correctly reproduced by the model (coefficient of  $-0.18$ ,  $p < 2 \times 10^{-6}$ ), in agreement with a previous study of these effects [Berlemont and Nadal, 2019]. Quite remarkably, we observe an effect of confidence on response times in the network, with negative slopes (coefficient of  $-0.02$ ,  $p = 0.005$ ) as in the experiment.

**Analysis of the underlying neural dynamics** To understand how the neural dynamics leads to these confidence-specific sequential effects, we make an analysis of the dynamics similar to the one done in Berlemont and Nadal [2019] for the analysis of post-error effects in the same neural model. We illustrate this analysis in Figure 8. On each panel, we compare the *mean* neural dynamics for post-low and post-high confidence trials (respectively red and blue lines). Without loss of generality, we assume that the previous decision was a *C* grating. We first note that the relaxation dynamics between two consecutive trials are different, resulting in different starting points for the next trial, from post-low and post-high confidence trials. Panel (A) corresponds to the case where the new stimulus is also *C* oriented ("repeated" case), at low strength level. The ending points of the relaxations fall into the correct basin of attraction. Because the post-high confidence relaxation lies deeper into the basin of attraction than the one of post-low trials, the subsequent dynamics will be faster for post-high confidence trials in this case. In panel (B) we represent the case, still at low stimulus strength, where the stimulus orientation of the new stimulus is the opposite ("alternated" case) to the one corresponding to the previous decision (hence an *AC* grating). Both dynamics lie close to the basin boundary of the two attractors, thus the dynamics are slow and there is no significant difference between post-low and post-high confidence trials. In panels (C) and (D) we represent the same situations as panels (A) and (B), respectively, but for high strength levels (easy trials). The ending points of the relaxations are far from the boundary of the basins of attraction, whatever the grating presented. The response times for post-high and post-low confidence trials are thus similar. This analysis shows that the non-linearity of the network dynamics is responsible for the considered sequential effect. Indeed, in the absence of non-linearity, the repeated and alternated case would compensate each other and there would be no specific effect related to the basin boundaries.

We now qualitatively confront the outcomes of the above analysis with the experimental data. To do so, we group the response times according to the same cases as previously: high and low stimulus strength, repeated or alternated trials. We compare post high and low confidence trials in each case, using a t-test [Fay and Proschan, 2010]. We find that mean response times between post low and high confidence trials are different in the low orientation stimuli and repeated case (t-test,  $p = 0.044$ ), but they are identical in the low orientation stimuli and alternated case, high orientation stimuli and alternated case, low orientation stimuli and repeated case they are identical (respectively  $p = 0.90$ ,  $p = 0.70$ ,  $p = 0.23$ ).



**Figure 8: Non linear dynamics in post-low and post-high confidence trials.** Phase-plane trajectories (in log-log plot, for ease of viewing) of the post low and high confidence trials. We assume that the previous decision was decision C. The axes represent the losing neural pool  $S_L$  and the winning neural pool  $S_W$  at the previous trial. The blue color codes for post-high confidence trials, and the red one for post-low confidence. Panels (A) and (B): Repeated and alternated case for low orientation stimuli; Panels (C) and (D): Repeated and alternated case for high orientation stimuli. In order to compare the decision times, the dynamics starting at the onset of the next stimulus is followed during 200ms, as if there were no decision threshold. The actual decision occurs at the crossing of the dashed gray line, indicating the threshold.

This is in accordance with the outcomes of the above analysis based on the non-linear dynamics.

The model reproduces sequential effects correlated with repetition and confidence, and we have shown that these effects result from the intrinsic nonlinear network dynamics. However the model does not reproduce the correlations of the response time with the previous response time (coefficient of  $-0.02$ ,  $p = 0.24$ ). This allows us to distinguish the effects that can be explained by the intrinsic dynamics of the attractor network, and the ones that would require the implementation of other cognitive processes.

## DISCUSSION

**Modeling confidence.** Dynamical models of decision making implement in different ways the same qualitative idea: decision between two categories is based on the competition between units collecting evidences in favor of one or the other category (or with a single unit whose activity represent the difference between the categorical evidences). Apart from very few works [Rolls et al., 2010a,b], authors propose that behavioral confidence can be modeled as a function of the balance of evidence [Vickers, 1979 (reeditited in 2014, Kepecs et al., 2008, Moreno-Bote, 2010, Pleskac and Busemeyer, 2010, Zylberberg et al., 2012, Kiani et al., 2014)]. We discuss the various approaches in light of our results, first by comparing the different models for confidence, then the specific effects of confidence on sequential decision-making.

Bayesian inference models compute confidence using drift-diffusion models (DDM) extensions based on decision variable balance [Vickers, 1979 (reeditited in 2014, Kepecs et al., 2008, Moreno-Bote, 2010, Zylberberg et al., 2012)], possibly with additional mechanisms - decision variable balance combined with response times [Kiani et al., 2014] or post-decisional deliberation [Pleskac and Busemeyer, 2010] (the dynamics continues after the decision, thus updating the balance of evidence). Similar studies have been made with independent race models (IRM) [Raab, 1962, Bogacz et al., 2006, Vickers, 1970, Merkle and Van Zandt, 2006]. These DDM or IRM models successfully account for various psychometric and chronometric specificities of human confidence. In DDMs, confidence based on decision variable balance predicts that confidence should deterministically decrease as a function of response times [Drugowitsch et al., 2012, Kiani and Shadlen, 2009]. However, Ratcliff and Starns [2009] have shown that the response times distributions strongly overlap across confidence levels. This property can be recovered making use of additional processes, such as with a two-stage drift-diffusion model [Pleskac and Busemeyer, 2010]. Yet, other effects remain unexplained within the framework of DDM. This is the case of early influence of sensory evidence on confidence [Zylberberg et al., 2012], as well as the fact that confidence is mainly influenced by evidence in favor of the selected choice [Zylberberg et al., 2012].

In order to model experiments in monkeys [Kiani and Shadlen, 2009], Wei and Wang [2015] make use of a ring attractor neural network with confidence computed from the balance of evidence. In the present paper, their approach has been extended to the case of a two-variables attractor network model, taking into account an inhibitory feedback allowing the network to engage in a sequence of trials [Berlemont and Nadal, 2019]. The reported confidence is modelled as a function of the difference in activity between the winning and loosing populations at the time of decision. The asymmetries mentioned above in the influence of evidence on confidence automatically arise in the accumulation of evidence in attractor neural networks [Wong et al., 2007]. We expect that these asymmetries will arise the same way on confidence, as we model confidence as a function of the balance of evidence.

Here we have shown, for the first time, that an attractor neural network can account for sequences of decisions and reproduce response times, accuracy and confidence individually for each participant. In addition, the model accounts for other effects reported in the literature, among which there were some effects that were believed to be specific signatures of Bayesian confidence [Sanders et al., 2016, Adler and Ma, 2018]: confidence increases with the orientation amplitude for correct trials but decreases for error trials.

**Confidence and Serial dependence.** The most common serial dependence effect in perceptual decision-making is the fact that the current stimulus appears more similar to recently seen stimuli than it really is [Cho et al., 2002, Fecteau and Munoz, 2003]. This effect can be observed for a large variety of perceptual features such as luminance [Fründ et al., 2014], orientation [Fischer and Whitney, 2014], direction of motion [Cho et al., 2002] or face identity [Lieberman et al., 2014]. A recent finding is that the magnitude of history biases increases when previous trials were faster and correct. Within the Signal Detection Theory framework, making use of the correlations between confidence, response time and accuracy, this effect is interpreted as an impact of confidence on the next decision [Braun et al., 2018]. By measuring directly the subjective confidence of the participants, recent studies confirm that confidence modulates the history biases [Desender et al., 2018a, Samaha et al., 2018, Desender et al., 2018b]. In our experiment, we observe that high confidence trials lead to faster subsequent choices in agreement with the above mentioned experimental studies.

On the theoretical side, the impact of confidence on response times of the subsequent trials has been investigated within the framework of DDMs [Desender et al., 2018a]. The trials are divided into two categories, subsequent to low or high

confidence trials, then a DDM is fitted separately on each type of trial. This amounts to assume that parameters (threshold and drift) are changed depending on the confidence level at the previous trial. Actually, without any change of parameters, the DDM or IRM models cannot account for the observed sequential effects, as discussed in Appendix B: the predicted sequential effect would be the opposite of the observed one.

In contrast to the DDM analysis, we calibrate the attractor model globally for each participant, and reproduce the sequential effects with a unique set of model parameters. This can be understood as resulting from the intrinsic nonlinear dynamics, as discussed in the Results section (see Figure 8). Hence, the attractor network model does not only account for the relationship between confidence, response times and accuracy, but also reproduces the influence of confidence on serial dependence.

**Decision and non decision times.** Human studies commonly report right-skewed response times distribution [Ratcliff, 1978, Luce et al., 1986, Ratcliff and Rouder, 1998]. Such long right tails are well captured by drift-diffusion models [Ratcliff and Rouder, 1998, Usher and McClelland, 2001] - and this is generally considered as a strong evidence in favor of the accumulation of evidence mechanism. However, with trained subjects, the right-skew is less pronounced and the response times distribution can be accurately reproduced by a Gaussian distribution [Peirce, 1873]. In contrast to human studies, experiments in monkeys do not show such long right tails in the response times histograms [Ditterich, 2006]. When assuming a constant value for the non-decision time, attractor neural network models cannot account for the right-skewed distributions, but accurately reproduce the shape of the distributions in monkeys experiments [Wang, 2008]. In accordance with these results, in this work we have shown that for the range of parameters we considered, the decision time distribution generated by the neural network can be approximated by a Gaussian distribution.

Within the neural attractor framework, the experimentally observed long right tails can thus be understood as originating only from the non-decision times. Here we have proposed an estimation of the distribution of these non-decision time allowing to fit the empirical response times distributions.

One should note that, even in the case of an analysis of experimental data within the DDM framework, the estimated non-decision times are not necessarily given by a constant value, but may show a distribution with a strong right skew, as shown by Verdonck and Tuerlinckx [2016]. These findings combined with ours suggest that the question of the origin of the long right tails in human response times has to be reconsidered.

To conclude, in this work, we designed a specific experiment in order to study confidence with human participants. We fitted a neural attractor network model specifically to each participant in order to describe their behavioral results in continuous sequences of perceptual decisions: response times, accuracy and confidence. Finally, we have shown that the impact of confidence on sequential effects is well described by the intrinsic nonlinear dynamics of the network.

## MATERIAL AND METHODS

### PARTICIPANTS

Nine participants (7 Females, Mean Age = 27.3, SD = 5.14) have been recruited from the Laboratoire de Psychologie cognitive et de Psycholinguistique's database (LSCP, DEC, ENS-EHESS-CNRS, PSL, Paris, France). Every subject had normal or corrected-to-normal vision. We obtained written informed consent from every participant who received a compensation of 15 euros for their participation. The participants performed three sessions on three distinct days in the same week for a total duration of about 2h15. The experiment followed the ethics requirements of the Declaration of Helsinki (2008) and has been approved by the local Ethics Committee. Three participants were excluded. Two of the excluded participants did not complete correctly the experiment and one exhibited substantially asymmetric performance (98% of correct responses for an angle of  $0.2^\circ$ , but 18% at  $-0.2^\circ$  degree). As a result, we analyzed data from 6 participants.

### STIMULI AND TASKS

The stimuli were generated using Matlab [MATLAB, 2016] along with the Psychophysics Toolbox [Kleiner et al., 2007]. They were displayed on a monitor at 57.3 cm of the participants' head. The participants performed the experiment in a quiet and darkened experimental room. Their heads were stabilized thanks to a chin-rest. Trials began with the presentation of a black fixation point (duration = 200 ms). Then the stimulus for the primary decision task was presented, consisting in a circular grating (diameter =  $4^\circ$ , Tukey window, 2 cycles per degree, Michelson contrast = 89%, duration = 100 ms, phase randomly selected at each trial). The grating had eight possible orientations with respect to the vertical meridian, and participants were asked to categorize them as clockwise or anti-clockwise with respect to the vertical meridian by pressing the right-arrow or left-arrow. Participants had been instructed to respond as follows: "You have to



respond quickly but not at the expense of precision. After 1.5 s the message, "Please answer", will appear on the screen. It would be really ideal, if you would answer before this message appears."

Trials were of three types, grouped in *pure* block, *feedback* block and *confidence* block (see below). Participants performed three sessions on three distinct days. Each session (45 min) consisted in three runs, each run being composed of one exemplar of each of the three type of block, in a random order. Before starting the experiment, participants performed a short training block of each type, with easier orientations than in the main experiment.

**Pure block** In this block, participants waited 300 ms after each decision, before the black fixation point appears. The stimulus appeared 200 ms after this fixation point. The eight possible orientations for the circular grating were  $[-1.6^\circ, -0.8^\circ, -0.5^\circ, -0.2^\circ, 0.2^\circ, 0.5^\circ, 0.8^\circ, 1.6^\circ]$  and a stimulus was chosen randomly among them with the following weights: [0.05, 0.1, 0.15, 0.2, 0.2, 0.15, 0.1, 0.05].

**Feedback block** In this block, 200 ms after the decision, the participants received an auditory feedback (during 200 ms) about the correctness of the decision they just made. The black fixation dot appeared 100 ms after this feedback and a new trial began. The orientations of the circular gratings were chosen randomly from  $[-1.6^\circ, -0.8^\circ, -0.2^\circ, 0.2^\circ, 0.8^\circ, 1.6^\circ]$  with the following weights [ 0.12, 0.18, 0.2, 0.2, 0.18, 0.12 ].

**Confidence block** In the confidence block, participants had to evaluate the confidence on the orientation task 200 ms after the decision. To perform this task they had to move a slider on a 10-points scale, from *pure guessing* to *certain to be correct*. Importantly, the initial position of the slider was chosen randomly for each trial. Participants moved the slider to the left by pressing the "q" key, and to the right with the "e" key. We ask the following kind of confidence judgment to the participants: one extreme of the scale is "pure guess", the other is "absolutely certain". They confirmed the choice of the value of confidence by pressing the space bar. The participants had the choice to indicate that they had made a "motor mistake" during the orientation task. For this they had to press a key with a red sticker instead of responding on the confidence scale. After the choice of confidence, the participants had to wait 300 ms before the black fixation dot appears. After the fixation dot the stimulus appeared 200 ms later. The orientations of the circular gratings were the same as in the feedback block.

Accuracy is higher and response times are slower in confidence blocks than in pure blocks, an effect already observed in previous studies ( **J-R. Martin, J. Sackur, *personnal communication, April, 2018***). To test this effect on accuracy we ran a binomial regression of responses with fixed factors of orientation and type of block (pure or confidence), the interaction between these factors and a random participant intercept. The orientation coefficient was 2.15 (SD = 0.17,  $z = 12.44$  and  $p < 10^{-16}$ ); there was no effect of block type ( $p = 0.385$ ). But we found a significant orientation by block type interaction (value of 0.55, SD = 0.08,  $z = 6.97$  and  $p = 3 \cdot 10^{-12}$ ), to the effect that participants were more accurate in confidence blocks than in no-confidence blocks. In a similar way, we test the effect on response times by using a mixed effect regression with the same factors and intercept as for the accuracy (only on the absolute value of the orientation). We found that the orientation coefficient (value of  $-0.08$ , SD = 0.013 and  $p = 0.0006$ ) and the block type coefficient (value of 0.095, SD = 0.028 and  $p = 0.011$ ) were significant, meaning that participant are slower in the confidence block. Moreover, the slope by block type interaction with orientation was also significant (value of  $-0.028$ , SD = 0.010 and  $p = 0.031$ ), meaning that the difference between the two types of blocks is more important at low orientation. Surprisingly, we find that performance and response time across participants are identical in the feedback and pure blocks (no statistically significant difference). The participants were highly trained in the orientation discrimination task.

## STATISTICAL ANALYSES

We used RStudio [RStudio Team, 2015] with the package *lme4* [Bates et al., 2015] to perform a linear mixed effects analysis Gelman and Hill [2007] of the history biases of the reaction times. The linear mixed effects model (LMM) we consider assumes that the logarithm of the response time at trial  $n$ ,  $RT_n$ , is a linear combination of factors as follows:

$$\ln(RT_n) = a_{0,p} + a_{1,p}|\theta| + a_{2,x_{\text{repetition}}} + a_{3,p} \ln(RT_{n-1}) + a_{4,\text{Conf}_{n-1}} \quad (1)$$

with  $x_{\text{repetition}}$  a binary variable taking the value 1 if the correct choice for the current trial is a repetition of the previous choice (and 0 otherwise),  $\theta$  the orientation of the Gabor (in degree),  $RT_{n-1}$  the response times of the previous trials, and  $\text{Conf}_{n-1}$  the confidence of the previous trial coded as 0 for *low* and 1 for *high*. The subscript  $p$  in a coefficient (e.g  $a_{0,p}$ ) indicates that for this parameter we allow for a random slope per participant.

We compared this LMM to other ones that do not include all these terms, using the ANOVA function (with the *lme4* package Bates et al. [2015]) that performs model comparison based on the Akaike and Bayesian Information Criteria (AIC and BIC) [Bates et al., 2014]. As we can note in Table 1 the LMM from Eq. 1 is preferable in all cases.

$a_{0,p} + a_{1,p} \theta  + a_2x_{\text{repetition}} + a_{3,p} \ln(RT_{n-1}) + a_4\text{Conf}_{n-1}$	Df	AIC	BIC	LogLik.	p value
$a_{0,p} + a_{1,p} \theta  + a_2x_{\text{repetition}} + a_{3,p} \ln(RT_{n-1})$	12	-335	-254	180	
$a_{0,p} + a_{1,p} \theta  + a_2x_{\text{repetition}}$	11	-324	-249	173	0.0003
$a_0 + a_1 \theta  + a_2x_{\text{repetition}} + a_3 \ln(RT_{n-1}) + a_4 \ln(\text{Conf}_{n-1})$	7	-4	-42	9	<2e-16
$a_0$	7	-225	-177	119	<2e-16
	3	-475	-495	-234	<2e-16

**Table 1: LMM tests on Data, models comparison.** The first row gives the tests for the LMM from Eq. 1. The p-values are for the tests based on BIC and AIC [Bates et al., 2014] between the LMM from Eq. 1 and the one of the corresponding row.

The analysis of the linear mixed effect model applied to the experimental data are described in the Results section. However, we summarize them in the two following tables for clarity purposes:

	Estimate	Std. Error	df	t-value	Pr	
$a_{0,p}$	5.428	1.466e-01	9.0	37.038	$6.92 \cdot 10^{-11}$	***
$a_{1,p}$	-0.1027	0.02390	9.0	-4.296	0.002001	**
$a_2$	$-3.402 \cdot 10^{-2}$	$2.840 \cdot 10^{-3}$	$8.472 \cdot 10^3$	-11.978	$< 2 \cdot 10^{-16}$	***
$a_{3,p}$	$1.517 \cdot 10^{-1}$	$1.651 \cdot 10^{-2}$	7.0	9.187	$3.23 \cdot 10^{-5}$	***
$a_4$	$-2.063 \cdot 10^{-2}$	$5.969 \cdot 10^{-3}$	$5.537 \cdot 10^3$	-3.456	0.000553	***

**Table 2: Results of the application of the LMM from Eq. 1 on the experimental data.** We note \*\* for  $p < 0.005$  and \*\*\* for  $p < 0.001$ .

	Estimate	Std. Error	df	t value	Pr	
$a_{0,p}$	5.999	0.08032	4.229	74.690	$9.22 \cdot 10^{-8}$	***
$a_{1,p}$	-0.01744	$5.551 \cdot 10^{-4}$	2.886	-31.420	$9.47 \cdot 10^{-5}$	***
$a_2$	-0.1814	$8.133 \cdot 10^{-3}$	$4.822 \cdot 10^3$	-22.301	$< 2 \cdot 10^{-16}$	***
$a_{3,p}$	-0.02075	$1.545 \cdot 10^{-2}$	4.628	-1.343	0.24139	
$a_4$	-0.02324	$8.336 \cdot 10^{-3}$	$4.847 \cdot 10^3$	-2.788	0.00533	**

**Table 3: Results of the application of the LMM from Eq. 1 on the data from the neural network simulations.** We note \*\* for  $p < 0.005$  and \*\*\* for  $p < 0.001$ .

To perform the comparison between the experimental data and the results of Figure 8 we first transform the response times of each participant using the z-score [Kreyszig, 1979]. This allows us to study all participants together as the response times are now normalized.

#### ATTRACTOR NEURAL NETWORK MODEL

We consider the decision-making recurrent network model governed by local excitation and feedback inhibition introduced in Berlemont and Nadal [2019] as an extension of the model proposed by Wong and Wang [2006].

Within a mean-field approach, Wong and Wang [2006] have derived a reduced firing-rate model of the full biophysical models of spiking neurons introduced and studied in Compte et al. [2000]. This reduced model is composed of two interacting neural pools which faithfully reproduces not only the behavioral behavior of the full model, but also the dynamics of the neural firing rates and of the output synaptic gating variables. The details can be found in Wong and Wang [2006] (main text and Supplementary Information). This model and its variants are used as proxies to simulate the full spiking network and for getting mathematical insights [Wong and Wang, 2006, Engel and Wang, 2011, Miller and Katz, 2013, Deco et al., 2013, Engel et al., 2015, Berlemont and Nadal, 2019]. The model variant that we consider here takes into account a corollary discharge [Sommer and Wurtz, 2008, Crapse and Sommer, 2009]. This results in an inhibitory current injected into the neural pools just after a decision is made, making the neural activities relax towards a low activity, neutral, state, therefore allowing the network to deal with consecutive sequences of decision making trials. The full details can be found in Berlemont and Nadal [2019]. We remind here the equations and parameters with notation adapted to the present study.

The model consists of two competing units, each one representing an excitatory neuronal pool, selective to one of the two categories,  $C$  or  $AC$ . The dynamics is described by a set of coupled equations for the synaptic activities  $S_C$  and  $S_{AC}$  of



the two units  $C$  and  $AC$ :

$$i \in \{C, AC\}, \quad \frac{dS_i}{dt} = -\frac{S_i}{\tau_S} + (1 - S_i) \gamma f(I_{i,tot}) \quad (2)$$

The synaptic drive  $S_i$  for pool  $i \in \{C, AC\}$  corresponds to the fraction of activated NMDA conductance, and  $I_{i,tot}$  is the total synaptic input current to unit  $i$ . The function  $f$  is the effective single-cell input-output relation [Abbott and Chance, 2005], giving the firing rate as a function of the input current:

$$r_i = f(I_{i,tot}) = \frac{aI_{i,tot} - b}{1 - \exp[-d(aI_{i,tot} - b)]} \quad (3)$$

where  $a, b, d$  are parameters whose values are obtained through numerical fit. The total synaptic input currents, taking into account the inhibition between populations, the self-excitation, the background current and the stimulus-selective current, can be written as:

$$I_{C,tot} = J_{C,C}S_C - J_{C,AC}S_{AC} + I_{stim,C} + I_{noise,C} + I_{CD}(t) \quad (4)$$

$$I_{AC,tot} = J_{AC,AC}S_{AC} - J_{AC,C}S_C + I_{stim,AC} + I_{noise,AC} + I_{CD}(t) \quad (5)$$

with  $J_{i,j}$  the synaptic couplings. The minus signs in the equations make explicit the fact that the inter-units connections are inhibitory (the synaptic parameters  $J_{i,j}$  being thus positive or null). The term  $I_{stim,i}$  is the stimulus-selective external input. The form of this stimulus-selective current is:

$$I_{stim,i} = J_{ext} (1 \pm c_\theta) \quad (6)$$

with  $i = C, AC$ . The sign,  $\pm$ , is positive when the stimulus favors population  $C$ , negative in the other case. Here the parameter  $J_{ext}$  combines a synaptic coupling variable and the global strength of the signal (which are parametrized separately in the original model [Wong and Wang, 2006, Berlemont and Nadal, 2019]). The quantity  $c_\theta$ , between 0 and 1, characterizes the stimulus strength in favor of the actual category, here an increasing function of the (absolute value of) the stimulus orientation angle,  $\theta$ .

In addition to the stimulus-selective part, each unit receives individually an extra noisy input, fluctuating around the mean effective external input  $I_0$ :

$$\tau_{noise} \frac{dI_{noise,i}}{dt} = -(I_{noise,i}(t) - I_0) + \eta_i(t) \sqrt{\tau_{noise} \sigma_{noise}} \quad (7)$$

with  $\tau_{noise}$  a synaptic time constant which filter the white-noise.

After each decision, a corollary discharge under the form of an inhibitory input is sent to both units until the next stimulus is presented:

$$I_{CD}(t) = \begin{cases} 0 & \text{during stimulus presentation} \\ -I_{CD,max} \exp(-(t - t_D)/\tau_{CD}) & \text{after the decision time, } t_D \end{cases} \quad (8)$$

This inhibitory input, delivered between the time of decision and the presentation of the next stimulus, allows the network to escape from the current attractor and engage in a new decision task [Berlemont and Nadal, 2019].

**Confidence modeling** Within the various decision making modelling frameworks, similar proposals have been made to model the neural correlate of the behavioral confidence level. In race models [Raab, 1962], which have equal number of accumulation variables and stimulus categories, as in attractor network models, the balance of evidence at the time of perceptual decisions has been used to model the neural correlate of the behavioral confidence [Vickers and Packer, 1982, Smith and Vickers, 1988, Wei and Wang, 2015]. This balance of evidence is given by the absolute difference between the activities of the category specific units at the time of decision. Here, we follow Wei and Wang [2015], considering that confidence is obtained as a function  $f$  of the difference in neural pools activities,  $\Delta r = |r_C - r_{AC}|$ .

In our experiment, the subjects expressed their confidence level by a number on a scale from 0 to 9. In order to match the neural balance of evidence with the confidence reported by the subject, we map the balance of evidence histogram onto the behavioral confidence histogram, a procedure called *histogram matching* [Gonzalez et al., 2002]. Note that the mapping is here from a continuous variable to a discrete one (taking integer values from 0 to 9).

## FITTING PROCEDURE

We perform a model calibration in order to fit the behavioral data of each participant. More precisely, we calibrate the model by fitting, for each participant, both the mean response times and the accuracies for each orientation, this separately for each block. We note that we only fit the means, which in particular implies that the fits do not take into account the

Parameter	Value	Parameter	Value
a	270 Hz/nA	$\sigma_{noise}$	0.02 nA
b	108 Hz	$\tau_{noise}$	2 mS
d	0.154 s	$I_0$	0.3255 nA
$\gamma$	0.641	$J_{ext}$	0.182 nA
$J_{C,C} = J_{AC,AC}$	0.2609 nA	$J_{C,AC} = J_{AC,C}$	0.0497 nA
$\tau_S$	100 ms		

Table 4: Numerical values of the model parameters taken from Berlemont and Nadal [2019] (common to all participants).

Parameter	Status
$I_{CD,max}$	common to all participants
$\tau_{CD}$	common to all participants
threshold $z$	specific to each participant
$c_\theta$ , for $\theta = \{0.2^\circ, 0.5^\circ, 0.8^\circ, 1.6^\circ\}$	specific to each participant

Table 5: List of parameters subject to calibration, separately for the pure and confidence blocks.

serial dependencies. Doing so, any sequential effects that will arise in the model will result from the intrinsic dynamics of the network, and not from a fitting procedure of these effects.

The model parameters values are those used in Berlemont and Nadal [2019], and reproduced in Table 4, except for the few parameters listed in Table 5, that is  $I_{CD,max}$ ,  $\tau_{CD}$ ,  $c_\theta$  and  $z$ , which are chosen in order to fit the data. The two parameters  $I_{CD,max}$  and  $\tau_{CD}$  are imposed common to all participants (joint optimization). The parameters  $c_\theta$  (one for each orientation value) and  $z$  are optimized across subjects and blocks.

The observed response time is the sum of a decision time and of a non decision time. Assuming no correlation between these two times, the mean non decision time is thus independent of the orientation. For comparing data with model simulations (which only gives a decision time) at any given orientation  $\theta$ , we first subtract to the mean response time the mean response time averaged over all orientations (this for both data and simulations). We calibrate the model parameters so as to fit these centered mean response times. This will provide a fit of the mean response times (at each angle) up to a global constant, which is the mean non decision time (the modeling of the non decision time distribution is presented in the next Section).

For each participant, and each block, we thus consider the cost function:

$$\begin{aligned} \text{Cost function} = & \alpha \frac{1}{m} \sum_{\theta} ([\langle RT \rangle_{network}(\theta) - \langle RT \rangle_{network}] - [\langle RT \rangle_{data}(\theta) - \langle RT \rangle_{data}])^2 \\ & + \frac{1}{n} \sum_{\theta} ((\langle accuracy \rangle_{network}(\theta) - \langle accuracy \rangle_{data}(\theta))^2 \end{aligned} \quad (9)$$

where the sums are over the orientation values,  $\theta = \{0.2^\circ, 0.5^\circ, 0.8^\circ, 1.6^\circ\}$ , the brackets  $\langle \dots \rangle$  design averages (as detailed below), and the normalization factors  $n$  (for response times) and  $m$  (for the accuracy) are given by

$$\begin{aligned} m &= \max_{\theta} ([\langle RT \rangle_{network}(\theta) - \langle RT \rangle_{network}] - [\langle RT \rangle_{data}(\theta) - \langle RT \rangle_{data}])^2 \\ n &= \max_{\theta} [(\langle accuracy \rangle_{network}(\theta) - \langle accuracy \rangle_{data}(\theta))^2 \end{aligned}$$

In these expression,  $\langle RT \rangle_{data}(\theta)$  denotes the mean experimental response time obtained by averaging over all trials at the orientations  $\pm\theta$ ,  $\langle RT \rangle_{data}$  is the average over all orientations;  $\langle RT \rangle_{network}(\theta)$  and  $\langle RT \rangle_{network}$  are the corresponding averages obtained from the model simulations. The coefficient  $\alpha$  denotes the relative weight given to the response time and accuracy cost terms. We present the results obtained when taking  $\alpha = 2$ , but it should be noted that the choice of this parameter does not impact drastically the fitted parameters. Finally we add a soft constraint on  $c_\theta$  so that this value does not diverge when the participant accuracy is close to 100%. Note that this constraint does not affect the results of the fitting procedure, but is necessary in order to compute the confidence interval of the fitted parameters.

For each subject, we minimize this cost function with respect to the choice of  $c_\theta$  and  $z$ , making use of a Monte Carlo Markov Chain fitting procedure, coupled to a subplex procedure [Rowan, 1990]. This method is adapted to handle simulation based models with stochastic dynamics [Bogacz and Cohen, 2004]. Finally,  $I_{CD,max}$  and  $\tau_{CD}$  are fitted using

a grid search algorithm as they have less influence on the cost function. In the model, the parameter  $c$  represents the stimulus ambiguity, which we expect here to be a monotonous function of the amplitude of the angle,  $\theta$ . When allowed to be independent parameter values for each value of the orientation,  $\theta = \{0.2^\circ, 0.5^\circ, 0.8^\circ, 1.6^\circ\}$ , we find that the  $c_\theta$  values can be approximated by a linear or quadratic function of  $\theta$  depending on the participant. We performed an AIC test [Akaike, 1992] between the linear and quadratic fit in order to choose which function to use for each participant. These approximations reduce the number of free parameters.

Block	Parameter	Participant 1	Participant 2	Participant 3	Participant 4	Participant 5	Participant 6
Pure Block	$z(Hz)$	10.78	12.69	14.07	12.80	10.05	12.89
	$\Delta z(Hz)$	(-0.7,+1.75)	(-2.1,+0.175)	(-1.92,+1.75)	(-0.1,+2.275)	(-1.8,2.1)	(-1.05,+2.45)
Confidence Block	$z(Hz)$	13.08	13.70	14.95	12.96	12.55	14.65
	$\Delta z(Hz)$	(-0.18,+2.28)	(-1.93,+0.18)	(-2.45,+1.40)	(-0.17,+2.1)	(-0.35,+2.8)	(-2.1,+0.53)

Table 6: Threshold parameter for each participant after fit of the mean accuracy and response times of the pure and confidence blocks. The ranges  $\Delta z$  correspond to one sigma deviation of the likelihood with respect to the corresponding parameter (see the Material and Methods Section).

	Type of fit	Pure Block	Confidence Block
Participant 1	Quadratic	$c_\theta = 554.8 * \theta - 1444 * \theta^2$	$c_\theta = 650 * \theta + 1.3 \times 10^4 * \theta^2$
Participant 2	Quadratic	$c_\theta = 729.6 * \theta - 9819 * \theta^2$	$c_\theta = 1056 * \theta - 1.4 \times 10^4 * \theta^2$
Participant 3	Linear	$c_\theta = 524.4 * \theta$	$c_\theta = 634.6 * \theta$
Participant 4	Quadratic	$c_\theta = 269.8 * \theta - 577.6 * \theta^2$	$c_\theta = 190 * \theta + 2.17 \times 10^4 * \theta^2$
Participant 5	Quadratic	$c_\theta = 387.6 * \theta + 4188 * \theta^2$	$c_\theta = 182.4 * \theta + 4.9 \times 10^4 * \theta^2$
Participant 6	Linear	$c_\theta = 551 * \theta$	$c_\theta = 1030 * \theta$

Table 7: Calibration of the stimulus strength parameter: best fit for each participant.  $\theta$ , in radian, stands for the absolute value of the orientation.

In order to obtain a confidence interval for the different parameters, we used the likelihood estimation of confidence interval for Monte-Carlo Markov Chains method [Ionides et al., 2016]. The confidence interval on the parameters is thus the 70% confidence interval, assuming a Gaussian distribution of the cost function. This provides an approximation of the reliability of the parameters values found. In order to assess the reliability of this method we checked that the threshold  $z$  and stimulus strength  $c_\theta$  parameters have an almost non-correlated influence onto the cost function.

The results of the calibrating procedure are summarized in Tables 6 and 7, with  $I_{CD,max} = 0.033$  nA and  $\tau_{CD} = 150$  ms. We note that we use the linear or quadratic approximation for  $c_\theta$  instead of the exact values given by the MCMC procedure. Thus, the fits on Figure 2 are not the exact ones, but this approximation has the advantage of giving us the correspondence between all stimuli orientation and  $c_\theta$ .

#### ESTIMATING THE NON-DECISION TIME

The above fitting procedure calibrates the mean response times up to a global constant, corresponding to the mean non decision time. As explained in the main text, we can go beyond and actually model the non-decision time distribution.

The non-decision time is considered to be due to encoding and motor execution [Luce et al., 1986]. Most model-based data analysis of response time distributions assume a constant non-decision time [Ratcliff and Rouder, 1998, Usher and McClelland, 2001, Wong and Wang, 2006, Brown and Heathcote, 2008]. However, Ratcliff [2013] have shown that fitting data originating from a skewed distribution under the assumption of a nonskewed non-decision time distribution is cause for bias in the parameter estimates if the model for non-decision time is not correct. Recently, Verdonck and Tuerlinckx [2016] proposed a mathematical method to fit a non-parametrical non-decision time distribution. Analyzing various experimental data with this method within the framework of drift-diffusion models, they find that strongly right skewed non-decision time distributions are common.

In this paper we make the hypothesis that the non-decision time distributions are ex-Gaussian distributions [Grushka, 1972], whose parameters are inferred from the data making use of the deconvolution method introduced in Verdonck and Tuerlinckx [2016] and detailed in Appendix B. We present in Figure 3 the fits of the response time distributions and the inferred non decision time distributions.

We compare the fit with the attractor neural network with the fit with the Usher and McClelland model [Usher and McClelland, 2001]. The equations of the model are the following:

$$\begin{aligned}\tau dx_1 &= -kx_1 dt - \beta f(x_2) dt + I_1 + \sigma \mu_1(t) \\ \tau dx_2 &= -kx_2 dt + \beta f(x_1) dt + I_2 + \sigma \mu_2(t)\end{aligned}$$

with  $\mu_i(t)$  a white-noise process and  $I_i$  the input current to the system. The external input is defined as  $I_i = 0.5 \pm c_\theta$ , with  $c_\theta$  the strength per angle as in the attractor neural network.  $\sigma = 0.4$  denotes the strength of the noise,  $k$  the relaxation strength,  $\tau = 0.1$  the relaxation time and  $\beta$  the inhibitory term. Finally, the function  $f$  is a sigmoidal function of gain  $G = 0.4$  and half-activity offset  $d = 0.5$ ,  $f(x_i) = 1/[1 + \exp(-G(x_i - d))]$ . The dynamics occurs until a threshold  $z$  is reached for one of the two units. It should be noted that, despite the non-linearity, the Usher-McClelland model is closer to drift-diffusion models than to biophysical attractor model (this because the only non-linearity is in the interaction between units). Reductions to one-dimensional drift diffusion models can be made in various ranges of parameters [Bogacz et al., 2006]. In order to fit this model to the experiments, we apply the same procedure as for our attractor network model.

Parameter	Subject 1	Subject 2	Subject 3	Subject 4	Subject 5	Subject 6
$z$	1.0	1.0	1.0	1.4	1.3	1.3
$\beta$	0.25	0.10	0.18	0.10	0.15	0.12
$k$	0.15	0.18	0.18	0.11	0.11	0.14
$c_{0.2}$	0.02	0.04	0.02	0.02	0.02	0.04
$c_{0.8}$	0.15	0.12	0.07	0.08	0.14	0.132
$c_{1.6}$	0.23	0.20	0.17	0.225	0.295	0.235

Table 8: Calibration of the Usher-McClelland model: Parameters for each participant after fit of the mean accuracy and response times of the confidence block.

## ACKNOWLEDGMENTS

We are grateful to Laurent Bonnasse-Gahot for useful discussions and suggestions. We thank Pascal Mamassian, Vincent de Gardelle and Xiao-Jing Wang for stimulating discussions. We thank Isabelle Brunet for her help in recruiting the participants and organizing the experimental sessions. KB acknowledges a fellowship from the ENS Paris-Saclay.

## REFERENCES

- L. Abbott and F. S. Chance. Drivers and modulators from push-pull and balanced synaptic input. *Progress in brain research*, 149:147–155, 2005.
- W. T. Adler and W. J. Ma. Comparing bayesian and non-bayesian accounts of human confidence reports. *PLoS computational biology*, 14(11):e1006572, 2018.
- H. Akaike. Information theory and an extension of the maximum likelihood principle. In *Breakthroughs in statistics*, pages 610–624. Springer, 1992.
- J. Y. Angela and J. D. Cohen. Sequential effects: superstition or rational behavior? In *Advances in neural information processing systems*, pages 1873–1880, 2009.
- J. V. Baranski and W. M. Petrusic. The calibration and resolution of confidence in perceptual judgments. *Perception & psychophysics*, 55(4):412–428, 1994.
- D. Bates, M. Mächler, B. Bolker, and S. Walker. Fitting linear mixed-effects models using lme4. *arXiv preprint arXiv:1406.5823*, 2014.
- D. Bates, M. Mächler, B. Bolker, and S. Walker. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1):1–48, 2015. doi: 10.18637/jss.v067.i01.
- J. M. Beck, W. J. Ma, R. Kiani, T. Hanks, A. K. Churchland, J. Roitman, M. N. Shadlen, P. E. Latham, and A. Pouget. Probabilistic population codes for bayesian decision making. *Neuron*, 60(6):1142–1152, 12 2008. ISSN 0896-6273. doi: 10.1016/j.neuron.2008.09.021. URL <http://https://doi.org/10.1016/j.neuron.2008.09.021>.

- K. Berlemont and J.-P. Nadal. Perceptual decision-making: Biases in post-error reaction times explained by attractor network dynamics. *Journal of Neuroscience*, 39(5):833–853, 2019. ISSN 0270-6474. doi: 10.1523/JNEUROSCI.1015-18.2018. URL <http://www.jneurosci.org/content/39/5/833>.
- R. Bogacz and J. D. Cohen. Parameterization of connectionist models. *Behavior Research Methods, Instruments, & Computers*, 36(4):732–741, 2004.
- R. Bogacz, E. Brown, J. Moehlis, P. Holmes, and J. D. Cohen. The physics of optimal decision making: a formal analysis of models of performance in two-alternative forced-choice tasks. *Psychological review*, 113(4):700, 2006.
- J. J. Bonaiuto, A. de Berker, and S. Bestmann. Response repetition biases in human perceptual decisions are explained by activity decay in competitive attractor models. *eLife*, 5:e20047, 2016.
- A. Braun, A. E. Urai, and T. H. Donner. Adaptive history biases result from confidence-weighted accumulation of past choices. *Journal of Neuroscience*, pages 2189–17, 2018.
- S. D. Brown and A. Heathcote. The simplest complete model of choice response time: Linear ballistic accumulation. *Cognitive psychology*, 57(3):153–178, 2008.
- R. Y. Cho, L. E. Nystrom, E. T. Brown, A. D. Jones, T. S. Braver, P. J. Holmes, and J. D. Cohen. Mechanisms underlying dependencies of performance on stimulus history in a two-alternative forced-choice task. *Cognitive, Affective, & Behavioral Neuroscience*, 2(4):283–299, 2002.
- F. R. Clarke, T. G. Birdsall, and W. P. Tanner Jr. Two types of roc curves and definitions of parameters. *The Journal of the Acoustical Society of America*, 31(5):629–630, 1959.
- A. Compte, N. Brunel, P. S. Goldman-Rakic, and X.-J. Wang. Synaptic mechanisms and network dynamics underlying spatial working memory in a cortical network model. *Cerebral Cortex*, 10(9):910–923, 2000.
- T. B. Crapse and M. A. Sommer. Frontal eye field neurons with spatial representations predicted by their subcortical input. *Journal of Neuroscience*, 29(16):5308–5318, 2009.
- G. Deco, E. T. Rolls, L. Albantakis, and R. Romo. Brain mechanisms for perceptual and reward-related decision-making. *Progress in Neurobiology*, 103:194–213, 2013.
- K. Desender, A. Boldt, T. Verguts, and T. H. Donner. Post-decisional sense of confidence shapes speed-accuracy tradeoff for subsequent choices. *bioRxiv*, page 466730, 2018a.
- K. Desender, P. R. Murphy, A. Boldt, T. Verguts, and N. Yeung. A post-decisional neural marker of confidence predicts information-seeking. *bioRxiv*, page 433276, 2018b.
- L. Ding and J. I. Gold. Neural correlates of perceptual decision making before, during, and after decision commitment in monkey frontal eye field. *Cerebral Cortex*, 22(5):1052–1067, 2011.
- J. Ditterich. Evidence for time-variant decision making. *European Journal of Neuroscience*, 24(12):3628–3641, 2006.
- J. Drugowitsch, R. Moreno-Bote, A. K. Churchland, M. N. Shadlen, and A. Pouget. The cost of accumulating evidence in perceptual decision making. *Journal of Neuroscience*, 32(11):3612–3628, 2012.
- J. Drugowitsch, R. Moreno-Bote, and A. Pouget. Relation between belief and performance in perceptual decision making. *PloS one*, 9(5):e96511, 2014.
- T. A. Engel and X.-J. Wang. Same or different? a neural circuit mechanism of similarity-based pattern match decision making. *Journal of Neuroscience*, 31(19):6982–6996, 2011.
- T. A. Engel, W. Chaisangmongkon, D. J. Freedman, and X.-J. Wang. Choice-correlated activity fluctuations underlie learning of neuronal category representation. *Nature communications*, 6:6454, 2015.
- M. O. Ernst and M. S. Banks. Humans integrate visual and haptic information in a statistically optimal fashion. *Nature*, 415(6870):429, 2002.
- M. P. Fay and M. A. Proschan. Wilcoxon-mann-whitney or t-test? on assumptions for hypothesis tests and multiple interpretations of decision rules. *Statistics surveys*, 4:1, 2010.
- J. H. Fecteau and D. P. Munoz. Exploring the consequences of the previous trial. *Nature Reviews Neuroscience*, 4(6):435, 2003.
- S. W. Fernberger. Interdependence of judgments within the series for the method of constant stimuli. *Journal of Experimental Psychology*, 3(2):126, 1920.
- J. Fischer and D. Whitney. Serial dependence in visual perception. *Nature neuroscience*, 17(5):738, 2014.
- S. M. Fleming, R. S. Weil, Z. Nagy, R. J. Dolan, and G. Rees. Relating introspective accuracy to individual differences in brain structure. *Science*, 329(5998):1541–1543, 2010.



- I. Fründ, F. A. Wichmann, and J. H. Macke. Quantifying the effect of intertrial dependence on perceptual decisions. *Journal of vision*, 14(7):9–9, 2014.
- S. J. Galvin, J. V. Podd, V. Drga, and J. Whitmore. Type 2 tasks in the theory of signal detectability: Discrimination between correct and incorrect decisions. *Psychonomic Bulletin & Review*, 10(4):843–876, 2003.
- E. S. Geller and C. P. Whitman. Confidence ill stimulus predictions and choice reaction time. *Memory & cognition*, 1(3): 361–368, 1973.
- A. Gelman and J. Hill. Data analysis using regression and hierarchical/multilevel models. *New York, NY: Cambridge*, 2007.
- C. M. Glaze, J. W. Kable, and J. I. Gold. Normative evidence accumulation in unpredictable environments. *Elife*, 4: e08825, 2015.
- J. I. Gold, C.-T. Law, P. Connolly, and S. Bennur. The relative influences of priors and sensory evidence on an oculomotor decision variable during perceptual learning. *Journal of neurophysiology*, 100(5):2653–2668, 2008.
- R. C. Gonzalez, R. E. Woods, et al. Digital image processing, 2002.
- D. Griffin and A. Tversky. The weighing of evidence and the determinants of confidence. *Cognitive psychology*, 24(3): 411–435, 1992.
- E. Grushka. Characterization of exponentially modified gaussian peaks in chromatography. *Analytical Chemistry*, 44(11): 1733–1738, 1972.
- M. N. Hebart, Y. Schriever, T. H. Donner, and J.-D. Haynes. The relationship between perceptual decision variables and confidence in the human brain. *Cerebral Cortex*, 26(1):118–130, 2014.
- E. L. Ionides, C. Breto, J. Park, R. A. Smith, and A. A. King. Monte carlo profile confidence intervals. *arXiv preprint arXiv:1612.02710*, 2016.
- J. Jaramillo, J. F. Mejias, and X.-J. Wang. Engagement of pulvino-cortical feedforward and feedback pathways in cognitive computations. *Neuron*, 101(2):321–336, 2019.
- A. Kepecs and Z. F. Mainen. A computational framework for the study of confidence in humans and animals. 367(1594): 1322–1337, 2012. ISSN 0962-8436. doi: 10.1098/rstb.2012.0037.
- A. Kepecs, N. Uchida, H. A. Zariwala, and Z. F. Mainen. Neural correlates, computation and behavioural impact of decision confidence. *Nature*, 455(7210):227, 2008. ISSN 1476-4687. doi: 10.1038/nature07200.
- R. Kiani and M. N. Shadlen. Representation of confidence associated with a decision by neurons in the parietal cortex. *Science*, 324(5928):759–764, 2009. ISSN 0036-8075. doi: 10.1126/science.1169405.
- R. Kiani, L. Corthell, and M. N. Shadlen. Choice certainty is informed by both evidence and decision time. *Neuron*, 84 (6):1329–1342, 2014.
- M. Kleiner, D. Brainard, D. Pelli, A. Ingling, R. Murray, C. Broussard, et al. What’s new in psychtoolbox-3. *Perception*, 36(14):1, 2007.
- Y. Komura, A. Nikkuni, N. Hirashima, T. Uetake, and A. Miyamoto. Responses of pulvinal neurons reflect a subject’s confidence in visual categorization. *Nature neuroscience*, 16(6):749, 2013.
- E. Kreyszig. Advanced engineering mathematics. fourth edi, 1979.
- A. Lak, G. M. Costa, E. Romberg, A. A. Koulakov, Z. F. Mainen, and A. Kepecs. Orbitofrontal cortex is required for optimal waiting based on decision confidence. *Neuron*, 84(1):190–201, 2014.
- D. Laming. Choice reaction performance following an error. *Acta Psychologica*, 43(3):199–224, 1979.
- D. A. Leopold, M. Wilke, A. Maier, and N. K. Logothetis. Stable perception of visually ambiguous patterns. *Nature neuroscience*, 5(6):605, 2002.
- A. Liberman, J. Fischer, and D. Whitney. Serial dependence in the perception of faces. *Current Biology*, 24(21):2569–2574, 2014.
- R. D. Luce et al. *Response times: Their role in inferring elementary mental organization*. Number 8. Oxford University Press on Demand, 1986.
- P. Mamassian. Visual confidence. *Annual Review of Vision Science*, 2(1):1–23, 2015. ISSN 2374-4642. doi: 10.1146/annurev-vision-111815-114630. URL <http://dx.doi.org/10.1146/annurev-vision-111815-114630>.
- S. Massoni. Emotion as a boost to metacognition: How worry enhances the quality of confidence. *Consciousness and cognition*, 29:189–198, 2014.

- MATLAB. *version R2016a*. The MathWorks Inc., Natick, Massachusetts, 2016.
- M. Mazurek, J. Roitman, J. Ditterich, and M. Shadlen. A role for neural integrators in perceptual decision making. *Cereb Cortex*, 13(11):1257–1269, 2003. ISSN 1047-3211. doi: 10.1093/cercor/bhg097.
- E. C. Merkle and T. Van Zandt. An application of the poisson race model to confidence calibration. *Journal of Experimental Psychology: General*, 135(3):391, 2006.
- F. Meyniel, D. Schlunegger, and S. Dehaene. The sense of confidence during probabilistic learning: A normative account. *PLoS computational biology*, 11(6):e1004305, 2015a.
- F. Meyniel, M. Sigman, and Z. F. Mainen. Confidence as bayesian probability: from neural origins to behavior. *Neuron*, 88(1):78–92, 2015b.
- P. Miller and D. B. Katz. Accuracy and response-time distributions for decision-making: linear perfect integrators versus nonlinear attractor-based neural circuits. *Journal of computational neuroscience*, 35(3):261–294, 2013.
- R. Moreno-Bote. Decision confidence and uncertainty in diffusion models with partially correlated neuronal integrators. *Neural computation*, 22(7):1786–1811, 2010.
- C. S. Peirce. On the theory of errors of observation. *Report of the Superintendent of the United States Coast Survey Showing the Progress of the Survey During the Year 1870*, pages 220–224, 1873.
- C. S. Peirce and J. Jastrow. On small differences in sensation. 1884.
- T. J. Pleskac and J. R. Busemeyer. Two-stage dynamic signal detection: a theory of choice, decision time, and confidence. *Psychological review*, 117(3):864, 2010.
- D. H. Raab. Division of psychology: Statistical facilitation of simple reaction times. *Transactions of the New York Academy of Sciences*, 24(5 Series II):574–590, 1962.
- R. Ratcliff. A theory of memory retrieval. *Psychological review*, 85(2):59, 1978.
- R. Ratcliff. Parameter variability and distributional assumptions in the diffusion model. *Psychological review*, 120(1): 281, 2013.
- R. Ratcliff and J. N. Rouder. Modeling response times for two-choice decisions. *Psychological Science*, 9(5):347–356, 1998.
- R. Ratcliff and J. J. Starns. Modeling confidence and response time in recognition memory. *Psychological review*, 116 (1):59, 2009.
- E. T. Rolls, F. Grabenhorst, and G. Deco. Choice, difficulty, and confidence in the brain. *Neuroimage*, 53(2):694–706, 2010a.
- E. T. Rolls, F. Grabenhorst, and G. Deco. Decision-making, errors, and confidence in the brain. *Journal of neurophysiology*, 104(5):2359–2374, 2010b.
- T. Rowan. *The subplex method for unconstrained optimization*. PhD thesis, Ph. D. thesis, Department of Computer Sciences, Univ. of Texas, 1990.
- RStudio Team. *RStudio: Integrated Development Environment for R*. RStudio, Inc., Boston, MA, 2015. URL <http://www.rstudio.com/>.
- J. Samaha, M. Switzky, and B. R. Postle. Confidence boosts serial dependence in orientation estimation. *bioRxiv*, page 369140, 2018.
- J. I. Sanders, B. Hangya, and A. Kepecs. Signatures of a statistical computation in the human sense of confidence. *Neuron*, 90(3):499–506, 2016.
- A. K. Seth. Post-decision wagering measures metacognitive content, not sensory consciousness. *Consciousness and cognition*, 17(3):981–983, 2008.
- J. D. Smith, W. E. Shields, and D. A. Washburn. The comparative psychology of uncertainty monitoring and metacognition. *Behavioral and brain sciences*, 26(3):317–339, 2003.
- P. L. Smith and D. Vickers. The accumulator model of two-choice discrimination. *Journal of Mathematical Psychology*, 32(2):135–168, 1988.
- M. A. Sommer and R. H. Wurtz. Visual perception and corollary discharge. *Perception*, 37(3):408–418, 2008.
- A. E. Urai, A. Braun, and T. H. Donner. Pupil-linked arousal is driven by decision uncertainty and alters serial choice bias. *Nature communications*, 8:14637, 2017.
- M. Usher and J. L. McClelland. The time course of perceptual choice: the leaky, competing accumulator model. *Psychological review*, 108(3):550, 2001.



- S. Verdonck and F. Tuerlinckx. Factoring out nondecision time in choice reaction time data: Theory and implications. *Psychological review*, 123(2):208, 2016.
- D. Vickers. Evidence for an accumulator model of psychophysical discrimination. *Ergonomics*, 13(1):37–58, 1970.
- D. Vickers. *Decision processes in visual perception*. Academic Press, 1979 (reedited in 2014).
- D. Vickers and J. Packer. Effects of alternating set for speed or accuracy on response time, accuracy and confidence in a unidimensional discrimination task. *Acta psychologica*, 50(2):179–197, 1982.
- X.-J. Wang. Probabilistic decision making by slow reverberation in cortical circuits. *Neuron*, 36(5):955–968, 2002.
- X.-J. Wang. Decision Making in Recurrent Neuronal Circuits. *Neuron*, 60:215–234, 2008. ISSN 0896-6273. doi: 10.1016/j.neuron.2008.09.034.
- Z. Wei and X.-J. Wang. Confidence estimation as a stochastic process in a neurodynamical system of decision making. *Journal of neurophysiology*, 114(1):99–113, 2015.
- F. Wilcoxon. Individual comparisons by ranking methods. 1(6):80, 1945. ISSN 0099-4987. doi: 10.2307/3001968.
- K.-F. Wong and X.-J. Wang. A recurrent network mechanism of time integration in perceptual decisions. *Journal of Neuroscience*, 26(4):1314–1328, 2006.
- K.-F. Wong, A. C. Huk, M. N. Shadlen, and X.-J. Wang. Neural circuit dynamics underlying accumulation of time-varying evidence during perceptual decision making. *Frontiers in Computational Neuroscience*, 1:6, 2007.
- N. Yeung and C. Summerfield. Metacognition in human decision-making: confidence and error monitoring. *Phil. Trans. R. Soc. B*, 367(1594):1310–1321, 2012.
- A. Zylberberg, P. Bartfeld, and M. Sigman. The construction of confidence in a perceptual decision. *Frontiers in integrative neuroscience*, 6:79, 2012.

## A ESTIMATION OF THE NON-DECISION TIME DISTRIBUTION

We consider that the nondecision time distribution, noted  $\rho_{NDT}$ , is described by an exponentially modified Gaussian (EMG) distribution [Verdonck and Tuerlinckx, 2016]:

$$\rho_{NDT}(t) = \frac{\lambda_{NDT}}{2} \exp\left(\frac{\lambda_{NDT}}{2}(2\mu + \lambda_{NDT} \sigma_{NDT}^2 - 2t)\right) \operatorname{erfc}\left(\frac{\mu + \lambda_{NDT} \sigma_{NDT}^2 - t}{\sqrt{2}\sigma_{NDT}}\right) \quad (10)$$

with  $\operatorname{erfc}$  the complementary error function. For such distribution, the mean non decision time,  $\langle NDT \rangle$ , is given by

$$\langle NDT \rangle = \mu_{NDT} + \frac{1}{\lambda_{NDT}}. \quad (11)$$

As we assume no correlation between response and non decision times, the total (observed) response time distribution,  $\rho_{data}$ , can be written as the convolution of the decision time distribution,  $\rho_{decision}$ , with the nondecision time distribution,  $\rho_{NDT}$ :

$$\rho_{data}(t) = \rho_{decision}(t) * \rho_{NDT}(t) = \int_0^t \rho_{decision}(t-u) \rho_{NDT}(u) du \quad (12)$$

(with  $*$  standing for the convolution operation). If the decision time distribution is Gaussian, the resulting total distribution is an EMG distribution [Grushka, 1972].

Using maximum likelihood estimation, for each subject we fit the empirical response time distribution (all orientations together) by an EMG distribution with parameters  $\mu_{data}, \lambda_{data}, \sigma_{data}$  (and we thus have  $\langle RT \rangle_{data} = \mu_{data} + 1/\lambda_{data}$ ). For what concerns the model, we find that the decision time distribution of the attractor neuronal network is well fitted by a Gaussian distribution with parameters  $\langle RT \rangle_{network}, \sigma_{network}$ . We thus model the decision time distribution by the one provided by the attractor network whose parameters have been calibrated as explained above. Hence, we identify the mean and variance of the decision time distribution with the ones of the network:  $\mu_{decision} = \langle RT \rangle_{network}, \sigma_{decision} = \sigma_{network}$ .

Taking the characteristic function of both sides of Equation 12, we get:

$$\left(1 - \frac{it}{\lambda_{data}}\right)^{-1} \exp\left(i\mu_{data}t - \frac{1}{2}\sigma_{data}^2t^2\right) = \exp\left(i\mu_{decision}t - \frac{1}{2}\sigma_{decision}^2t^2\right) \left(1 - \frac{it}{\lambda_{NDT}}\right)^{-1} \exp\left(i\mu_{NDT}t - \frac{1}{2}\sigma_{NDT}^2t^2\right)$$

We can then identify the terms on both sides of this equations, which, given the use for the decision parameters of the ones of the attractor network, gives the equations

$$\lambda_{NDT} = \lambda_{data}, \tag{13}$$

$$\langle NDT \rangle = \mu_{NDT} + \frac{1}{\lambda_{NDT}} = \langle RT \rangle_{data} - \langle RT \rangle_{network} \tag{14}$$

and

$$\sigma_{NDT}^2 = \sigma_{data}^2 - \sigma_{network}^2. \tag{15}$$

from which we can compute the non-decision time distribution parameters,  $\lambda_{NDT}, \mu_{NDT}, \sigma_{NDT}$ .

## B EXPECTED SEQUENTIAL EFFECT IN AN IRM

In this Appendix we show that, without any change of parameters, the independent race model (IRM) cannot account for the observed sequential effects.

First, we present the equations governing the independent race models (IRM) used in Fig. 9. During the accumulation of evidence the equations of evolution are:

$$dx_i = I_0(1 \pm c)dt + \sigma\nu_i(t) \tag{16}$$

where  $\nu_i(t), i = \{C, AC\}$ , are white noise processes. The first race that reaches a threshold  $z$  (or  $-z$ ) is the winning race. The confidence in the decision is modelled as a monotonic function of the balance of evidence  $|z - x_{losing}|$  [Vickers, 1979 (reedited in 2014, Drugowitsch et al., 2014, Mamassian, 2015, Wei and Wang, 2015)].

We extend the IRM in order to deal with sequences of trials. To do so, we allow for a relaxation dynamics between trials, in a way analogous to the relaxation dynamics in the attractor network model. Hence, after a decision is made, both units receive a non specific inhibitory input leading to a relaxation until the next stimulus is presented (see Appendix B - Figure 9). Within this extended IRM framework, we study how the sequential effects would be correlated with confidence in an IRM model with a fixed set of parameters.

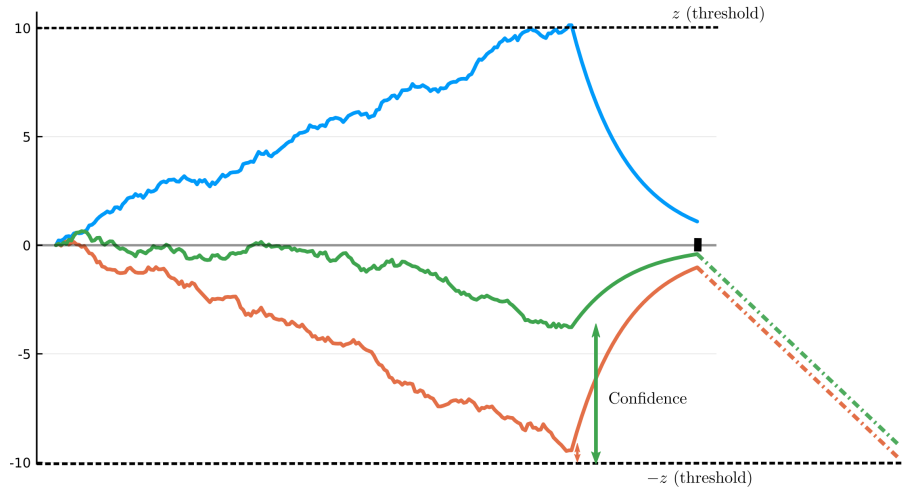


Figure 9: **Schematic dynamics of a race model with a relaxation mechanism.** The upper and bottom dash lines correspond to the two opposite decision thresholds. The blue trajectory is a typical winning race. The black rectangle on the x-axis denotes the onset of the next stimulus, hence the end of the relaxation period. The green and orange trajectories are the losing races in two trials with different confidence outcomes. The green and orange dashed lines represent the

*mean dynamics of these two races during the presentation of the next stimulus.*

Since in the IRM there is no interaction between the two races, the relaxation of the winning race is the same in both low and high confidence trials. However, the ending point of the relaxation following a decision is closer to the base-line (0 line) for a high confidence trial than when it comes to a trial with low confidence trial (Figure Appendix B - Figure 9). For the next trial, if the winning race is the same as previously, then the mean response times are identical in low and high confidence cases. However, if the opposite decision is made, the response time in the post-low confidence case is faster than the one in the post-high confidence case, as we can observe with the mean race shown in Appendix B - Figure 9. This behavior is in contradiction with the experimental data for which we observe the opposite effect (see Material and Methods Section, Table 2).

This conclusion applies more generally to any race-type model without interactions between units. Race models with interactions have also been considered [Bogacz et al., 2006], but such models are thus more similar to attractor networks, yet with less biophysical foundation.