1 **Separable neural signatures of confidence during perceptual decisions**

2 Balsdon, T.[1]*, Mamassian, P.[1]**, and Wyart, V.[2]**

3     1. Laboratoire des Systèmes Perceptifs (CNRS UMR 8248), DEC, ENS, PSL University, 75005 Paris,
4        France

5     2. Laboratoire de Neurosciences Cognitives et Computationelles (Inserm U960), DEC, ENS, PSL
6        University, 75005, Paris, France

7 * Corresponding author

8 ** Equal contributors

## 9 Abstract

10 Perceptual confidence is an evaluation of the validity of perceptual decisions. While there is behavioural
11 evidence that confidence evaluation differs from perceptual decision-making, disentangling these two
12 processes remains a challenge at the neural level. Here we examined the electrical brain activity of human
13 participants in a protracted perceptual decision-making task where observers tend to commit to perceptual
14 decisions early whilst continuing to monitor sensory evidence for evaluating confidence. Premature decision
15 commitments were revealed by patterns of spectral power overlying motor cortex, followed by an
16 attenuation of the neural representation of perceptual decision evidence. A distinct neural representation
17 was associated with suboptimalities affecting confidence reports, with sources localised in the superior
18 parietal and orbitofrontal cortices. In agreement with a dissociation between perception and confidence,
19 these neural resources were recruited even after observers committed to their perceptual decisions, and
20 thus delineate an integral neural circuit for the computation of confidence. [148 words]

## 21 Introduction

22 Whilst perception typically feels effortless and automatic, it requires probabilistic inference to resolve the
23 uncertain causes of essentially ambiguous sensory input (Helmholtz, 1856). Human observers are capable of
24 discriminating which perceptual decisions are more likely to be correct using subjective feelings of
25 confidence (Pollack and Decker, 1958). These feelings of perceptual confidence have been associated with
26 metacognitive processes (Fleming and Daw, 2017) that enable self-monitoring for learning (Veenman,
27 Wilhelm, & Beishuizen, 2004) and communication (Bahrami et al., 2012; Frith, 2012). We are only just
28 beginning to uncover the complex functional role of metacognition in human behaviour, and outline the
29 computational and neural processes that enable metacognition. The study of perceptual confidence offers
30 promising insight into metacognition, because one can use our detailed knowledge of perceptual processes
31 to isolate factors which affect the computation of perceptual confidence.

32 At the computational level, perceptual decisions are described by sequential sampling processes (Vickers,
33 1970; Ratcliff, 1978; Pleskac and Busemeyer, 2010), in which noisy samples of evidence are accumulated

34    over time, until there is sufficient evidence to commit to a decision. Perceptual confidence tends to reflect

35    the quantity and quality of evidence used to make the perceptual decision (Vickers, 1979; Kepecs et al.,

36    2008; Moreno-Bote, 2010). In this way, perceptual confidence is necessarily tethered to decision evidence:

37    more evidence for the perceptual decision yields greater perceptual accuracy, and therefore higher

38    confidence. This makes it difficult to dissociate what processes could be specifically involved in the

39    computation of confidence beyond the underlying perceptual processes. Indeed, confidence (or a non-

40    human primate proxy for confidence) can be reliably predicted from the firing rates of neurons coding the

41    perceptual decision itself (Kiani and Shadlen, 2009), suggesting that confidence may be a direct by-product

42    of perceptual processing.

43    However, a large body of behavioural studies suggest that confidence is affected by additional sources of

44    noise that do not influence perceptual decisions (Bang, Shekhar and Rahnev, 2019; Shekhar and Rahnev,

45    2020). And conversely, the precision of perceptual confidence can be boosted by integrating additional

46    information, such as decision time (Kiani, Corthell, and Shadlen, 2014) or continued evidence accumulation

47    after the observer commits to a perceptual decision (Baranski and Petrusic, 1994; Pleskac and Busemeyer,

48    2010). Together these factors mean that the same perceptual decision evidence can lead to different levels of

49    confidence, explaining the diverse range of confidence precision displayed by human observers, and

50    suggesting essential differences in the processes for perceptual and confidence decisions. Moreover,

51    evidence suggesting that confidence precision is correlated across different tasks (such as memory and

52    perception; Mazancieux et al., 2018) further calls into question whether confidence is a mere consequence of

53    perceptual processes, or rather, recruits specialised metacognitive resources that operate across cognition,

54    incurring similar suboptimalities in evaluating any cognitive process.

55    In this experiment we aimed to delineate the neural processes contributing to perception and confidence,

56    using electroencephalography (EEG). We exploited a protracted decision-making task in which the evidence

57    presented to the observer can be carefully controlled. On each trial, the observer was presented with a

58    sequence of visual stimuli, oriented Gabor patches, which offer a specific amount of evidence towards the

59    perceptual decision. The orientations were sampled from one of two overlapping circular Gaussian

60    distributions, and the observer was asked to categorise which distribution the orientations were sampled

61    from. We manipulated the amount of evidence presented such that the observer tends to covertly commit to

62    their perceptual decision before evidence presentation has finished, whilst continuing to monitor ongoing

63    evidence for assessing their confidence (Balsdon et al., 2020). These covert decisions were evident from

64    behaviour and computational modelling, and we show similarities between the neural processes of decision-

65    making across conditions of immediate and delayed response execution.

66    Human behaviour was compared to an optimal observer who perfectly accumulates all the presented

67    evidence for perceptual decisions and confidence evaluation. The optimal observer must accurately encode

68    the stimulus orientation, the decision update relevant for the categorisation, and add this to the accumulated

69    evidence for making the perceptual decision. We uncovered dynamic neural representations of these

70    variables, and examined how the precision of these representations fluctuate with behavioural

71  suboptimalities. We found two distinct representations of the accumulated evidence where imprecision in

72  the representation was related to suboptimal behaviour in the perceptual decisions and confidence

73  evaluations respectively. The noise contributing to the imprecision of the confidence representation was

74  localised to the Superior Parietal and Orbitofrontal cortices. Whilst the perceptual representation was

75  attenuated following covert decisions, the confidence representation continued to reflect evidence

76  accumulation. This is consistent with a neural circuit that can be recruited for confidence evaluation

77  independently of perceptual processes, providing empirical evidence for the theoretical dissociation

78  between perception and confidence.

# Results

## The computational architecture of perceptual confidence

81  Human observers (N = 20) performed two versions of the task whilst EEG was recorded. Across the two

82  tasks, 100 predefined sequences of oriented Gabors were repeated for each observer, with stimuli presented

83  as described in **Figure 1a**. In the Free task, the sequence continued until observers entered their perceptual

84  decision (**Figure 1b**), indicating which category (**Figure 1d**) they thought the orientations were sampled

85  from. Observers were instructed to enter their response as soon as they 'felt ready', on three repeats of each

86  predefined sequence (300 trials in total). In the Replay task (**Figure 1c**), observers were shown a specific

87  number of samples and could only enter their response after the response cue. After entering their

88  perceptual (Type-I) decision, they made a confidence (Type-II) evaluation, how confident they were that

89  their perceptual decision was correct, on a 4-point scale. Importantly, the number of samples shown in the

90  Replay task was manipulated relative to the Free task, in three intermixed conditions: in the Less condition,

91  they were shown two fewer than the minimum they had chosen to respond to over the three repeats of that

92  predefined sequence in the Free task; in the Same condition they were shown the median number of

93  samples; and in the More condition, four more than the maximum (**Figure 1e**). The variability across repeats

94  in the Free task means that in the More condition, observers were show at least four additional stimuli, but
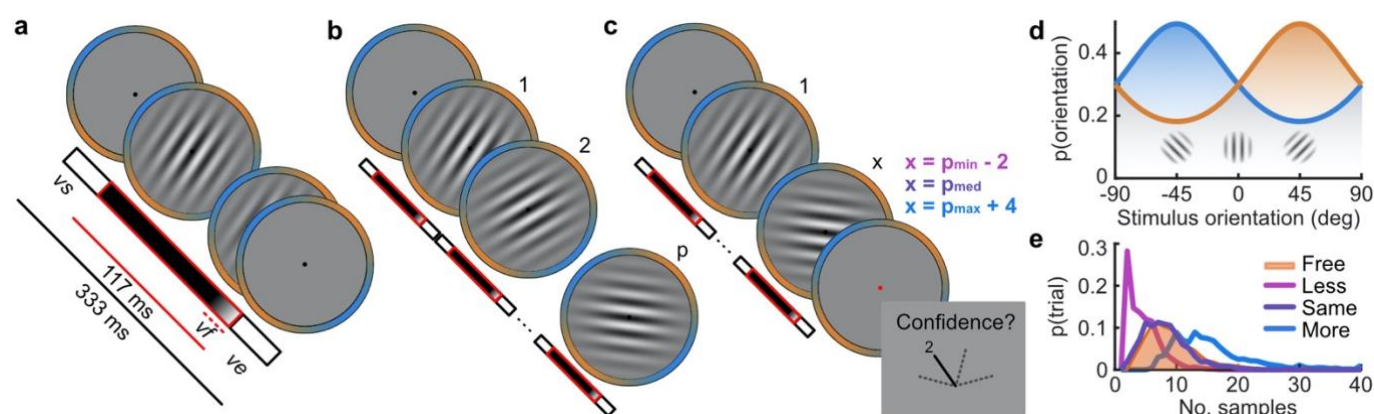
95  often more than that.

96



97  ***Figure 1. Procedure. a)*** *Stimulus presentation: stimuli were presented at an average rate of 3 Hz, but with*

98  *variable onset and offset (* $vs \in [83, 133]$ *ms,* $vs_s + ve_{s-1} \geq 216$ *ms; see **Methods**). Stimuli were presented*

3

99 *within a circular annulus which acted as a colour guide for the category distributions. The colour guide and the*

100 *fixation point were present throughout the trial.* **b)** *Free task: on each trial observers were presented with a*

101 *sequence of oriented Gabors, which continued until the observer entered their response (or 40 samples were*

102 *shown). 100 sequences were predefined and repeated three times.* **c)** *Replay task: The observer was presented*

103 *with a specific number of samples and could only enter their response after the cue (fixation changing to red).*

104 *The number of samples (x) was determined relative to the number the observer chose to respond to on that*

105 *same sequence in the Free task (p). There were three intermixed conditions, Less (x = $p_{min}$ – 2; where $p_{min}$ is the*

106 *minimum p of the three repeats), Same (x = $p_{med}$; where $p_{med}$ is the median p) and More (x = $p_{max}$ + 4; where $p_{max}$*

107 *is the maximum p of the three repeats of that predefined sequence).* **d)** *Categories were defined by circular*

108 *Gaussian distributions over the orientations, with means -45° (blue) and 45° (orange), and concentration* $\kappa$ =

109 **0.5**. *The distributions overlapped such that an orientation of 45° was most likely drawn from the orange*

110 *distribution but could also be drawn from the blue distribution with lower likelihood.* **e)** *Distributions of the*

111 *number of samples per trial in the Free task, and Replay task conditions (over all observers).*

112 Based on previous findings (Balsdon et al., 2020) we expected observers to prematurely commit to

113 perceptual decisions in the More condition, whilst continuing to monitor sensory evidence for evaluating

114 their confidence. Replicating previous results (Balsdon et al., 2020), we found that perceptual decision

115 sensitivity was significantly decreased in the Less condition compared to those same ($p_{min}$) trials in the Free

116 task (Wilcoxon sign rank $Z$ = 3.88, $p$ < 0.001, Bonferroni corrected for three comparisons), there was no

117 significant difference for the Same condition ($Z$ = 1.21, $p$ = 0.23, uncorrected), nor the More condition ($Z$ = -

118 1.53, $p$ = 0.13, uncorrected; despite at least an additional four samples being presented compared to the $p_{max}$

119 trials; **Figure 2a**). In addition, reaction times in the More condition were significantly decreased compared

120 to the Same condition (on average, 60 ms faster; $Z(19)$ = 2.58, $p$ = 0.010; **Figure 2b**).

121 This lack of substantial increase in performance in the More condition could be the result of either a

122 performance ceiling effect or a premature commitment to the perceptual decision. The former explanation

123 reflects a limitation of the perceptual evidence accumulation process, whereas the latter refers to an active

124 mechanism that ignores the final sensory evidence. We compared these two hypotheses using a

125 computational modelling approach (Balsdon et al., 2020; see **Methods**). Specifically, we compared a model

126 in which performance in the More condition is limited by the suboptimalities evident from the Same and the

127 Less conditions (inference noise, and temporal integration bias, see **Supplementary Note 1**), to a model in

128 which performance could be impacted by a covert bound at which point observers commit to a decision

129 irrespective of additional evidence. Cross-validated model comparison provided significant evidence that

130 observers were implementing a covert bound (mean relative increase in model log-likelihood = 0.048,

131 bootstrapped $p$ = 0.001, **Figure 2c**). The winning model provided a good description of the data (red open

132 markers in **Figure 2a**).

133 In contrast to what we found for the perceptual decision, there was no evidence that observers were

134 implementing a covert bound on confidence: Implementing the same bound as the perceptual decision did

135 not improve the fit (relative improvement with bound = -0.007, bootstrapped $p$ = 0.11, uncorrected) and an

4

136   independent bound actually significantly *reduced* the fit compared to continued accumulation (relative

137   improvement = -0.014, *p* = 0.022, Bonferroni corrected for two comparisons; **Figure 2c**). We obtained

138   further distinctions between perceptual and confidence processes through computational modelling:

139   additional noise was required to explain the confidence ratings, along with a separate temporal bias. The

140   best description of both perceptual and confidence responses was provided by a partially dissociated

141   computational architecture (full details in **Supplementary Note 1**), where perceptual and confidence

142   decisions are based on the same noisy representation of the sensory evidence, but confidence accumulation

143   incurs additional noise and can continue after the completion of perceptual decision processes (**Figure 2d**).

144   These computational differences between perceptual decisions and confidence evaluations result in

145   deviations between the inference errors associated with perceptual and confidence decisions (see

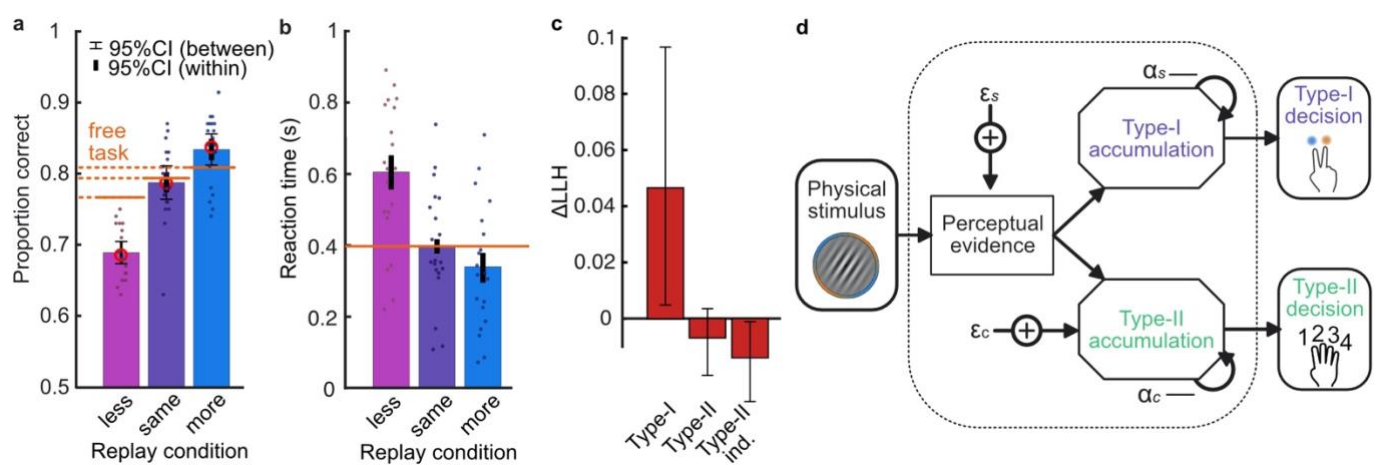146   **Supplementary Note 2** for model simulations).



147

148   *Figure 2. Behaviour and computational modelling. a) Proportion correct in each condition of the Replay*

149   *task, relative to the Free task (orange horizontal lines). Individual data are shown in scattered points, error*

150   *bars show 95% between- (thin) and 95% within- (thick) subject confidence intervals. Open red markers show*

151   *the model prediction. b) Median reaction time (from response cue) in each condition of the Replay task, error*

152   *bars show 95% within-subject confidence intervals, the orange horizontal line shows the non-decision time*

153   *estimated in the Free task, based on computational modelling. Individual data are shown in scattered points. c)*

154   *Difference in log-likelihood of the models utilising a covert bound relative to the models with no covert bound.*

155   *On the left, the model fitting perceptual decisions only. The middle bar shows the difference in log-likelihood of*

156   *the fit to confidence ratings with identical perceptual and confidence bounds. The right bar shows the difference*

157   *in log-likelihood of the fit to confidence ratings of the model with an independent bound for confidence evidence*

158   *accumulation. Error bars show 95% between-subject confidence intervals. d) Computational architecture of*

159   *perceptual and confidence decisions. Perceptual (Type-I) and confidence (Type-II) decisions accumulate the*

160   *same noisy perceptual evidence, but confidence is affected by additional noise ($\varepsilon_c$) and a separate temporal bias*

161   *($\alpha_c$). This partial dissociation allows Type-II accumulation to continue after the observer has committed to a*

162   *perceptual decision.*

5

## Model-free EEG analysis

We first examined EEG amplitude modulations around the time of the response: the CPP (Central-Parietal Positivity; O'Connell et al., 2012) and the LRP (Lateralised Readiness Potential; Deecke et al., 1976). There were significant differences in the CPP and the LRP between the More and the Less conditions of the Replay task (for the CPP, from -500 prior to the response, the largest cluster showing $t_{ave}(19) = -2.85$, $p_{cluster} = 0.006$; for the LRP, from just after the response, the first cluster from 32 to 196 ms; $t_{ave}(19) = -3.57$, $p_{cluster} < 0.002$; **Figure 3**, left). There were also differences based on perceptual decision accuracy (for CPP, the main cluster emerges from -156 ms to 592 ms around the response; $t_{ave}(19) = 4.38$, $p_{cluster} < 0.002$; and LRP from -84 ms to 652 ms around the response, with the largest difference just after the response, $t_{ave}(19) = 2.81$, $p_{cluster} < 0.002$; **Figure 3**, middle). There was no significant difference in the LRP between trials with high confidence (ratings of 3 and 4) and low confidence (ratings of 1 and 2), but a substantial difference was observed in the CPP (from 250 ms prior to the response; $t_{ave}(19) = 4.46$, $p_{cluster} < 0.002$; **Figure 3**, right), in line with previous findings (e.g. Herding et al., 2019). These modulations are consistent with the differences in the underlying accumulated evidence driving observers' responses. We aimed to more closely examine the neural processes underlying these broad effects on EEG amplitude, especially with respect to the distinctions between perceptual decision-making and confidence evaluation, as identified by the computational model of behaviour: perceptual decision processes can conclude prior to the confidence evaluation processes, and rely on a representation of the evidence that incurs distinct inference errors.
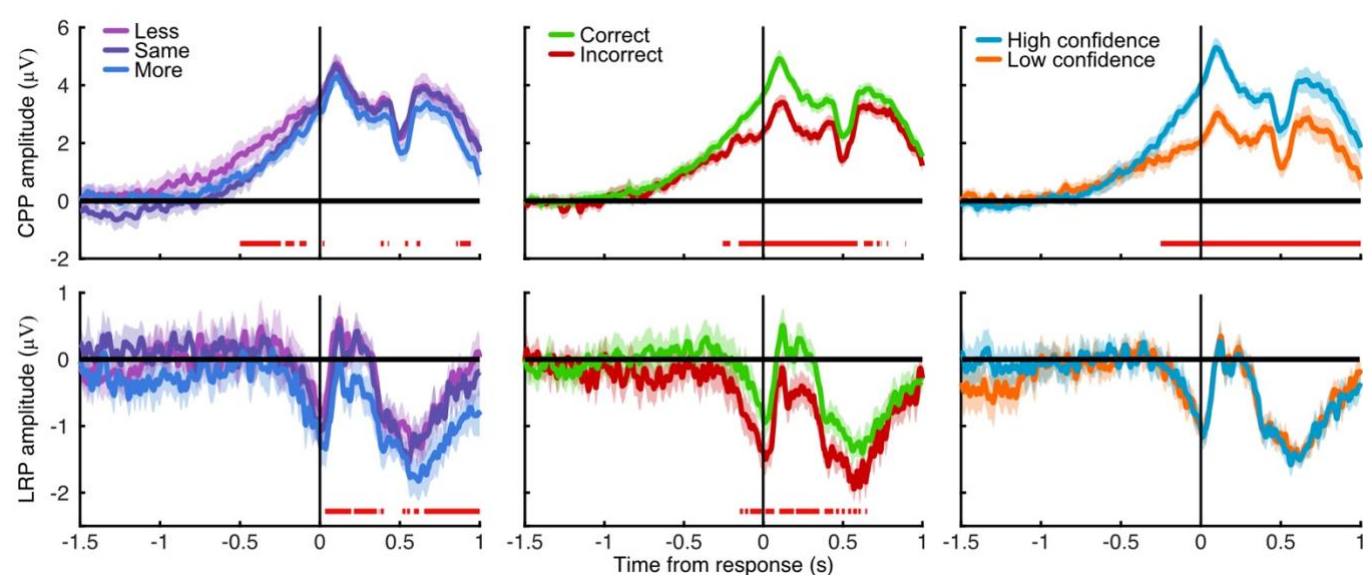


*Figure 3. Amplitude modulations with task variables.* *Central-Parietal Positivity (CPP, top) and Lateralised Readiness Potential (LRP, bottom), within condition (Less/Same/More; left), correct and incorrect perceptual responses (middle), and high and low confidence ratings (right). Vertical black lines mark the time of the response, red horizontal lines mark cluster corrected significant differences. Shaded regions show 95% within-subject confidence intervals.*

**EEG signatures of premature perceptual decision commitment**

We examined the neural signatures of perceptual decision commitment using a linear discriminant analysis of the spectral power of band-limited EEG oscillations (see **Methods**). A classifier wad trained to discriminate observers' perceptual decisions based on the spectral power in 8 to 32 Hz frequency bands at time-points leading up to the response in the Free task (**Supplementary Note 4**). This classifier was then tested across time in each condition of the Replay task, to trace the progression of perceptual decision-making in comparison to the Free task (where decisions are directly followed by response execution). There were opposite asymmetries in the cross-classification of the Less and the More conditions (**Figure 4a**). Statistical comparison revealed substantial clusters of significant differences (**Figure 4b**): Training around -0.78 to 0.44 s from the time of the response in the Free task led to significantly better accuracy testing in the More condition than in the Less condition, prior to when the response was entered (for the cluster testing at -2.5 to -1.6 s $Z_{ave}$ = 2.04, $p_{cluster}$ = 0.002; testing at -1.5 to -1 s, $Z_{ave}$ = 1.95, $p_{cluster}$ = 0.01; testing at -0.8 to -0.3, $Z_{ave}$ = 2.32, $p_{cluster}$ < 0.001). This pattern of findings suggests that observers were not only committing to their perceptual decision early, but already preparing their motor response, which would explain the faster reaction times in the More condition (**Figure 2b**).

We found that the accumulated evidence over all samples could predict the strength of the neural signature of response execution (mean $\beta$ = 0.11, $t(19)$ = 3.89, $p$ < 0.001; **Figure 4c**). For the Same and Less conditions, the weight on the accumulated evidence appeared to decrease as evidence was accumulated to samples further prior from the response. But, in the More condition, the evidence accumulated up to four samples prior to the response still predicted the classifier response ($t(19)$ = 3.81, $p$ = 0.001). This difference between conditions over samples is evidenced by a significant interaction based on a repeated measures ANOVA ($F(8,152)$ = 2.429, $p$ = 0.05, after Bonferroni correction for three comparisons). Leading up to the response, the accumulated evidence becomes increasingly predictive of the strength of the neural signature of response execution, except in the More condition, where this prediction is already accurate up to four samples prior to the response: After committing to a perceptual decision, the observer's perceptual response is no longer influenced by additional evidence.
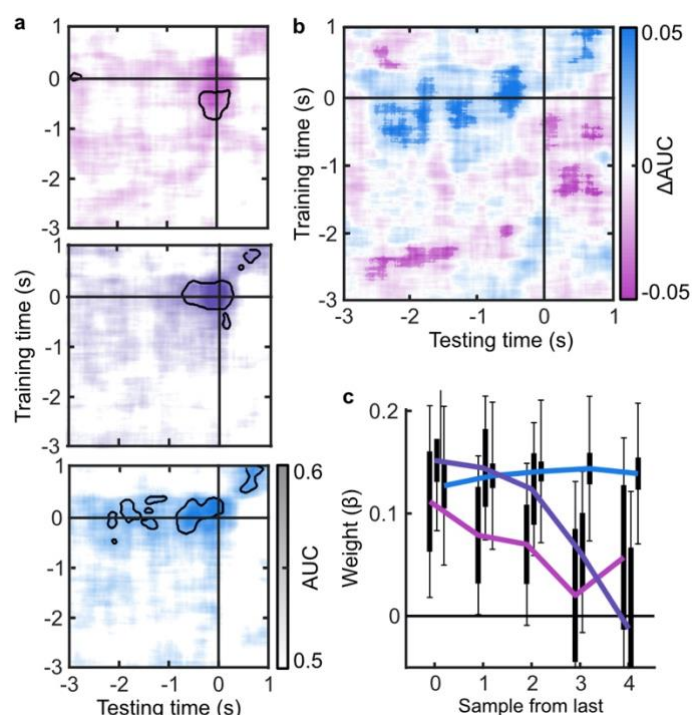
7

*Figure 4. EEG signatures of premature perceptual decisions. a) Classifier AUC training at each time-point in the Free task and testing across time in the Less (top), Same (middle), and More (bottom) conditions of the Replay task. Black contours encircle regions where the mean is 3.1 standard deviations from chance (0.5; 99% confidence). b) Difference in AUC between the More and Less conditions. Cluster corrected significant differences are highlighted. c) The relationship between the evidence accumulated up to n samples prior to the response cue and the strength of the neural signature of response execution in each condition. Error bars show 95% within- (thick) and between-subject (thin) confidence intervals.*

## Representations of decision evidence in EEG signals

To perform this task the optimal observer must encode the orientation of the stimulus, estimate the decision update based on the categories, and add this to the accumulated evidence for discriminating between the categories (Wyart et al., 2012; Wyart et al., 2015). We examined the neural representation of these optimal variables using a regression analysis with the EEG signals (evoked response, bandpass filtered between 1 and 8 Hz, see **Methods**). **Figure 5a** shows the time course of the neural coding of stimulus orientation, momentary decision update, and accumulated evidence, locked to stimulus onset. The representations of these variables showed distinct time courses and relied on distinct patterns of EEG activity over scalp topography (**Figure 5b**). There was a transient representation of stimulus orientation localised over occipital electrodes. The representation of the momentary decision update was maintained for a longer duration, initially supported by occipital electrodes, then increasingly localised over central-parietal electrodes. The representation of the accumulated evidence was sustained even longer and relied on both frontal and occipital electrodes.

The precision with which the EEG representations reflect the optimal decision variables can be compared with observers' suboptimal inference, based on whether the observers' behavioural responses matched those of an optimal observer. For each variable, we estimated the representation precision separately for epochs leading to suboptimal behavioural responses, and responses that matched those of the optimal observer (Replay task epochs only; **Figure 5c**; **Supplementary Note 3**). For perceptual decisions, the optimal observer responds with the correct category. For confidence evaluations, the optimal observer gives high confidence on trials with greater than the median evidence (over all trials) for their perceptual response. The precision of the representation of stimulus orientation did not significantly vary with behavioural suboptimalities. The representation precision of the momentary decision update showed a significant effect of perceptual decision suboptimality from 380 to 468 ms ($F_{avg}(1,19) = 7.97$, $p_{cluster} = 0.008$)

8

237  and a significant interaction between perceptual and confidence suboptimalities from 396 to 468 ms

238  ($F_{avg}(1,19) = 6.66$, $p_{cluster} = 0.022$) and from 716 to 856 ms ($F_{avg}(1,19) = 10.75$, $p_{cluster} < 0.001$). The largest

239  effects were seen in the representation precision of the accumulated evidence. Representation precision was

240  significantly reduced in epochs leading to suboptimal perceptual decisions from 108 ms post stimulus onset

241  to the end of the epoch ($F_{avg}(1,19) = 13.65$, $p_{cluster} <0.001$). In addition, there was a significant interaction with

242  suboptimal confidence from 696 to 836 ms ($F_{avg}(1,19) = 8.72$, $p_{cluster} = 0.005$). The precision of the EEG

243  representations showed distinct associations with the suboptimality of behavioural responses.

244  The presence of a covert bound implies that, after the observer commits to a decision, they no longer

245  incorporate additional evidence for that decision. We should therefore see significant decreases in the

246  precision of representations that specifically contribute to perceptual evidence accumulation. Indeed, the

247  precision of the early representation of accumulated evidence was significantly attenuated for the last four

248  samples of the More condition (in which observers were likely to have already committed to a decision),

249  compared to the last four samples of the Less condition (where observers were unlikely to have committed

250  to a decision; from the start of the epoch to 424 ms, **Figure 5d**; $t_{avg}(19) = -5.19$, $p_{cluster}<0.001$). These

251  differences in representation precision were not present for the encoding of stimulus orientation, nor the

252  decision update, nor was the decreased precision evident in a comparison of the first four samples

253  (**Supplementary Note 5**). Together, these comparisons suggest that different aspects of these evolving EEG

254  representations of decision variables are related to the neural processes for perception and confidence.



255

256  *Figure 5. Representation of decision variables. a) Representation precision (Fischer transformed correlation*

257  *coefficient, z) of stimulus orientation (blue, left), momentary decision update (green, middle), and accumulated*

258  *decision evidence (purple, right). The encoded variables are shown in the insets (the accumulated evidence is*

259  *the cumulative sum of the momentary evidence signed by the response, only one example sequence is shown).*

260  *Shaded regions show 95% between-subject confidence intervals. b) Relative electrode representation precision*

9

261 *over three characteristic time windows (100 – 200 ms, left; 400 – 600 ms, middle; and 600 – 800 ms, right).* **c)**

262 *Representation precision for epochs leading to optimal and suboptimal perceptual (T1) and confidence (T2)*

263 *responses. Lighter lines show perceptual decisions that match the optimal response, dashed lines show*

264 *suboptimal confidence ratings. Dashed red horizontal lines show significant interactions between perceptual*

265 *and confidence suboptimality. The light red horizontal line shows the significant effect of suboptimal perception*

266 *and the dark red horizontal line shows the significant effect of suboptimal confidence. Shaded regions show*

267 *95% within-subject confidence intervals.* **d)** *Difference in decoding precision between the More and the Less*

268 *conditions for epochs corresponding to the last four samples of the trial. The purple horizontal line shows the*

269 *significant difference in decoding of accumulated evidence.*

270 **Neural processes for confidence**

271 The analysis above shows that at certain times there was on average more noise affecting the EEG

272 representation of accumulated evidence on epochs leading to suboptimal behavioural responses. We

273 examined whether this increase in noise was due to a systematic change in the representation that could be

274 functionally related to the inference suboptimalities affecting observers' decision-making and confidence

275 evaluation. Cluster modelling with multivariate Bayesian scan statistics (Neill, 2011; Neill, 2019) was used to

276 isolate contiguous signals in space (electrode location) and time where imprecision in the representation of

277 accumulated evidence was associated with behavioural suboptimalities beyond what could be explained by

278 deviations in measurement noise alone (see **Supplementary Note 6** for further details). For perceptual

279 decision-making, signals were initially clustered over posterior electrodes, becoming dispersed over more

280 anterior electrodes late in the epoch (**Figure 6a**, top). For confidence, we found two co-temporal clusters in

281 posterior and anterior electrodes emerging from 668 ms to 824 ms from stimulus onset (**Figure 6a**,

282 bottom). We combined the signals from the two confidence clusters to estimate the confidence

283 representation of accumulated evidence (**Figure 6b**, dark green bar). We used this representation to

284 estimate the single-sample inference error of the observer, based on the deviation of the effective (noisy)

285 value from the predicted value, given the representation and the true value presented to the observer.

286 We compared the inference error estimated from the confidence representation to the inference error

287 estimated from the computational model of behaviour. There was a significant correlation with the error

288 estimated from the model of confidence ratings (mean $z = 0.05$, $t(19) = 5.12$, $p < 0.001$), and this correlation

289 was significantly greater than the error estimated from the model of perceptual decisions alone ($t(19) =$

290 $2.62$, $p = 0.017$; see **Supplementary Note 7**). This suggests that the noise present in this cluster-wide

291 representation specifically contributes to suboptimal confidence ratings over and above perceptual noise.

292 Moreover, the precision of the confidence representation persisted through the last four samples of the More

293 condition (**Figure 6b**), as expected of a signal that continues to process evidence for confidence after

294 perceptual decision commitment. In contrast, the early posterior representation found to be relevant for

295 perceptual decision-making did show attenuation for the last four samples of the More condition (a repeated

296 measures ANOVA revealed a significant interaction between cluster and condition for decoding precision in

297 the last four samples, $F(1,19) = 32.00$, $p = 0.001$, Bonferroni corrected for three comparisons; **Figure 6b**),

298 and the perceptual representation error was unrelated to suboptimal confidence ratings (in fact the

299 evidence was in favour of the null hypothesis; summed log likelihood ratio = -1176). These results are

300 consistent with dissociable stages of neural processing for confidence evaluation and perceptual decision-

301 making.

302 Greater error in the confidence representation of accumulated evidence was associated with greater model

303 estimated inference error and suboptimal behavioural confidence evaluations. We examined the sources of

304 the EEG representation error by comparing 'Noise Min' and 'Noise Max' epochs (the top and bottom quartile

305 of epochs sorted by the confidence representation precision). The presented sensory evidence was similar

306 across these epochs (see **Supplementary Note 7**), but the additional noise in the Noise Max epochs pushes

307 the represented evidence further from the mean, and should therefore correspond to a greater absolute

308 normalised signal. We estimated the sources of activity in the Noise Min and Noise Max epochs using a

309 template brain (**Figure 6c**; see **Methods**) and tested for differences in the rectified normalised current

310 density in ROIs defined based on the previous literature (**Figure 6d**; Graziano, Parra, and Sigman, 2015;

311 German and Philiastides, 2018; Herding et al., 2019, see **Supplementary Note 9**). As expected, Noise Max

312 epochs showed a greater increase in current density power over time. Significant differences first emerged

313 in the Superior Parietal cortex (**Figure 6e**; 276 - 304 ms; $t_{avg}(19) = 2.37$, $p_{cluster} = 0.016$, re-emerging at 596 –

314 748 ms; $t_{avg}(19) = 2.53$, $p_{cluster} = 0.016$; and 912 ms; $t_{avg}(19) = 2.50$, $p_{cluster} = 0.014$), and then in the

315 Orbitofrontal cortex (OFC; 516 – 556 ms; $t_{avg}(19) = 2.30$, $p_{cluster} = 0.022$, re-emerging at 660 – 772 ms; $t_{avg}(19)$

316 $= 2.79$, $p_{cluster} = 0.032$, and 824 – 1000 ms; $t_{avg}(19) = 2.60$, $p_{cluster} = 0.022$). No differences in the rostral Middle

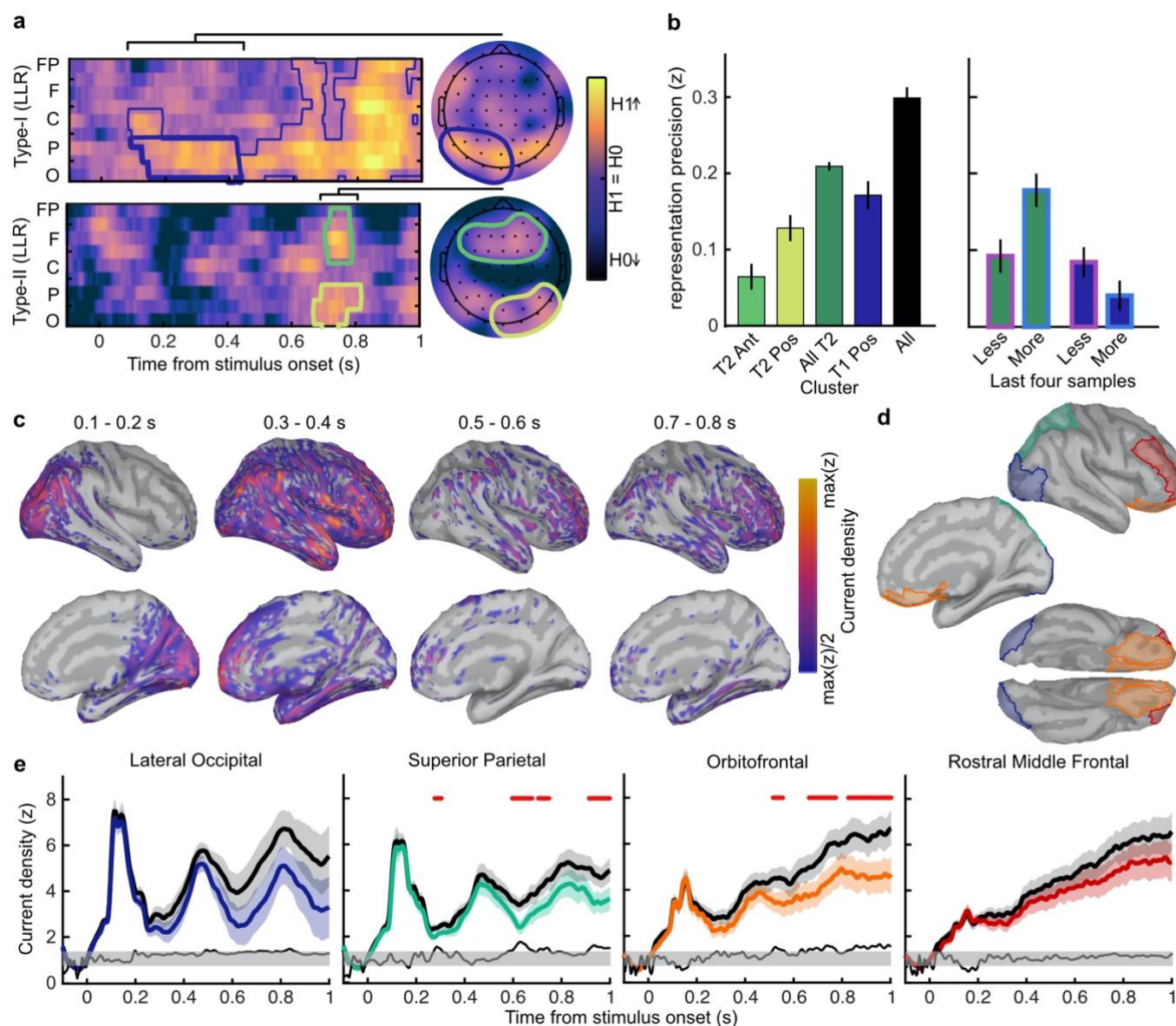317 Frontal cortex nor Lateral Occipital cortex survived cluster correction.

318

*Figure 6. Clusters of behaviourally relevant representations and their sources. a) Log likelihood ratio (LLR) of the data given the hypothesis that decoding precision varies with behavioural suboptimalities, against the null hypothesis that decoding precision varies only with measurement noise. Perceptual (Type-I) behaviour is shown on top and confidence (Type-II) behaviour is shown on the bottom. Clusters where the log posterior odds ratio outweighed the prior are circled, only the bold area of the perceptual cluster was further analysed. Time series (left) show the maximum LLR of electrodes laterally, with frontal polar electrodes at the top descending to occipital electrodes at the bottom. Scalp maps (right) show the summed LLR over the indicated time windows. b) Left: representation precision (z) training and testing on signals within the clusters. Colours correspond to the circles in a), with the dark green bar showing the combined decoding precision of the anterior and posterior confidence clusters, and the black bar showing the combined representation precision of all clusters. Right: Representation precision of the last four samples in the Less and the More conditions for the combined confidence representation and the perceptual representation. Error bars show 95% within-subject confidence intervals. c) Average rectified normalised current density in Noise Min epochs for the corresponding time windows, filtered above the half-maximum amplitude. d) ROIs (defined by mindBoggle coordinates; Klein et al., 2017): Lateral Occipital cortex (blue); Superior Parietal cortex (green); Orbitofrontal cortex (orange);*

334 *and Rostral Middle Frontal cortex (red).* **e)** *ROI time series for Noise Max (black) and Noise Min (coloured)*
335 *epochs, taking the average rectified normalised current density (z) across participants. Shaded regions show*
336 *95% within-subject confidence intervals, red horizontal lines indicate cluster corrected significant differences.*
337 *Standardised within-subject differences are traced above the x-axis, with the shaded region marking z = 0 to z =*
338 *1.96 (95% confidence).*

## Discussion

340 We examined the dynamic neural signals associated with suboptimal accumulation of evidence for
341 evaluating confidence in perceptual decisions. Observers were required to integrate evidence over multiple
342 samples provided by a sequence of visual stimuli. When observers were unable to control the amount of
343 evidence they were exposed to, they employed a covert decision bound, committing to decisions when they
344 had enough evidence, even if stimulus presentation continued. We had previously shown evidence for this
345 premature decision commitment based on behaviour and computational modelling (Balsdon, Wyart and
346 Mamassian, 2020). We replicated these results here, and further examined the neural signatures of covert
347 decision making. We found that the distribution of spectral power associated with preparing a motor
348 response in the Free task (where the response is entered as soon as the decision is made) could be used to
349 accurately predict responses in the More condition of the Replay task over 1 s prior to when the response
350 was entered, and with significantly greater sensitivity than in the Less condition (when observers were
351 unlikely to have committed to a decision early). This suggests that covert decisions could trigger the motor
352 preparation for pressing the response key. Moreover, the strength of the eventual motor response signal
353 could be predicted by earlier decision evidence in the More condition, as if observers are maintaining some
354 representation of the decision evidence whilst waiting to press the response key.

355 Based on the evoked representation of accumulated evidence, perceptual decision accuracy relied on a flow
356 of information processing from early Occipital and Parietal signals, which then spread through to anterior
357 electrodes. When observers committed to perceptual decisions prematurely, only the early part of the
358 representation of accumulated evidence was attenuated. This selective dampening of the representation of
359 accumulated evidence following premature decision commitment delineates which computations are
360 devoted solely to the perceptual decision process, and which computations reflect the input to the decision
361 process: The representations of stimulus orientation and decision update (Wyart et al., 2012; Wyart et al.,
362 2015; Weiss et al., 2019), which are necessary input for the perceptual decision, did not substantially change
363 after committing to a perceptual decision. This initial perceptual processing stage likely remained important
364 for the continued accumulation of evidence for evaluating confidence (even after the completion of
365 perceptual decision processes), though it could also be that these processes are automatically triggered by
366 stimulus onset irrespective of whether the evidence is being accumulated for decision-making.

367 Confidence should increase with increasing evidence for the perceptual decision. It is therefore unsurprising
368 that the neural correlates of confidence magnitude have found similar EEG markers as those related to the
369 accumulation of the underlying perceptual decision evidence: the P300 (Gherman and Philiastides, 2015;

13

370  Desender et al., 2016; Desender et al., 2019; Zakrzewski et al., 2019; Rausch et al., 2020); and Central

371  Parietal Positivity (CPP; Boldt et al., 2019; Herding et al., 2019, indeed we show a similar effect in **Figure 3**).

372  The analysis presented in this manuscript targeted confidence precision rather than confidence magnitude,

373  by assessing confidence relative to an optimal observer who gives high confidence ratings on trials where

374  the evidence in favour of the perceptual choice is greater than the median across trials. We isolated part of

375  the representation of accumulated evidence where greater error in the representation was followed by

376  suboptimal confidence ratings, and showed that this was also associated with greater error estimated by the

377  computational model fit to describe confidence behaviour.

378  The precision of the confidence representation was found to be disrupted by noise localised to the Superior

379  Parietal and Orbitofrontal cortices. This is not at odds with the previous literature: The difference in

380  Superior Parietal cortex could be linked with findings from electrophysiology that suggest that confidence is

381  based on information coded in Parietal cortex, where the underlying perceptual decision evidence is

382  integrated (Kiani et al., 2009; Rutishauser et al., 2018; though at least a subset of these neurons reflect

383  bounded accumulation, which is in contrast with the continued confidence accumulation described in this

384  experiment; Kiani, Hanks, and Shadlen, 2007). Early electrophysiological investigation into the function of

385  the Orbitofrontal cortex revealed neural coding associated with stimulus value (Thorpe, Rolls, and

386  Maddison, 1983), which has since been linked with a confidence-modulated signal of outcome-expectation

387  (Kepecs et al., 2008; and in human fMRI; Rolls, Grabenhorst, and Deco, 2010) and recently, shown to be

388  domain-general (single OFC neurons were associated with confidence in both olfactory and auditory tasks;

389  Masset et al., 2020). The source localisation analysis therefore connects previous findings, indicating

390  confidence feeds off an evidence accumulation process, culminating in higher-order brain areas that appear

391  to function for guiding outcome-driven behaviour based on decision certainty.

392  These neural signatures of confidence evidence encoding were present throughout the process of making a

393  perceptual decision. This is in line with more recent evidence suggesting that confidence could be computed

394  online, alongside perceptual evidence accumulation (Zizlsperger et al., 2014; Gherman and Philiastides,

395  2015; Balsdon et al., 2020), as opposed to assessing the evidence in favour of the perceptual decision only

396  after committing to that decision. Computational model comparison supported this interpretation, showing

397  the best description of confidence behaviour was an accumulation process that was partially dissociable

398  from perceptual evidence accumulation (**Supplementary Note 1**; replicating our previous analysis, Balsdon

399  et al., 2020). This partial dissociation mediates the ongoing debate between single-channel (for example,

400  Maniscalco and Lau, 2016) and dual-channel (for example, Charles, King, and Deheane 2014) models, as it

401  constrains confidence by perceptual suboptimalities, at the same time as allowing additional processing to

402  independently shape confidence. The combination of this partial dissociation and online monitoring could

403  allow for metacognitive control of perceptual evidence accumulation, to flexibly balance perceptual accuracy

404  against efficiency by bounding perceptual evidence accumulation according to contemporaneous confidence.

405  Using this protocol, we were able to delineate two distinct representations of accumulated evidence which

406  correspond to perceptual decision-making and confidence evaluations. These neural representations were

partially dissociable in that the perceptual representation neglected additional evidence following premature decision commitment whilst the confidence representation continued to track the updated evidence independently of decision commitment. This partial dissociation validates the predictions of the computational model and provides a framework for the cognitive architecture underlying the distinction between perception and confidence. That the neural resources involved in the confidence representation can be recruited independently of perceptual processes implies a specific neural circuit for the computation of confidence, a necessary feature of a general metacognitive mechanism flexibly employed to monitor the validity of any cognitive process.

# Methods

## Participants

A total of 20 participants were recruited from the local cognitive science mailing list (RISC) and by word of mouth. No participant met the pre-registered (https://osf.io/346pe/?view_only=ddbc092996f34438964cf45a239498bb) exclusion criteria of chance-level performance or excessive EEG noise. Written informed consent was provided prior to commencing the experiment. Participants were required to have normal or corrected to normal vision. Ethical approval was granted by the INSERM ethics committee (ID RCB: 2017-A01778-45 Protocol C15-98).

## Materials

Stimuli were presented on a 24" BenQ LCD monitor running at 60 Hz with resolution 1920x1080 pixels and mean luminance 45 cd/m². Stimulus generation and presentation was controlled by MATLAB (Mathworks) and the Psychophysics toolbox (Brainard, 1997; Pelli, 1997; Kleiner et al., 2007), run on a Dell Precision M4800 Laptop. Observers viewed the monitor from a distance of 57 cm, with their head supported by a chin rest. EEG data were collected using a 64-electrode BioSemi ActiveTwo system, run on a dedicated mac laptop (Apple Inc.), with a sample rate of 512 Hz. Data were recorded within a shielded room.

## Stimuli

Stimuli were oriented Gabor patches displayed at 70% contrast, subtending 4 dva and with spatial frequency 2 cyc/deg. On each trial a sequence of stimuli was presented, at an average rate of 3 Hz, with the stimulus presented at full 70% contrast for a variable duration between 50 and 83 ms, with a sudden onset, followed by an offset ramp over two flips, where the stimulus contrast decreased by 50% and 75% before complete offset. Stimulus onset timing was jittered within the stimulus presentation interval such that the timing of stimulus onset was irregular but with at least 216 ms between stimuli. These timings and stimulus examples are shown in **Figure 1a**.

On each trial the orientations of the presented Gabors were drawn from one of two circular Gaussian (Von Mises) distributions centred on +/- 45° from vertical (henceforth referred to as the 'orange' and 'blue' distributions respectively), with concentration κ = 0.5 (shown in **Figure 1d**). Stimuli were displayed within an annular 'colour-guide' where the colour of the annulus corresponds to the probability of the orientation under each distribution, using the red and blue RGB channels to represent the probabilities of each

15

443  orientation under each distribution. Stimuli were presented in the centre of the screen, with a black central
444  fixation point to guide observers' gaze.

## Procedure

446  The task was a modified version of the weather prediction task (Knowlton et al., 1996; Drugowitsch et al.,
447  2016). Throughout the experiment, the observer's perceptual task was to categorise which distribution the
448  stimulus orientations were sampled from. They were instructed to press the 'd' key with their left hand (of a
449  standard querty keyboard) for the blue distribution and the 'k' key with their right hand for the orange
450  distribution. There were two variants of the task: The Free task and the Replay task. The trials were
451  composed of three repetitions of 100 predefined sequences of up to 40 samples (50 trials from each
452  distribution) for each observer (300 trials per task).

453  In the 'Free' task, observers were continually shown samples (up to 40) until they entered their response.
454  They were instructed to enter their response as soon as they 'feel ready' to make a decision, with emphasis
455  on both accuracy (they should make their decision when they feel they have a good chance of being correct)
456  and on time (they shouldn't take too long to complete each trial). A graphical description of this task is
457  shown in **Figure 1b**.

458  After completing the Free task, observers then completed the Replay task. In this task they were shown a
459  specific number of samples and could only enter their response after the sequence finished, signalled by the
460  fixation point turning red. The number of samples was determined based on the number observers chose to
461  respond to in the Free task. There were three intermixed conditions: In the Less condition observers were
462  shown two fewer samples than the minimum they had chosen to respond to on that predefined sequence in
463  the Free task; In the Same condition observers were shown the median number of samples from that
464  predefined sequence; in the More condition observers were shown four additional samples compared to the
465  maximum number they chose to respond to on that sequence in the Free task. After entering their
466  perceptual (Type-I) response, observers were cued to give a confidence rating (Type-II decision). The
467  confidence rating was given on a 4-point scale where 1 represents very low confidence that the perceptual
468  decision was correct, and 4, certainty that the perceptual decision was correct. The rating was entered by
469  pressing the 'space bar' when a presented dial reached the desired rating.  The dial was composed of a black
470  line which was rotated clockwise to each of 4 equidistant angles (marked 1 - 4) around a half circle, at a rate
471  of 1.33 Hz. The dial started at a random confidence level on each trial and continued updating until a rating
472  was chosen. A graphical description of this task is shown in **Figure 1c**.

473  Prior to commencing the experimental trials, participants were given the opportunity to practice the
474  experiment and ask questions. They first performed 20 trials of a fixed number of samples with only the
475  perceptual decision, with feedback after each response as to the true category. They then practiced the
476  Replay task with the confidence rating (and an arbitrary number of samples). Finally, they practiced the Free
477  task, before commencing the experiment with the Free task.

16

478 **Analysis**

479 **Behaviour**

480 Perceptual (Type-I) decisions were evaluated relative to the category the orientations were actually drawn

481 from. Performance is presented as proportion correct, whilst statistical analyses were performed on

482 sensitivity (d'). Confidence was evaluated relative to an optimal observer who gives high confidence when

483 the log-likelihood of the chosen category, based on the presented orientations, is above the median across

484 trials, and low confidence on trials with less than the median log-likelihood. More broadly, confidence should

485 increase with increasing evidence in favour of the perceptual decision, see **Supplementary Note 3**.

486 **Computational modelling**

487 Computational modelling followed the same procedure as Balsdon, Wyart, and Mamassian (2020). The

488 model parametrically describes suboptimalities relative to the Bayesian optimal observer. The Bayesian

489 optimal observer knows the category means, $\mu_1 = -\frac{\pi}{4}, \mu_2 = \frac{\pi}{4}$, and the concentration, $\kappa = 0.5$, and takes the

490 probability of the orientation $\theta_n$ (at sample $n$) given each category $\psi$ ($\psi = 1$ or $\psi = 2$):

$$p(\theta_n \mid \psi) = \frac{e^{\kappa \cos (2(\theta_n - \mu_\psi))}}{\pi I_0(\kappa)} \tag{1}$$

491 Where $I_0(\cdot)$ is the modified Bessel function of order 0. The optimal observer then chooses the category

492 $\psi$ with the greatest posterior probability over all samples for that trial, $T$ ($T$ varies from trial to trial). Given a

493 uniform category prior, $p(\psi) \propto \frac{1}{2}$, and perfect anticorrelation in $p(\theta_n \mid \psi)$ over the categories, the log

494 posterior is proportional to the sum of the difference in the log-likelihood for each category ($\ell_n = \ell_{n,1} -$

495 $\ell_{n,2}$):

$$z = \sum_{n=1}^{T} \ell_n \tag{2}$$

496 Where:

$$\ell_{n,\psi} = \log p(\theta_n \mid \psi) = \kappa \cos (2(\theta_n - \mu_\psi)) + const. \tag{3}$$

497 Such that the Bayesian optimal decision is 1 if $z > 0$ and 2 if $z \leq 0$.

498 The suboptimal observer suffers inaccuracies in the representation of each evidence sample, captured by

499 additive independent identically distributed (i.i.d) noise, $\varepsilon_n$. The noise is Gaussian distributed with zero

500 mean, and the degree of variability parameterised by $\sigma$, the standard deviation:

$$\varepsilon_n \sim N(0, \sigma^2) \tag{4}$$

501 The evidence over samples is also imperfectly accumulated, incurring primacy or recency biases

502 parameterised by $\alpha$, the weight on the current accumulated evidence compared to the new sample ($\alpha > 1$

503 creates a primacy effect). By the end of the trial, the weight on each sample $n$ is equal to:

$$v_n = \alpha^{T-n} \qquad \qquad \textbf{(5)}$$

504    Where $T$ is the eventual total samples on that trial and $n \in [1, T]$.

505    In the Free task the observer responds when accumulated evidence reaches a bound, $\Lambda$. The optimal

506    observer sets a constant bound on proportion correct over sequence length, which is an exponential function

507    on the average evidence over the samples accumulated. The human observer can set the scale, $b$, and the rate

508    of decline, $\lambda$, of the bound suboptimally, resulting in:

$$\Lambda_{n+} = n \times \left( a + be^{-\frac{n}{\lambda}} \right) \qquad \qquad \textbf{(6)}$$

509    for the positive decision bound (the negative bound, $\Lambda_{n-} = -\Lambda_{n+}$). The likelihood $f(n)$ of responding at

510    sample n was estimated by computing the frequencies, over 1000 samples from $\varepsilon_n$ (Monte Carlo simulation),

511    of first times where the following inequality is verified:

$$\left| \sum_{n=1}^{N} (\ell_n + \varepsilon_n) \cdot v_n \right| > \Lambda_n \qquad \qquad \textbf{(7)}$$

512    The response time, relative to reaching the decision bound, is delayed by non-decision time for executing the

513    motor response, which is described by a Gaussian distribution of mean, $\mu_U$, and variance, $\sigma_U^2$.

### Model fitting

515    Parameters were optimised to minimise the negative log-likelihood of the observer making response $r$ on

516    sample $n$ on each trial for each participant using Bayesian Adaptive Direct Search (Acerbi and Ma, 2017). The

517    log-likelihoods were estimated using Monte Carlo Simulation, with the sensitivity of this approach being

518    addressed in previous work (Balsdon et al., 2020). The full model was simplified using a knock-out

519    procedure based on Bayesian Model Selection (Rigoux et al., 2014) to fix the bias (exceedance probability =

520    0.93) and lapse (exceedance probability >0.99) parameters (not described above, see **Supplementary Note**

521    **1**).

522    In the Replay task, confidence ratings were fit using the same model described above, but with additional

523    criteria determining confidence ratings, described by three bounds on the confidence evidence,

524    parameterised in the same manner as the decision bound. These models were then used to simulate the

525    internal evidence of each observer from sample to sample, and the error compared to the ideal evidence

526    (uncorrupted by suboptimalities, see **Supplementary Note 2**).

### EEG pre-processing

528    EEG data were pre-processed using the PREP processing pipeline (Bigdely-Shamlo, et al., 2015),

529    implemented in EEGlab (v2019.0, Delorme & Makeig, 2004) in MATLAB (R2019a, Mathworks). This includes

530    line noise removal (notch filter at 50 Hz and harmonics) and re-referencing (robust average re-reference on

531    data detrended at 1 Hz). The data were then filtered to frequencies between 0.5 and 80 Hz, and down-

532    sampled to 256 Hz. Large epochs were taken locked to each stimulus (-500 to 1500 ms) and each response (-

533    5000 to 1500 ms). Independent Components Analysis was used to remove artefacts caused by blinks and

534 excessive muscle movement identified using labels with a probability greater than 0.35 from the ICLabel

535 project classifier (Swartz Centre for Computational Neuroscience).

**Classical analyses**

537 We present several 'classical' comparisons, examining the effect of confidence on EEG amplitude

538 (microvolts). In **Figure 3**, we show Central Parietal Positivity (CPP; O'Connell et al., 2012; average amplitude

539 of electrodes CP1, CP2, and CPz over response locked epochs), and the Lateralised Readiness Potential (LRP;

540 difference in microvolts between the average of electrodes [C1, C3], and [C2, C4], signed by response hand;

541 Deecke et al., 1976). In all cases, the group average within condition over the 100 ms prior to the first

542 stimulus of each trial was used as a baseline, the data were otherwise unfiltered except for the pre-

543 processing.

**Response classification analysis**

545 The power spectrum across frequency tapers from 1 to 64 Hz with 25% spectral smoothing was resolved

546 using wavelet convolution implemented in FieldTrip (Oostenveld et al., 2011). The epochs were then clipped

547 at -3 to 1 s around the time of entering the perceptual response.  Linear discriminant analysis was

548 performed to classify perceptual responses, using 10-fold cross validation, separately on each taper at each

549 time-point. An analysis of the frequencies contributing to accurate classification at the time of the response

550 revealed significant contributions from 8 to 26 Hz (**Supplementary Note 4**). We therefore continued by

551 using the power averaged across these frequency bands to train and test the classifier. Classifier accuracy

552 was assessed using the area under the receiver operating characteristic curve (AUC). At the single-trial level,

553 the probability of the response based on the classifier was computed from the relative normalised Euclidean

554 distance of the trial features from the response category means in classifier decision space.

**Encoding Variable Regression**

556 We used a linear regression analysis to examine the EEG correlates of different aspects of the decision

557 evidence (encoding variables) in epochs locked to stimulus onset. Regularised ridge regression (ridge $\lambda = 1$)

558 was used to predict the encoding variables based on EEG data, over 10-fold cross validation. The precision of

559 the representation of each encoding variable was computed within each observer by taking the Fisher

560 transform of the correlation coefficient (Pearson's r) between the encoded variable and predicted variable.

561 To maximise representation precision, the data were bandpass filtered (1 – 8 Hz) and decomposed into real

562 and imaginary parts using a Hilbert Transform (**Supplementary Note 5**). For each time point, the data from

563 all electrodes were used to predict the encoded variable. The temporal generalisation of decoding weights

564 was examined by training at one time point and testing at another. The contribution of information from

565 signals at each electrode was examined by training and testing on the signals at each electrode at each time

566 point (further details in **Supplementary Note 5**).

567 Behaviourally relevant signals were isolated by comparing representation precision at each time point and

568 electrode for epochs leading to optimal and suboptimal perceptual and confidence responses. Cluster

569 modelling was used to isolate contiguous signals where the log posterior odds were in favour of the

570   alternative hypothesis that representation precision was affected by inference noise beyond what could be

571   explained by measurement noise alone (**Supplementary Note 6**). New regression weights were then

572   calculated on signals from the entire cluster and representation errors calculated as the difference of the

573   predicted variable from the expected value given the representation precision.

574   **Source Localisation**

575   Identifying the clusters of signals associated with confidence processes offers relatively poor spatial and

576   temporal (given the bandpass filter; de Cheveigné, and Nelken, 2019) resolution for identifying the source of

577   the suboptimalities affecting confidence ratings. Source localisation was therefore performed, using

578   Brainstorm (Tadel et al., 2011). The forward model was computed using OpenMEEG (Gramfort et al., 2010;

579   Kybic et al., 2005) and the ICBM152 anatomy (Fonov et al., 2011; 2009). Two conditions were compared,

580   Noise Min and Noise Max, which corresponded to quartiles of epochs sorted by representation error in the

581   confidence clusters (see **Supplementary Note 7** for more details). Cortical current source density was

582   estimated from the average epochs using orientation-constrained minimum norm imaging (Baillet, Mosher,

583   and Leahy, 2001). ROIs in the Lateral Occipital, Superior Parietal, Rostral Middle Frontal (including dlPFC),

584   Medial Orbitofrontal, and rostral Anterior Cingulate Cortex, were defined using MindBoggle coordinates

585   (Klein et al., 2017). Statistical comparisons were performed on the bilateral ROI time series (using cluster

586   correction and a minimum duration of 20 ms), computed over separate conditions on rectified normalised

587   subject averages (low-pass filtered at 40 Hz).

588

589 **References**

590 Acerbi, L., & Ma, W. J. Practical Bayesian optimization for model fitting with Bayesian adaptive direct search.
591     In *Advances in Neural Information Processing Systems,* December 2017; 1836-1846

592 Bahrami, B., Olsen, K., Bang, D., Roepstorff, A., Rees, G., & Frith, C. What failure in collective decision-making
593     tells us about metacognition. *Philosophical Transactions of the Royal Society B: Biological*
594     *Sciences*, 2012; **367**(1594), 1350-1365.

595 Baillet, S., Mosher, J. C., & Leahy, R. M. Electromagnetic brain mapping. *IEEE Signal Processing*
596     *Magazine*, 2001; **18**(6), 14-30.

597 Balsdon, T., Wyart, V., & Mamassian, P. Confidence controls perceptual evidence accumulation. *Nature*
598     *Communications,* 2020; **11**(1), 1-11

599 Bang, J. W., Shekhar, M., & Rahnev, D. Sensory noise increases metacognitive efficiency. *Journal of*
600     *Experimental Psychology: General,* 2019; **148**(3), 437.

601 Baranski, J. V., & Petrusic, W. M. The calibration and resolution of confidence in perceptual
602     judgments. *Perception & Psychophysics*, 1994; **55**(4), 412-428.

603 Bigdely-Shamlo, N., Mullen, T., Kothe, C., Su, K. M., & Robbins, K. A. The PREP pipeline: standardized
604     preprocessing for large-scale EEG analysis. *Frontiers in Neuroinformatics*, 2015; **9**, 16

605 Boldt, A., Schiffer, A. M., Waszak, F., & Yeung, N. Confidence predictions affect performance confidence and
606     neural preparation in perceptual decision making. *Scientific Reports*, 2019; **9**(1), 1-17.

607 Brainard, D. H. The psychophysics toolbox. *Spatial Vision*, 1997; **10**(4), 433-436.

608 Charles, L., King, J. R., & Dehaene, S. Decoding the dynamics of action, intention, and error detection for
609     conscious and subliminal stimuli. *Journal of Neuroscience*, 2014; **34**(4), 1158-1170.

610 de Cheveigné, A., & Nelken, I. Filters: when, why, and how (not) to use them. *Neuron*, 2019; **102**(2), 280-293.

611 Deecke, L., Grozinger, B., & Kornhuber, H. H. Voluntary finger movement in man: Cerebral potentials and
612     theory. *Biological Cybernetics*, 1976; **23**, 99–119.

613 Delorme, A., & Makeig, S. EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including
614     independent component analysis. *Journal of Neuroscience Methods*, 2004; **134**(1), 9-21.

615 Desender, K., Van Opstal, F., Hughes, G., & Van den Bussche, E. The temporal dynamics of metacognition:
616     Dissociating task-related activity from later metacognitive processes. *Neuropsychologia*, 2016; **82**, 54-
617     64.

618 Desender, K., Murphy, P., Boldt, A., Verguts, T., & Yeung, N. A postdecisional neural marker of confidence
619     predicts Information-Seeking in Decision-Making. *Journal of Neuroscience*, 2019; **39**(17), 3309-3319.

620 Drugowitsch, J., Wyart, V., Devauchelle, A. D., & Koechlin, E. Computational precision of mental inference as
621     critical source of human choice suboptimality. *Neuron*, 2016; **926**, 1398-1411

622 Fleming, S. M., & Daw, N. D. Self-evaluation of decision-making: A general Bayesian framework for
623     metacognitive computation. *Psychological Review*, 2017; **124**(1), 91.

624 Fonov VS, Evans AC, McKinstry RC, Almli CR, Collins DL. Unbiased nonlinear average age-appropriate brain
625     templates from birth to adulthood. *NeuroImage*, 2009; **47**, S102.

626 Fonov, V., Evans, A. C., Botteron, K., Almli, C. R., McKinstry, R. C., Collins, D. L., & Brain Development

627        Cooperative Group. Unbiased average age-appropriate atlases for pediatric studies. *Neuroimage*, 2011;
628          **54**(1), 313-327.

629    Frith, C. D. The role of metacognition in human social interactions. *Philosophical Transactions of the Royal*
630          *Society B: Biological Sciences*, 2012; **367**(1599), 2213-2223.

631    Gherman, S., & Philiastides, M. G. Human VMPFC encodes early signatures of confidence in perceptual
632          decisions. *eLife*, 2018; **7**, e38293.

633    Gherman, S., & Philiastides, M. G. Neural representations of confidence emerge from the process of decision
634          formation during perceptual choices. *NeuroImage*, 2015; **106**, 134-143.

635    Gramfort, A., Papadopoulo, T., Olivi, E., & Clerc, M. OpenMEEG: opensource software for quasistatic
636          bioelectromagnetics. *Biomedical Engineering Online*, 2010; **9**(1), 45.

637    Graziano, M., Parra, L. C., & Sigman, M. Neural correlates of perceived confidence in a partial report
638          paradigm. *Journal of Cognitive Neuroscience*, 2015; **27**(6), 1090-1103.

639    Helmholtz, H.L.F.v. *Treatise on Physiological Optics,* Thoemmes Press 1856.

640    Herding, J., Ludwig, S., von Lautz, A., Spitzer, B., & Blankenburg, F. Centro-parietal EEG potentials index
641          subjective evidence and confidence during perceptual decision making. *NeuroImage*, 2019; **201**,
642          116011.

643    Kepecs, A., Uchida, N., Zariwala, H. A., & Mainen, Z. F. Neural correlates, computation and behavioural impact
644          of decision confidence. *Nature*, 2008; **455**, 227–231.

645    Kiani, R., & Shadlen, M. N. Representation of confidence associated with a decision by neurons in the parietal
646          cortex. *Science*, 2009; **324**, 759–764.

647    Kiani, R., Corthell, L., & Shadlen, M. N. Choice certainty is informed by both evidence and decision
648          time. *Neuron*, 2014; **84**(6), 1329-1342.

649    Kiani, R., Hanks, T. D., & Shadlen, M. N. Bounded integration in parietal cortex underlies decisions even when
650          viewing duration is dictated by the environment. *Journal of Neuroscience*, 2008; **28**(12), 3017-3029.

651    Klein, A., Ghosh, S. S., Bao, F. S., Giard, J., Häme, Y., Stavsky, E., ... & Keshavan, A. Mindboggling morphometry
652          of human brains. *PLoS Computational Biology*, 2017; **13**(2), e1005350.

653    Kleiner, M., Brainard, D., & Pelli, D. What's new in Psychtoolbox-3? 2007.

654    Knowlton, B. J., Mangels, J. A., & Squire, L. R. A neostriatal habit learning system in humans. *Science*, 1996;
655          **273**(5280), 1399-1402

656    Kybic, J., Clerc, M., Abboud, T., Faugeras, O., Keriven, R., & Papadopoulo, T. A common formalism for the
657          integral formulations of the forward EEG problem. *IEEE Transactions on Medical Imaging*, 2005; **24**(1),
658          12-28.

659    Maniscalco, B., & Lau, H. The signal processing architecture underlying subjective reports of sensory
660          awareness. *Neuroscience of Consciousness*, 2016; **1**.

661    Masset, P., Ott, T., Lak, A., Hirokawa, J., & Kepecs, A. Behavior- and modality-general representation of
662          confidence in orbitofrontal cortex. *Cell*, 2020; **182**(1), 112-126.

663    Mazancieux, A., Fleming, S., Souchay, C., & Moulin, C. Retrospective confidence judgments across tasks:
664          domain-general processes underlying metacognitive accuracy. *BioRxiv* 2018.

665    Moreno-Bote, R. Decision confidence and uncertainty in diffusion models with partially correlated neuronal

666      integrators. *Neural Computation*, 2010; **22**, 1786–1811.

667 Neill, D. B. Fast Bayesian scan statistics for multivariate event detection and visualization. *Statistics in*
668      *Medicine*, 2011; **30**(5), 455-469.

669 Neill, D. B. Bayesian Scan Statistics. In: Glaz J., Koutras M. (eds) *Handbook of Scan Statistics*. Springer, New
670      York, NY. 2019.

671 O'Connell, R. G., Dockree, P. M., & Kelly, S. P. A supramodal accumulation-to-bound signal that determines
672      perceptual decisions in humans. *Nature Neuroscience*, 2012; **15**(12), 1729.

673 Oostenveld, R., Fries, P., Maris, E., & Schoffelen, J. M. FieldTrip: open source software for advanced analysis of
674      MEG, EEG, and invasive electrophysiological data. *Computational Intelligence and Neuroscience*, 2011.

675 Pelli, D. G.  The VideoToolbox software for visual psychophysics: Transforming numbers into movies. *Spatial*
676      *Vision*, 1997; **10**, 437-442.

677 Pleskac, T. J., & Busemeyer, J. R. Two-stage dynamic signal detection: a theory of choice, decision time, and
678      confidence. *Psychological Review*, 2010; **117**(3), 864.

679 Pollack, I., & Decker, L. R. Confidence ratings, message reception, and the receiver operating
680      characteristic. *The Journal of the Acoustical Society of America*, 1958; **30**(4), 286-292.

681 Ratcliff, R. A theory of memory retrieval. *Psychological Review*, 1987; **85**(2), 59.

682 Rausch, M., Zehetleitner, M., Steinhauser, M., & Maier, M. E. Cognitive modelling reveals distinct
683      electrophysiological markers of decision confidence and error monitoring. *NeuroImage*, 2020; **218**,
684      116963.

685 Rigoux, L., Stephan, K.E., Friston, K.J. & Daunizeau, J. Bayesian Model Selection for Group Studies Revisited.
686      *NeuroImage* 2014; **84,** 971-85.

687 Rolls, E. T., Grabenhorst, F., & Deco, G. Choice, difficulty, and confidence in the brain. *NeuroImage,* 2010;
688      **53**(2), 694-706.

689 Rutishauser, U., Aflalo, T., Rosario, E. R., Pouratian, N., & Andersen, R. A. Single-neuron representation of
690      memory strength and recognition confidence in left human posterior parietal cortex. *Neuron*, 2018;
691      **97**(1), 209-220.

692 Shekhar, M., & Rahnev, D. Sources of Metacognitive Inefficiency. *Trends in Cognitive Sciences*.

693 Tadel, F., Baillet, S., Mosher, J. C., Pantazis, D., & Leahy, R. M. Brainstorm: a user-friendly application for
694      MEG/EEG analysis. *Computational Intelligence and Neuroscience*, 2011.

695 Thorpe, S. J., Rolls, E. T., & Maddison, S. The orbitofrontal cortex: neuronal activity in the behaving
696      monkey. *Experimental Brain Research*, 1983; **49**(1), 93-115.

697 Veenman, M. V., Wilhelm, P., & Beishuizen, J. J. The relation between intellectual and metacognitive skills
698      from a developmental perspective. *Learning and Instruction*, 2004; **14**(1), 89-109.

699 Vickers, D. Evidence for an accumulator model of psychophysical discrimination. *Ergonomics*, 1970; **13**(1),
700      37-58.

701 Vickers, D. *Decision processes in visual perception.* New York, NY: Academic Press. 1979.

702 Weiss, A., Chambon, V., Drugowitsch, J., & Wyart, V. Interacting with volatile environments stabilizes hidden-
703      state inference and its brain signatures. *bioRxiv*, 2019

704 Wyart, V., De Gardelle, V., Scholl, J., & Summerfield, C. Rhythmic fluctuations in evidence accumulation during

705      decision making in the human brain. *Neuron*, 2012; **76**(4), 847-858.

706      Wyart, V., Myers, N. E., & Summerfield, C. Neural mechanisms of human perceptual choice under focused and

707      divided attention. *Journal of Neuroscience*, 2015; **35**(8), 3485-3498.

708      Zakrzewski, A. C., Wisniewski, M. G., Iyer, N., & Simpson, B. D. Confidence tracks sensory-and decision-related

709      ERP dynamics during auditory detection. *Brain and Cognition*, 2019; **129**, 49-58

710      Zizlsperger, L., Sauvigny, T., Händel, B., & Haarmeier, T. Cortical representations of confidence in a visual

711      perceptual decision. *Nature Communications*, 2014; **5**(1), 1-13.

712

1 Supplementary materials

## 2 **Supplementary Note 1**

### 3 **Computational Model fitting**

4 The computational model is described in full in the **Methods** section. Briefly, the model is based on the

5 Bayesian optimal observer with full knowledge of the category distributions (means $\mu_1$ and $\mu_2$,

6 concentration $\kappa$), and takes as evidence the difference in the log posterior probability ($\ell_n$) of each category

7 given the orientation ($\theta_n$)

$$\ell_n = \ell_{n,1} - \ell_{n,2} = \kappa \cos\big(2(\theta_n - \mu_1)\big) - \kappa \cos\big(2(\theta_n - \mu_2)\big) \tag{1}$$
$$= 2\,\kappa \sin(\mu_1 - \mu_2) \sin(2\theta_n - \mu_1 - \mu_2) = \sin(2\theta_n)$$

8 where chosen values ($\kappa = 0.5$, $\mu_1 = -\pi/4$, and $\mu_2 = \pi/4$) have been implemented in the last equation.

9 Whilst the optimal observer perfectly sums the evidence over each sample, the suboptimal human observer

10 accumulates evidence with some temporal integration bias, $\alpha$ (where $\alpha > 1$ creates a primacy effect, and

11 $\alpha < 1$, a recency effect), and incurs inference error (noise in the estimate of the true evidence)

12 parameterised by $\sigma$, the standard deviation of the Gaussian distribution from which each sample of noise, $\varepsilon_n$,

13 is drawn from. The human observer may also experience some response bias, $c$ (the tendency to choose one

14 category irrespective of the evidence), and incur lapses (pressing a random key), described by the lapse rate,

15 $l$. The accumulated evidence, $z$, up to sample $n$, is suboptimally accumulated by

$$z_n = \alpha z_{n-1} + \ell_n + \varepsilon_n \tag{2}$$

16 The observer then chooses category 1 if $z > c$, except on a proportion of trials, $l$, where the response is

17 randomly selected.

18 These four parameters were used to capture the differences in the human observers' responses (category

19 choice and confidence rating) from the optimal observer who perfectly integrates all evidence presented.

20 In the Free task, the model was designed not only to describe the category choice, but at which sample the

21 human observer chose to respond. This was achieved via a decision boundary, the nature of which has been

22 addressed in previous work (Balsdon, Wyart, and Mamassian, 2020). The boundaries, $\Lambda_{n+}$ and $\Lambda_{n-}$, follow

23 an exponential function on the average evidence over samples (which is a constant bound on the probability

24 of a correct response), described by three parameters: the minimum, $a$, the scale, $b$, and the rate of decline, $\lambda$

$$\Lambda_{n+} = n \times \left(a + be^{-\frac{n}{\lambda}}\right) \tag{3}$$

25 There is an optimal combination of these parameters to achieve any particular proportion correct across the

26 experiment, but the human observer may set their bound suboptimally. In addition, non-decision time (the

27    time from the last sample integrated to pressing the response key) was described by a Normal distribution

28    with mean $\mu_U$, and variance $\sigma_U^2$. Giving an additional five parameters for describing when the observer enters

29    their response.

30    We followed the same procedure as in Balsdon et al., 2020, involving four stages:

31    1. Reduce the number of free parameters with a knock-out procedure.
32    2. Compare (covert) Bound and No-bound models of the perceptual decision in the Replay task.
33    3. Identify any systematic differences in the parameters required to describe the confidence ratings,
34      compared to the perceptual decision, in order to discern the relationship between processes for
35      perceptual decisions and confidence.
36    4. Apply the same Bound vs. No-bound comparison for describing the confidence ratings.

37    The average parameter values and fit metrics for Stage 1. are shown in Table 1. According to this analysis,

38    the bias (c) and lapse rate (l) were fixed. There was some evidence the boundary minimum (a) could be fixed

39    in the Replay task, but the preference in the Free task was to leave it free to vary.

| Free Task | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Model** | $\sigma$ | $\alpha$ | $c$ | $\mu_U$ | $\sigma_U^2$ | $a$ | $b$ | $\lambda$ | $l$ | **LLH** | **$\Sigma$BIC** |
| Full | 0.83 | 0.98 | -0.04 | 425 | 0.52 | 0.10 | 6.04 | 1.93 | 0.016 | -734.91 | 30423.01 |
| $\alpha = 1$ | 0.83 | 1.00 | 0.00 | 430 | 0.50 | 0.13 | 6.61 | 2.03 | 0.014 | -734.97 | 30311.59 |
| $c = 0$ | 0.80 | 0.92 | 0.00 | 452 | 0.54 | 0.11 | 5.28 | 2.01 | 0.017 | -736.86 | 30387.02 |
| $\mu_U = 400$ | 0.76 | 0.94 | 0.00 | 400 | 0.52 | 0.09 | 5.52 | 2.23 | 0.016 | -739.77 | 30503.40 |
| $\sigma_U^2 = 1$ | 0.69 | 0.96 | -0.02 | 435 | 1.00 | 0.10 | 6.34 | 1.97 | 0.015 | -754.18 | 31079.84 |
| $a = 0.1$ | 0.77 | 0.92 | 0.03 | 417 | 0.52 | 0.10 | 5.78 | 2.20 | 0.016 | -735.48 | 30331.75 |
| $b = 5.5$ | 0.78 | 0.94 | 0.02 | 410 | 0.64 | 0.13 | 5.50 | 1.79 | 0.013 | -742.18 | 30599.67 |
| $l = 0.001$ | 0.82 | 0.98 | 0.01 | 400 | 0.48 | 0.10 | 4.77 | 2.22 | 0.001 | -730.66 | 30139.17 |
| $c = 0; l = 0.001$ | 0.79 | 0.94 | 0.00 | 397 | 0.51 | 0.10 | 4.52 | 2.26 | 0.001 | -732.66 | 30104.74 |
| $c = 0; l = 0.001; a = 0.1$ | 0.77 | 0.94 | 0.00 | 403 | 0.52 | 0.10 | 5.37 | 2.13 | 0.001 | -742.42 | 30381.13 |
| **Replay Task - no-bound** | | | | | | | | | | | |
| **Model** | $\sigma$ | $\alpha$ | $c$ | $\mu_U$ | $\sigma_U^2$ | $a$ | $b$ | $\lambda$ | $l$ | **LLH** | **$\Sigma$BIC** |
| Full | 0.47 | 0.90 | 0.05 | ~ | ~ | ~ | ~ | ~ | 0.012 | -81.13 | 3701.44 |
| $\alpha = 1$ | 0.56 | 1.00 | 0.10 | ~ | ~ | ~ | ~ | ~ | 0.012 | -92.21 | 4030.55 |
| $c = 0$ | 0.48 | 0.90 | 0.00 | ~ | ~ | ~ | ~ | ~ | 0.009 | -82.73 | 3651.38 |
| $l = 0.001$ | 0.50 | 0.91 | 0.06 | ~ | ~ | ~ | ~ | ~ | 0.001 | -82.05 | 3624.39 |
| $c = 0; l = 0.001$ | 0.51 | 0.90 | 0.00 | ~ | ~ | ~ | ~ | ~ | 0.001 | -83.64 | 3573.67 |
| **Replay task - bound** | | | | | | | | | | | |
| **Model** | $\sigma$ | $\alpha$ | $c$ | $\mu_U$ | $\sigma_U^2$ | $a$ | $b$ | $\lambda$ | $l$ | **LLH** | **$\Sigma$BIC** |
| Full | 0.44 | 0.87 | 0.10 | ~ | ~ | 0.17 | 8.68 | 11.71 | 0.012 | -79.81 | 3991.09 |
| $c = 0; l = 0.001$ | 0.48 | 0.88 | 0.00 | ~ | ~ | 0.13 | 8.58 | 15.55 | 0.001 | -82.22 | 3859.24 |
| $c = 0; l = 0.001; a = 0.1$ | 0.48 | 0.88 | 0.00 | | | 0.10 | 8.91 | 15.88 | 0.001 | -82.38 | 3751.55 |

40    ***Table S1. Average parameter values.*** *Table shows the average values and the sum of BIC across participants.*

41    *The large difference in the average loglikelihood (LLH) across tasks is due to the fact the Free task model was fit*

42    *to both when and what observers responded, whereas in the Replay task only the response was fit. Red values*

43    *show the fixed parameters. Colour code of the BIC column corresponds to the goodness of fit (the greener the*

44    *better).*

45     To compare the Bound and No-bound models in Stage 2. we used five-fold cross validation. The No-bound

46     model had two free parameters: $\alpha$ (temporal bias) and $\sigma$ (inference noise), which were fit to the Same and

47     Less conditions of the Replay task, but tested across all conditions. The Bound model had three free

48     parameters to describe the bound, with the inference noise and temporal bias parameters fixed to those fit

49     to the Same and Less conditions only. In this way, the no-bound model must account for the lack of increased

50     performance in the More condition with the suboptimalities present in the Same and Less conditions, whilst

51     the bound model can limit performance in the More condition in particular by stopping further evidence

52     accumulation. The results of this analysis are presented in the manuscript: the bound significantly improved

53     the fit, mean relative increase in model log-likelihood = 0.048, bootstrapped $p$ = 0.001, **Figure 2c** in the main
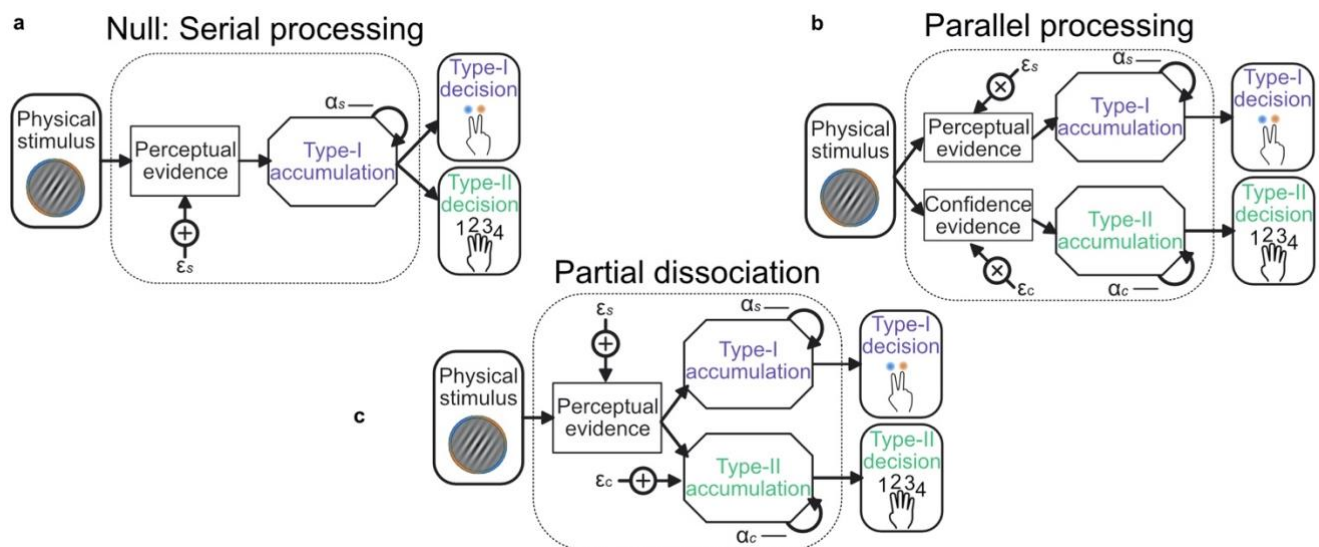
54     text.

55     Of additional interest is the pattern of parameters fit to each condition separately, when the model attempts

56     to explain behaviour without a bound. There was little difference in parameters fit to the Same and Less

57     conditions (mean $\sigma_S = 0.48$, $\sigma_L = 0.44$, $Z(19)$ = -1.46, $p$ = 0.15; $\alpha_S = 0.86$, $\alpha_L = 0.78$, $Z(19)$ = 1.38, $p$ = 0.17).

58     The inference noise fit to the More condition significantly increased from the Less condition ($\sigma_M = 0.55$,

59     $Z(19)$ = -2.61, $p_{bonf*4}$ = 0.036), but there was significantly reduced temporal integration bias ($\alpha_M = 0.93$,

60     $Z(19)$ = -2.50, $p_{bonf*4}$ = 0. 0496) suggesting observers' performance was worse than predicted by the Same

61     and Less conditions, and they were putting less weight on the more recent evidence. These differences in

62     parameters are consistent with the model trying to mimic bounded evidence accumulation without a bound,

63     providing additional support for the comparison described above.

64     Stage 3. of the model procedure was to account for the confidence ratings. We compared three processing

65     architectures that span the space from single-channel to dual-channel (Maniscalco and Lau, 2016). We took

66     as the null hypothesis a serial processing (single-channel) architecture in which the confidence ratings

67     (Type-II decisions) can be described by the exact same evidence as used to make the perceptual (Type-I)

68     decision. A weaker version of this null hypothesis is that the same suboptimal inference process is used for

69     both perception and confidence, but that the observer can commit to their perceptual decision whilst

70     continuing to monitor additional evidence for evaluating their confidence (a schematic of these processes is

71     shown in **Figure S1a**). The average parameter values are shown in **Table S2**, labelled 'Serial' and 'Serial

72     continued' respectively. Note the substantial increase in inference noise ($\sigma$) and reduction in temporal bias

73     ($\alpha$ is closer to 1) when attempting to describe both the perceptual decision and the confidence rating

74     compared to only the perceptual decision (**Table S1**, Replay task – bound, model *c = 0; l = 0.001*). This is

75     indicative of the difficulty of describing both perception and confidence with the same suboptimalities.

76     At the other extreme is the parallel processing (dual-channel) architecture, in which perception and

77     confidence are computed by independent resources, based on the same sensory input (**Figure S1b**, labelled

78     'Parallel' in **Table S2**). This is the most computationally expensive description, and provided a lack of

79     parsimony that was only surpassed by a model that attempted to describe confidence ratings with only the

80     inference noise evident from the perceptual decisions.

81    The intermediate models in this architectural space are the partial dissociation models (**Figure S1c**), which

82    suggest that confidence inherits the same noisy perceptual evidence as the perceptual decision, but may

83    incur some independent suboptimalities. We compared four versions of these models: same $\sigma$ (no additional

84    inference noise); accumulation noise (additional inference noise with each sample of evidence); read-out

85    noise (one additional sample of noise before the confidence response); and same $\alpha$ (the temporal bias

86    affecting the confidence accumulation is the same as that affecting the perceptual accumulation).

87    In all cases the models were fit to minimise the negative loglikelihood of both perceptual and confidence

88    decisions. The model comparison overwhelmingly favoured the partial dissociation models, and of these, the

89    best description was offered by a model with an independent temporal bias on the confidence evidence

90    accumulation, and additional noise at the read-out stage. We caution against interpreting this result as

91    meaning that there is no additional accumulation noise in the processing of confidence evidence, whilst the

92    models are very similar, it is possible that the read-out noise in this case can additionally capture some noise

93    in setting and maintaining bounds for assigning a rating to the confidence evidence.

94



95    *Figure S1. Schematic of possible relationships between perceptual (Type-I) and confidence (Type-II)*

96    *evidence accumulation. a) Same evidence accumulation processes: Type-I (perceptual) and Type-II*

97    *(confidence) decisions are different responses to the same evidence: each sample of perceptual evidence is*

98    *disrupted by a sample of sensory noise ($\varepsilon_s$) drawn from a zero-mean Gaussian with standard deviation $\sigma$, and*

99    *accumulated with a temporal bias described by $\alpha_s$. b) Parallel processing: Type-I and Type-II decisions rely on*

100    *entirely separate processing of the same physical stimulus: the confidence decision also incurs noise and*

101    *temporal integration bias (with subscript c), but these may vary independently of the perceptual processing*

102    *suboptimalities (subscript s). c) Partial dissociation: Type-I and Type-II decisions rely on partially dissociable*

103    *accumulation of the same evidence.*

104

105

| Model | $\sigma$ | $\alpha$ | $a$ | $b$ | $\lambda$ | $a_c$ | $b_c$ | $\lambda_{c1}$ | $\lambda_{c2}$ | $\lambda_{c3}$ | LLH | ΣBIC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Serial | 0.73 | 0.92 | 0.10 | 12.74 | 17.07 | 0.07 | 0.64 | 1.28 | 6.81 | 31.38 | -428.36 | 18275.28 |
| Serial continued | 0.67 | 0.91 | 0.13 | 9.60 | 17.98 | 0.06 | 0.53 | 0.66 | 7.08 | 30.41 | -424.88 | 18135.83 |
| Parallel | 0.76 | 0.90 | ~ | ~ | ~ | 0.01 | 0.58 | 0.18 | 7.51 | 30.68 | -437.25 | 18288.50 |
| Partial - same sigma | 0.00 | 0.89 | ~ | ~ | ~ | 0.06 | 0.47 | 1.03 | 6.77 | 25.92 | -446.41 | 18540.68 |
| Partial - accumulation noise | 0.45 | 0.91 | ~ | ~ | ~ | 0.03 | 0.58 | 0.50 | 7.71 | 31.03 | -421.59 | 17662.25 |
| Partial - read-out noise | 0.12 | 0.90 | ~ | ~ | ~ | 0.02 | 0.52 | 1.85 | 8.63 | 37.39 | -417.94 | 17516.29 |
| Partial - same alpha | 0.12 | 0.88 | ~ | ~ | ~ | 0.02 | 0.52 | 0.98 | 8.22 | 35.16 | -423.02 | 17605.29 |

106

***Table S2. Average parameter values for perceptual and confidence behaviour.*** *Bound parameters with subscript c describe the criteria for confidence ratings, which take the same form as the perceptual decision bound. They have the same minimum and scale, but different rates of decline, such that $\lambda_{c1}$ determines the upper bound on a confidence rating of 1, and the lower bound on a rating of 2. Apart from the 'Serial' and 'Serial continued' models, parameters for perceptual decisions were fixed to those fit in the winning perceptual decision model and the listed parameters affect only the confidence evidence accumulation.*

113 The model comparison of Stage 3. just described mainly assumed continued, unbounded accumulation of confidence evidence (with the exception of the strictly serial processing architecture). Stage 4. was to formally compare bounded and unbounded accumulation for confidence evaluations in the same manner as with the perceptual decisions. This time, two versions of the bound were compared: the same bound as perceptual evidence accumulation (the participant could close their eyes after committing to their perceptual decisions and their responses would not change); or an independent bound (the participant can continue to accumulate evidence for confidence decisions after the committing to the perceptual decision, but will eventually stop). As reported in the manuscript, neither bound improved the fit, if anything, adding the bound decreased the log-likelihood of the model (same bound: relative improvement with bound = -0.007, bootstrapped $p$ = 0.11, uncorrected; independent bound: relative improvement = -0.014, $p$ = 0.022, Bonferroni corrected for two comparisons; **Figure 2c,** in the main text). This reflects the fact that even a very high bound affects the shape of the accumulation trace, which will harm the fit when behaviour is not affected by a bound.

126 In summary, this computational modelling procedure suggests a partial dissociation in the processing for perception and confidence. In the Replay task, perceptual decisions were best described by bounded evidence accumulation, enabling observers to commit to decisions before the sequence of presented samples finishes. The confidence ratings required additional noise and reduced temporal integration bias compared to the suboptimalities affected the perceptual decisions. These differences were best described by the partial dissociation architecture where confidence received the same noise samples of evidence as the perceptual decision, though they are accumulated differently. In addition, model comparison suggested confidence evidence accumulation continued to the end of the sequence, even in cases of premature commitment to the perceptual decision. The results of these comparisons replicate the results of Balsdon et al. (2020), with the exception of the confidence noise comparison: here we find evidence in favour of read-out noise, whereas the previous analysis found the models indistinguishable.

## Supplementary Note 2

**Model Simulation**

The computational model comparison suggested a partial dissociation in the evidence used to make perceptual decisions and confidence evaluations. We compared the evidence underlying the observers' perceptual decisions and confidence ratings by simulating the winning computational model. For each trial, 10,000 samples of noise per decision update were randomly sampled from the Gaussian distribution describing the observer's inference noise. These were combined to give 10,000 simulated evidence traces per trial. The first 1,000 simulated evidence traces that agreed with the observer's response on that trial were taken to measure the median evidence trace (or, the process was repeated until 1,000 adequate simulated evidence traces were drawn, up to 100 repeats). **Figure S2a** demonstrates this process for one example trial of one observer. For the perceptual evidence (**Figure S2a**, left) simulated evidence traces that agreed with the observer's response are those that reach the respective decision bound before the opposing decision bound, or reach no bound but show evidence in favour of the response by the final sample. It was assumed that once the evidence reaches the bound, that evidence is maintained until the response. For the confidence evaluation (in the example, a confidence rating of 3), the final evidence had to be between the confidence rating bounds to agree with the observer's confidence decision (after the final sample of additional noise – which is why a few samples in **Figure S2a**, right, exceed the bounds). The median evidence was compared to the ideal evidence (green lines of **Figure S2a**).

The estimated inference error (used in **Supplementary Note 7**) scaled the difference between the median consistent evidence and the ideal evidence by the probability of the response given all samples, to estimate the relative deviation of the observers' internal evidence from the optimal observer's evidence. This estimate of the error is quite imprecise: the median trace tends to be quite close to the ideal, even though any one of the traces (which reflect much larger error) could have described the internal evidence of the observer. **Figure S2b** shows the predicted final accumulated evidence for the perceptual (Type-I) compared to the confidence (Type-II) decision for the same example observer. The evidence is strongly correlated but there are substantial deviations, because of the additional noise, different temporal bias, and continued accumulation for the confidence decision, especially in the More condition (light blue). The example observer is a more extreme case because of the relatively strong bound on perceptual evidence accumulation. The (Fisher transformed) correlation for each observer is shown in **Figure S2c**. For many observers there are substantial differences between the median simulated evidence consistent with the perceptual and confidence responses, meaning the simulated evidence could be useful in distinguishing representations important for perception vs. confidence.
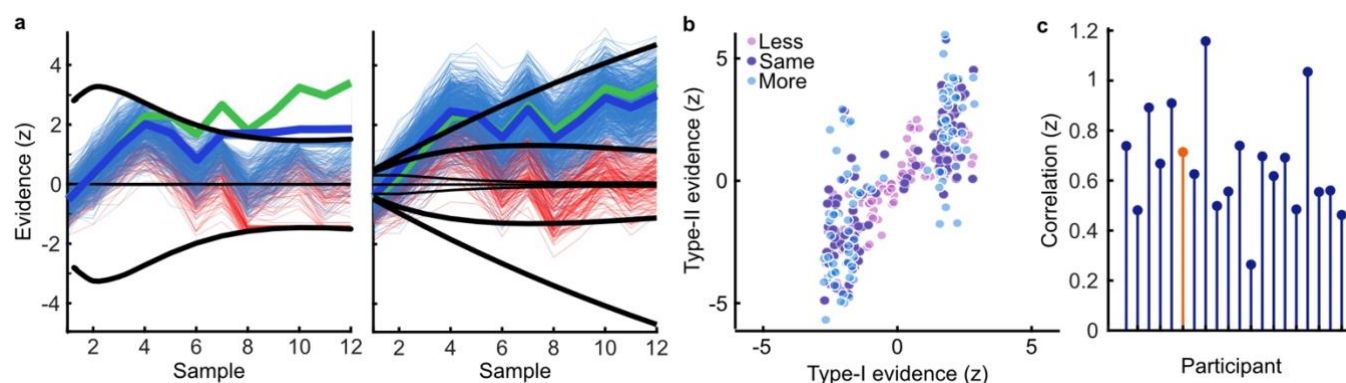
6

169

*Figure S2. Model simulation of accumulated evidence for perceptual and confidence decisions. a)*

*Example trial from one observer showing simulated evidence traces agreeing with the observer's response*

*(blue) and a sample of example traces which did not agree (red). The perceptual decision is shown on the left.*

*An evidence trace was taken to agree with the observer's decision if the corresponding bound was reached prior*

*to the opposing bound, or if no bound was reached but the final accumulated evidence was in favour of the*

*chosen option. The median evidence trace (thick blue line) was calculated assuming the evidence that reached*

*the bound early was maintained until the response was entered. For the confidence rating (right) we compared*

*the median evidence from traces where the final accumulator (plus one additional sample of noise) agreed with*

*the observer's confidence rating. We examined the difference from the ideal accumulated evidence (thick green*

*line) relative to the likelihood of the observers' rating given all simulated evidence traces. **b)** Median final*

*simulated accumulated evidence for the perceptual decision (abscissa), and the confidence decision (ordinate)*

*for all trials of the example observer, colours indicate the condition. **c)** Correlation (Fisher transformed z)*

*between perceptual and confidence evidence for each observer. The example observer is highlighted in orange.*

## Supplementary Note 3

### Confidence behaviour

Proportion correct increased with increasing confidence, reflecting the observers' ability to use their

confidence ratings to discriminate correct from incorrect responses (**Figure S3a**). Observers appeared to be

monitoring the decision evidence to make their confidence ratings, as opposed to some proxy for confidence

such as the number of samples they were shown (**Figures S3b** and **S3c**).

We required a single-trial measure of confidence precision for identifying the key neural processes

underlying the computation of confidence. To do so, we compared observers' responses to an optimal

observer. The optimal observer perfectly accumulates all presented evidence and assigns ratings to equally

partition the evidence for their perceptual decision. To simplify, we split trials by the median evidence for

the chosen category, where the optimal observer gives a high confidence rating (3 or 4) to those trials with

greater than the median evidence, and a low confidence rating (1 or 2) to those with less than the median

evidence. We labelled trials as 'suboptimal confidence' when the observer's confidence response disagreed

with the response of this optimal observer. This trial labelling is demonstrated for two example observers in

**Figure S3d**. We reasoned that on suboptimal confidence trials the internal evidence of the human observer

7

198 was less likely to be close to the optimal presented evidence, and the neural representation of the optimal

199 presented evidence should be less precise in neural circuits that actually represent this suboptimal

200 confidence evidence. That this measure of confidence precision does capture the suboptimalities in

201 confidence evaluation is confirmed by the significant increase in model estimated confidence error on

202 suboptimal confidence trials (Wilcoxon sign rank test: $Z(19) = 3.85$, $p < 0.001$; **Figure S3e**).

203 In this way, observers' confidence is assessed relative to a "super-ideal" observer, who has perfect access to

204 the presented evidence (Mamassian and de Gardelle, under review). Theoretically, observers' confidence

205 should be assessed relative to the internal evidence for their perceptual decision, that is, relative to the

206 evidence based on suboptimal inference (afflicted by noise and temporal integration biases). However, the

207 single-trial estimates of the internal evidence for perceptual decisions, based on model simulations, were

208 relatively imprecise (see **Supplementary Note 2**), and could also introduce systematic errors from the

209 model assumptions, making this estimate of the internal evidence unappealing for the purpose of assessing

210 confidence. Moreover, the goal of this measure was to compare observers' confidence ratings to the neural

211 representation of the accumulated evidence, which was also assessed relative to the optimal evidence. We

212 therefore chose to assess confidence ratings relative to the optimal observer in the same way that neural

213 responses were assessed relative to optimal, though this ignores the fact that some suboptimality is actually

214 inherited from perceptual decision processes.

215 A second important consideration with this measure is that it is affected by confidence bias. There are three

216 types of biases that could affect confidence ratings: first, a response bias to enter a certain response

217 irrespective of the evidence; second, a miscalibration bias such that ratings mean different things to different

218 observers (the same value of evidence will be given a rating of 4 for one observer and 3 for another, for

219 example); third, a miscalling bias such that perceptual evidence is relatively exaggerated or diminished in

220 the assessment of confidence. All these biases mean that the same internal perceptual evidence could result

221 in systematically different confidence ratings across observers, and observers could report on average

222 higher or lower confidence despite similar perceptual performance and precision in representing the

223 internal evidence for evaluating their confidence.

224 Taking an average proportion of suboptimal confidence ratings and comparing across observers would

225 result in observers of similar ability having different scores simply because of biases in how they implement

226 the confidence rating responses: greater biases will increase average proportion suboptimal. Importantly,

227 this single-trial measure of confidence was not used for this purpose. Rather, it was compared to neural

228 activity during the process of accumulating evidence for the perceptual decision and confidence evaluation.

229 We expect that biases that are not of interest for the computation of confidence (in particular, response bias

230 and miscalibration bias) are incorporated at a later stage, when the confidence evaluation is converted into a

231 rating for executing the response. The biases will only reduce the sensitivity with which a trial labelled as

232 suboptimal truly reflects internal evidence that differs from optimal, reducing our ability to identify neural

233 processes underlying confidence computation. This is simulated in **Figure S3f**, where a relative bias is

234 introduced by assessing human confidence ratings to a biased optimal observer (who responds on 65% of

8

235 trials with high confidence – making the human observers relatively more liberal, or 35% high confidence –

236 making the human observers more conservative). The general trend for the difference between confidence

237 ratings that match the (biased) optimal observer and those that are suboptimal remains the same, though
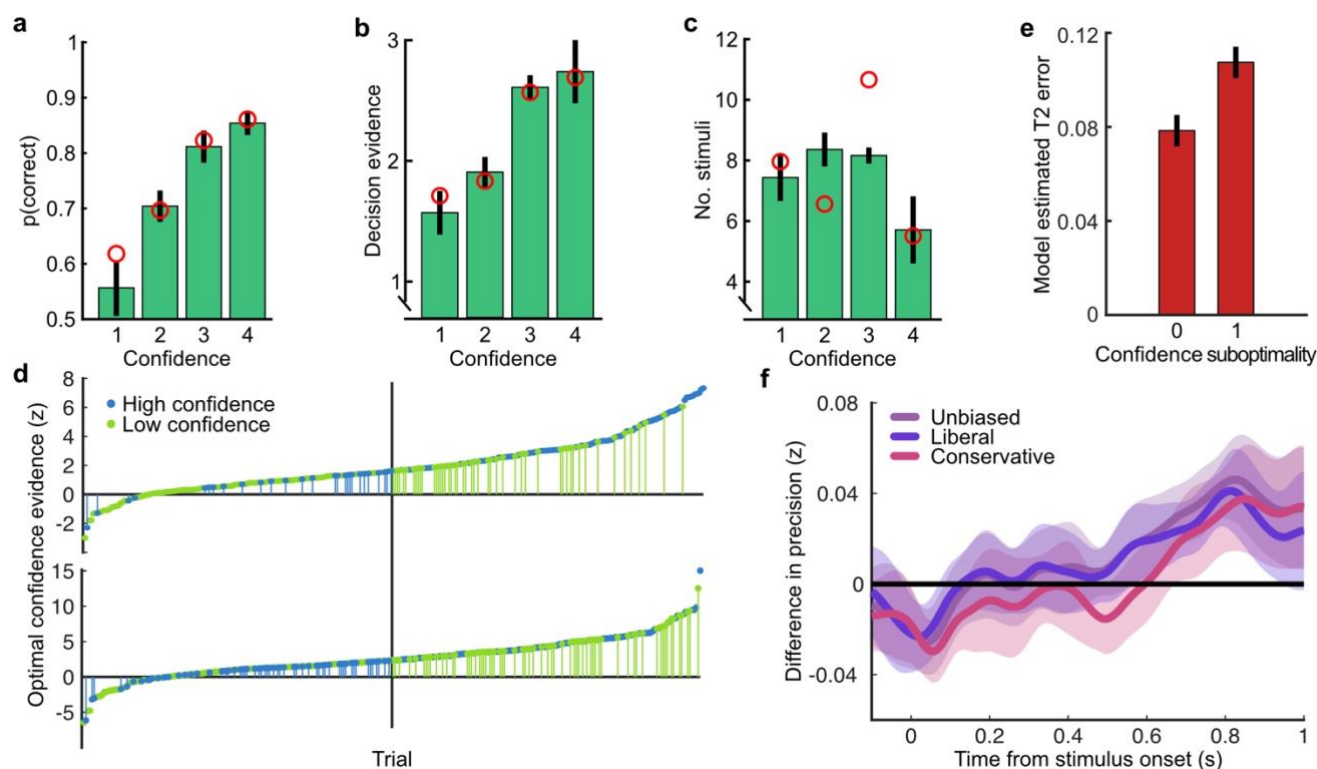
238 the bias reduces the difference.



239

240 **Figure S3. Confidence behaviour**. *a) proportion correct (in the perceptual decision) by confidence rating. b)*

241 *Decision evidence (based on the presented samples) by confidence rating. c) Number of samples presented by*

242 *confidence rating. In all plots, error bars show 95% within-subject confidence intervals. Red circles show the*

243 *predictions of the best fitting confidence model (**Supplementary Note 1**). d) Confidence responses of two*

244 *observers (top and bottom panels) on all trials sorted by the confidence evidence of the optimal observer. The*

245 *median confidence evidence (shown by a black vertical line) defines an optimal confidence observer whose*

246 *confidence above this median are rated high. Observers' high confidence ratings are shown in blue and low*

247 *confidence ratings in green. Suboptimal confidence ratings, where human and optimal confidence observers do*

248 *not match, are indicated with small vertical segments (green for Type-II misses and blue for Type-II false*

249 *alarms). Negative confidence evidence corresponds to incorrect perceptual decisions. The observer shown on*

250 *top clearly has fewer suboptimal responses compared with the observer below, and the frequency of suboptimal*

251 *responses decreases further from the median. e) Model estimated confidence error by confidence rating*

252 *suboptimality (0 = the observer's confidence rating was the same as the optimal observer, 1 = suboptimal*

253 *confidence rating). f) The effect of response bias on the analysis of suboptimal confidence in the EEG*

254 *representation of accumulated evidence. Observers' confidence ratings were compared to an unbiased optimal*

255 *observer (purple), and two biased (but otherwise optimal) observers, who respond with high confidence on 35%*

256 *and 65% of trials (making the human observers relatively more liberal and conservative with their response*

257 *strategy in comparison). Thick lines show the within-subject difference in precision (Fisher transformed*

258    *correlation) between trials where the human observers' confidence ratings correspond to the (un/biased)*

259    *optimal observer and suboptimal confidence ratings. Shaded regions show the 95% between-subject confidence*

260    *intervals on the difference.*

261    # Supplementary Note 4

262    ## Response classification

263    A linear discriminant analysis was used to classify the perceptual decision response based on the spectral

264    power of band-limited EEG signals in epochs locked to the time of the response. The spectral power across

265    frequency tapers from 1 to 64 Hz with 25% spectral smoothing was resolved using wavelet convolution

266    implemented in FieldTrip (Oostenveld et al., 2011). The epochs were then clipped at -3 to 1 s around the

267    time of entering the perceptual decision response. We first trained and tested at each frequency taper at

268    each time point in the Free task (**Figure S4a**). Classifier performance was measured as the area under the

269    curve (AUC). The power in frequency bands between 8 and 32 Hz yielded the most accurate classification

270    performance. The difference in the average power across these frequency bands between -0.5 and 0.5

271    seconds around the time of the response for right- and left-handed responses showed a clear lateralisation

272    over central and parietal electrodes (**Figure S4b**). Training and testing at each time point in each condition

273    of the Replay task showed a similar pattern to the Free task, with reliable classifier performance from

274    around -0.5 to 0.5 seconds around the response (**Figure S4c**). Training and testing within each condition of

275    the Replay task resulted in a larger between-subject error, likely because there are only 100 trials per

276    condition. In the main text, we present a cross-classification analysis where the classifier is trained on the

277    Free task, and tested on each condition in the Replay task, which more directly examines when the signals

278    relevant for entering a response (based on the Free task) emerge during the lead up to the response in each
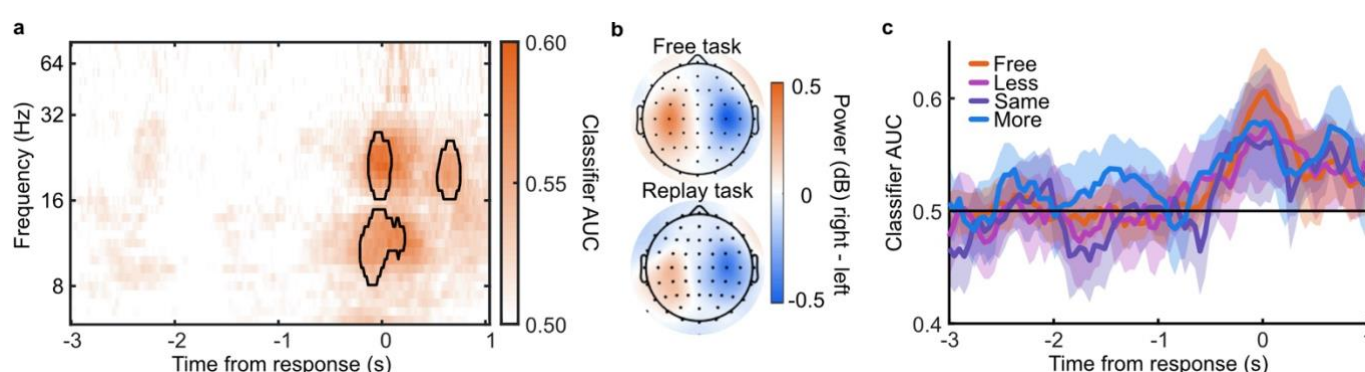
279    condition of the Replay task.

280



281    ***Figure S4. Response Classification analysis. a)*** *Classifier AUC training and testing at each time point*

282    *(abscissa) based on the power (dB) in each frequency band (ordinate). Clusters where average performance is*

283    *greater than 3.1 standard deviations (99% confidence) from baseline (0.5) are circled in black.* ***b)*** *Scalp map of*

284    *the difference in power for right- compared to left-handed responses averaged over 8 to 32 Hz and -0.5 to 05*

285    *seconds around the response.* ***c)*** *Classifier performance (AUC) training and testing at each time point, in each*

286    *condition of the Replay task and in the Free task.*

## Supplementary Note 5

**Encoding variable regression**

Linear regression was used to examine the representation of encoding variables in the EEG signals. First, regression weights ($\widehat{W}$) were computed using ridge regression of the encoding variables ($C$, an $n \times 1$ matrix) on the EEG signals ($D$, an $n \times m$ matrix, where $m$ is the number of EEG signals, and $n$, the number of epochs):

$$\widehat{W} = (D^T D + \lambda I)^{-1} D^T C \qquad \qquad \textbf{(4)}$$

The regularisation parameter, $\lambda$, was set to 1, where $I$ is the identity matrix. Weights were computed on 90% of the epochs, and used to predict the encoding variables on the other 10% (10-fold cross validation) simply as: $\hat{C} = D * \widehat{W}$. The precision of the prediction was calculated as the correlation between $\hat{C}$ and $C$, standardised using a Fisher transformation.

Three different encoding variables, $C_\theta$, $C_\ell$, and $C_z$, were examined (**Figure S5a**): the stimulus orientation ($C_\theta = \pi - |\theta_n|$), the momentary decision update ($C_\ell = |\ell_n| = \left| \kappa cos\big(2(\theta_n - \mu_1)\big) - \kappa cos\big(2(\theta_n - \mu_2)\big) \right|$), and the accumulated evidence ($C_z = z_n = \sum_{N=1}^{n} \ell_N$ , signed by the response). These variables are not entirely independent: There is a weak correlation between the stimulus orientation and the momentary decision update ($r = 0.03$), and a weak correlation between the momentary decision update and the accumulated evidence ($r = 0.09$). In addition, the accumulated evidence is strongly correlated over samples ($r = 0.92$ at n+1, and $r = 0.85$ at n+2). The cross-correlations are shown in **Figure S5c**.

The EEG signals in D were low-pass filtered and decomposed into real and imaginary parts using a Hilbert transform. Regression precision was first calculated using the signals from all electrodes ($m = 128$) separately for each time-point in the stimulus-locked epochs. Initial analysis showed a low-pass cut-off of 8 Hz was appropriate to decrease noise whilst maintaining precision (**Figure S5b**). The previous literature has shown similar results (Salvador et al., 2020).

Temporal generalisation of the representation of encoding variables was tested by computing weights at each time point and testing the predicted encoding variables across time (**Figure S5d**). Though the representation of the momentary decision update is maintained for a relatively longer duration than the representation of stimulus orientation, there is little temporal generalisation, suggesting the representation in the EEG signals evolves over time. This is also the case for the representation of accumulated evidence, however, there are also strong off-diagonals in the temporal generalisation matrix. This is likely because of the strong correlation across consecutive samples (**Figure S5c**).

The precision of the representation of accumulated evidence was compared across the Less and More conditions for the first four and the last four stimuli (**Figure S5e**). As reported in the main text, representation precision was substantially attenuated for the last four stimuli of the More condition. This was not the case for the first four samples, where decoding precision in the More condition was briefly (from

319    132 to 244 ms) greater than in the Less condition ($t_{ave}(19)$ = 3.67, $p_{cluster}$ < 0.001).

320    Given the sustained precision of decoding accumulated evidence over time, and the strong correlation

321    between consecutive samples, it is curious that the measured precision does drop to baseline at the start of

322    the epoch. That the same pattern is found when decoding sample n-1 and sample n+1 based on the epoch at

323    sample n (**Figure S5f**) suggests that the onset of the stimulus is disrupting the ongoing representation (or at

324    least, our ability to measure it). Furthermore, this decrease in performance is not seen in the temporal

325    generalisation matrix, where the off-diagonal is not aligned with the onset of successive samples (due to the

326    jitter in stimulus presentation timing). Comparing precision between groups of epochs where the timing of

327    the subsequent sample is aligned (**Figure S5g**; red 317 ms, green 333 ms, blue 350 ms) suggests there could

328    be an interaction between the timing of ongoing updates and the precision of the representation of the

329    accumulated evidence (but not the momentary decision update). This could be of interest for future

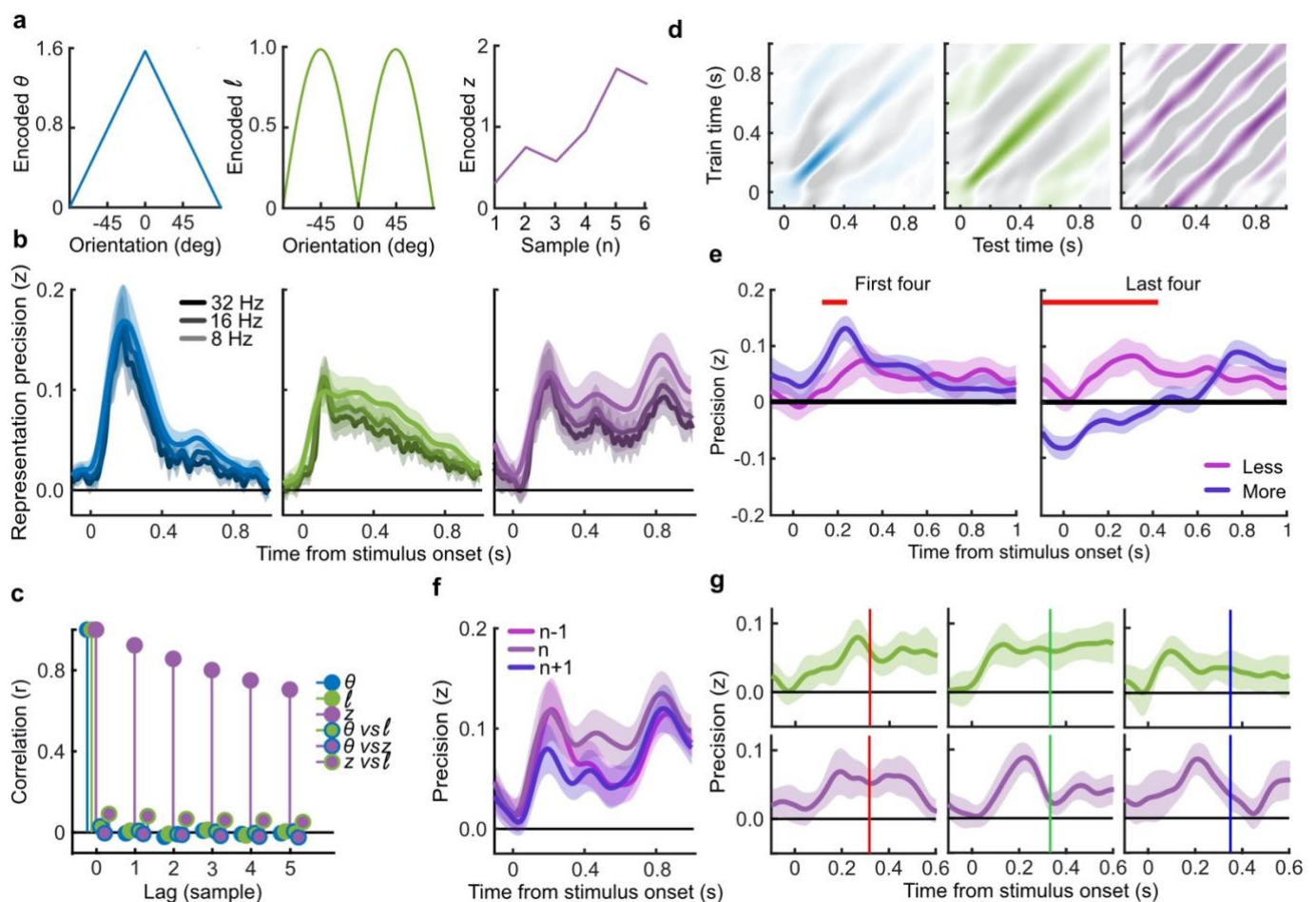330    research.



332    ***Figure S5. Encoding variable regression. a)** Encoded variables used to regress EEG signals. The encoded*

333    *orientation ($C_\theta$, left) and encoded momentary decision update ($C_\ell$, middle) were dependent on the orientation*

334    *presented to the observer. The encoded accumulated evidence ($C_z$) varied over all presented orientations in a*

335    *trial, the figure on the right shows only one example. **b)** Representation precision of encoding variables using*

336    *different low-pass filters. **c)** Cross correlation between encoding variables over consecutive samples. **d)***

337    *Temporal generalisation of representations: the regression weights were calculated on EEG signals at each time*

338 *point and precision was tested across time. Colour scales are relative to the maximal precision, with zero*

339 *precision in white and negative in grey (a sign flip of the regression weights).* ***e)*** *Representation precision of the*

340 *accumulated evidence for the first (left) and last (right) four stimuli of the Less and More conditions. Shaded*

341 *error bars show the 95% within subject confidence intervals, red horizontal bars mark cluster corrected*

342 *significant differences between conditions.* ***f)*** *Representation precision of the previous (n-1), current (n) and*

343 *future (n+1) accumulated evidence, based on the EEG signals locked to the current epoch.* ***g)*** *Representation*

344 *precision of the momentary decision update (top) and the accumulated evidence (bottom) for epochs separated*

345 *by the timing of the subsequent stimulus, shown in coloured bars (317 ms, red, left; 333 ms, green, middle; and*

346 *350 ms blue, right).*

## Supplementary Note 6

### Cluster modelling

349 Cluster modelling was used to isolate contiguous signals in space (electrode location) and time, where the

350 precision of the representation of accumulated evidence systematically varied with the suboptimalities

351 evident from behavioural responses. Suboptimal responses result from greater inference error, where the

352 internal representation of the accumulated evidence deviates further from the presented evidence, thus

353 neural signals that reflect the internal evidence of the observer should also deviate further from the optimal

354 evidence used in the regression. Clusters were isolated using a multivariate Bayesian scan statistic (Neill,

355 2011; Neill, 2019). This statistic was calculated based on the loglikelihood ratio of the alternative hypothesis

356 (that representation precision depends on the inference noise of the observer) against the null hypothesis

357 (that any difference in representation precision is due to measurement noise alone, which is independent

358 across epochs). It is assumed that the neural signals reflect the input (cumulative presented evidence) with

359 added measurement noise ($N_m$) and, when the neural signals are relevant for behaviour, inference noise ($N_i$)

360 that reflects the observers' suboptimal internal representation of the decision evidence:

$$Y_{out} = Y_{in} + N_i + N_m \qquad\qquad \textbf{(5)}$$

361 Where the two sources of noise are assumed to be gaussian distributed ($N(0, \sigma^2)$). The total measured

362 correlation ($r_T$) between $Y_{in}$ and $Y_{out}$ is a function of the additional noise (where $Y_{in}$ is normalised):

$$r_T = \frac{1}{\sqrt{2 + \sigma_i^2 + \sigma_m^2}} \qquad\qquad \textbf{(6)}$$

363 The observer makes suboptimal decisions when the inference noise pushes their internal representation of

364 the accumulated evidence further from the true value, resulting in a weaker correlation between the internal

365 representation and the presented evidence. Therefore, when we split based on behaviour, we expect that on

366 average there is greater inference noise on incorrect trials than correct trials. The correlation over all

367 samples can be described as:

13

$$r_T = \frac{1}{\sqrt{2 + p(I)\sigma_{iI}^2 + p(C)\sigma_{iC}^2 + \sigma_m^2}} \tag{7}$$

368 Where $p(I)$ is the observed probability of a suboptimal decision, and $p(C)$, a decision that corresponds to that
369 of the optimal observer. The null hypothesis is that the neural signal is not relevant for behaviour,
370 specifically, signals on suboptimal trials do not reflect additional inference noise. Any difference in the
371 correlation is due to variance in the measurement noise.

$$H_0: \sigma_{iI} = \sigma_{iC} = 0 \tag{8}$$

372 The alternative hypothesis is that the neural signals are relevant for behaviour, reflecting the greater noise
373 on trials where the observer makes a suboptimal decision.

$$H_1: \sigma_{iI} > \sigma_{iC} \text{ , or } \sigma_{iI}^2 = (\sigma_{iC}^2 - x) \text{ where } x > 0 \tag{9}$$

374 The difference in the inference noise is limited by the total variance:

$$p(I)(\sigma_{iI}^2) + p(C)(\sigma_{iI}^2 + x) = \frac{1}{r_T^2} - 2 - \sigma_m^2 \tag{10}$$

375 Solving for $\sigma_{iI}^2$ (since $p(C) + p(I) = 1$):

$$\sigma_{iI}^2 = \frac{1}{r_T^2} - 2 - \sigma_m^2 - p(C)x \tag{11}$$

376 If we consider the correlation between the neural representation and the presented evidence on trials with
377 optimal responses and suboptimal trials separately (for simplicity, let $R = \frac{1}{r_T^2}$):

$$r_I = \frac{1}{\sqrt{R - p(C)x}} \tag{12}$$

378

$$r_C = \frac{1}{\sqrt{R - p(C)x - x}} \tag{13}$$

379 Setting a uniform prior on the ratio of inference and measurement noise, results in a linearly descending
380 prior on $x$:

$$p(x) = \frac{R - 2 - p(I)x}{\int_0^{(R-2))/p(I)} R - 2 - p(I)x \ dx} \tag{14}$$

381 We actually measure the difference in the Fischer transform of the correlation:

14

$$z_C - z_I = 0.5\log\left(\frac{(1 + r_c)(1 - r_I)}{(1 - r_c)(1 + r_I)}\right) \tag{15}$$

382 Since $r_c$ and $r_I$ are independent of the assumed measurement noise, there is one $x$ that corresponds to a
383 measured difference $z_C - z_I$, given the overall correlation $r_T$.

384 For each participant, for each electrode, at each time-point, the prior on $\sigma_m^2$ for $H_0$ is calculated by permuting
385 the data labels (accurate vs inaccurate behavioural responses). The probability of the data given $H_0$ and $H_1$
386 are calculated as above and used to compute the loglikelihood ratio:

$$LLR = log\left(\frac{p(D|H_1)}{p(D|H_0)}\right) \tag{16}$$

387 The clusters are identified using the Fast Subset Sums procedure: The loglikelihood ratios are summed
388 across participants, for each electrode and time-point. We then find small clusters by thresholding the log
389 posterior odds ratio:

$$POR = LLR + log\left(\frac{p(H_1)}{p(H_0)}\right) \tag{17}$$

390 Where the prior $p(H_1)$ is set to 0.05. The cluster with the largest LLR (summed across electrodes and time
391 points) is then expanded by continuing to add the largest neighbour and the new log prior ($p(H_1)$ = 0.05/n),
392 where $n$ is the size of the cluster, whilst the POR remains in favour of $H_1$. This is repeated until all clusters
393 with evidence in favour of $H_1$ have been identified.

## Supplementary Note 7

### Estimating single-sample confidence inference error

396 We aimed to examine the neural processes that are important for the precise representation of the decision
397 evidence for computing confidence. To do so, we explored the source(s) of the noise affecting the neural
398 representation of the accumulated evidence in the clusters of signals identified as relevant for suboptimal
399 confidence evaluations. We used the representation error as an estimate of the inference error of the
400 observer: the absolute difference between the cluster predicted value and the expected value given the
401 cluster representation and the true value of accumulated evidence based on the orientations presented to
402 the observer. This estimate is likely substantially affected by measurement noise, in addition to the inference
403 error of the observer. However, we do not expect measurement noise to be systematically driven by a
404 specific source, especially not across subjects. Noise Min and Noise Max epochs were selected by taking the
405 top and bottom quartiles of epochs sorted by representation error.

406 A separate estimate of the inference error was obtained by simulating the computational model (**Figure S6a**
407 shows the process of obtaining these estimates and their mutual reliance on the input stimulus variables and
408 the behavioural output). This measure also has its drawbacks: It is relatively imprecise, given the large range

15

409    of errors that are consistent with the observers' behavioural responses (see **Supplementary Note 2**); and is

410    based on the assumptions of the model. By examining these two estimates, we avoid relying on the same set

411    of assumptions throughout the analysis. The correlation between these estimates suggests that they do tap

412    into the suboptimal inference of the observer.

413    We considered how the different measures vary across samples and by the division in Noise Min and Noise

414    Max epochs. **Figure S6b** shows the correlation of these measures, averaged across subjects. The average

415    absolute effect size of the within subject difference between different variables dividing trials by Noise Min

416    and Noise Max epochs is shown in **Figure S6c**. There was a larger effect on confidence inference error ($d$ =

417    0.06) than perceptual inference error ($d$ = 0.02), from the model estimate. There were some effects on

418    stimulus variables: a small effect of condition (More vs Less, $d$ = 0.03), a large effect on sample position in the

419    sequence (Noise Min epochs tended to correspond to earlier samples, $d$ = 0.2), and an effect on decision

420    update (Noise Min epochs tended to correspond to smaller momentary decision updates, $d$ = 0.08). The

421    effects on behaviour were largest for confidence accuracy ($d$ = 0.06), with limited effect on perceptual

422    accuracy ($d$ = 0.02) and confidence rating (Noise Min epochs were somewhat more associated with high

423    confidence ratings, $d$ = 0.03).
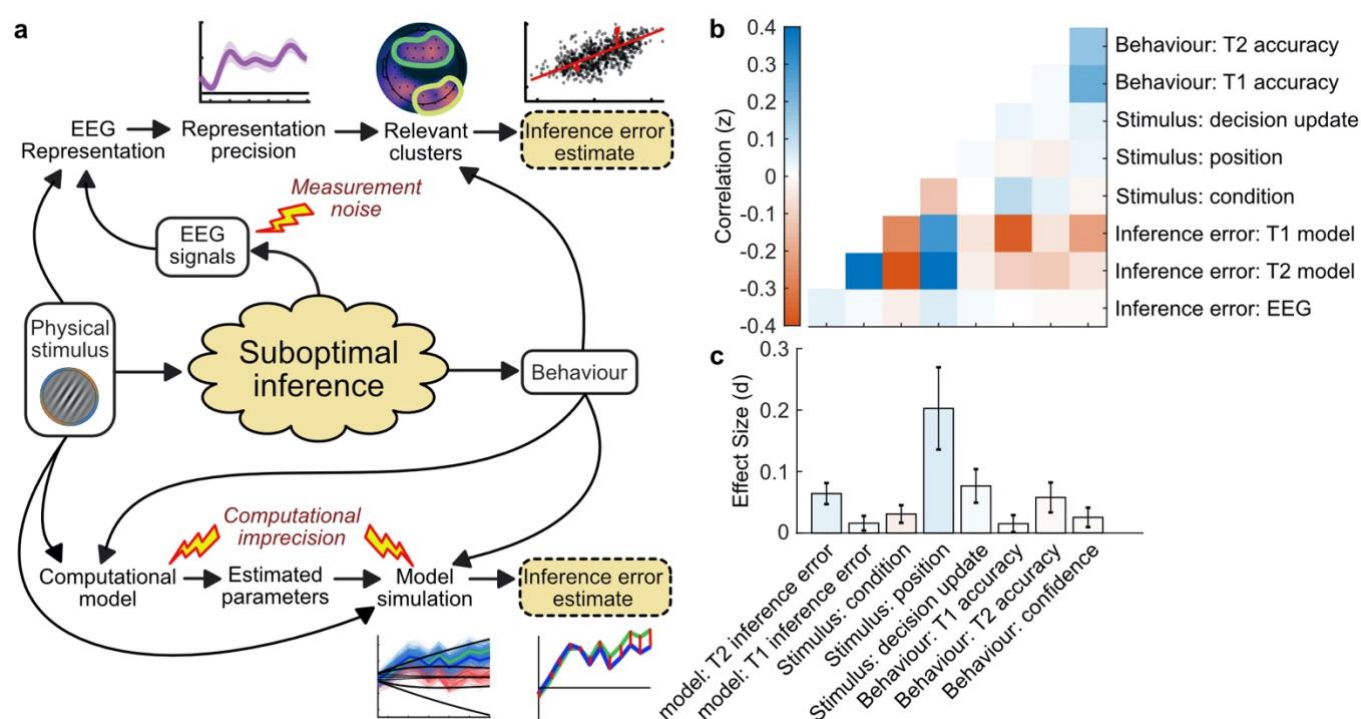


424

425    *Figure S6. Estimating inference error. a) Two approaches to estimate inference error. It is assumed the*

426    *observer's behaviour is based on a suboptimal inference over the physical stimulus. We do not have access to the*

427    *single-sample inference error, but can estimate it using the measured variables: the physical stimulus*

428    *properties, the behaviour, and the EEG signals. Two approaches are outlined: The EEG inference error estimate,*

429    *which relies on the error of the representation of the accumulated evidence, in clusters where the precision of*

430    *the representation is related to suboptimal behaviour; and the model error, which relies on simulating the*

431    *processing of the evidence based on the fitted model parameters, and taking the median of simulated traces*

432 *which concur with the observer's response. **b)** Correlation between variables measured from behaviour, the*

433 *stimulus input, and the estimated inference error. **c)** Effect size on the difference between Noise Min and Noise*

434 *Max epochs.*

## Supplementary Note 8

### Regions of interest

437 Regions of interest were selected based on the previous literature. Specifically, Herding et al. (2019) found

438 subjective evidence to modulate activity in the Superior Parietal Cortex; Gherman and Philiastides (2018)

439 found correlates of confidence encoding in the ventro-medial Prefrontal cortex (overlapping with the

440 MindBoggle Orbitofrontal Cortex coordinates), whilst Graziano et al., (2015) examined ROIs in the Anterior

441 Cingulate cortex, Orbitofrontal cortex, Temporal lobe, Posterior Parietal cortex, and Occipital cortex. We

442 chose to use ROIs defined by MindBoggle (Klein et al., 2017) that corresponded to similar regions: Lateral

443 Occipital cortex, Superior Parietal cortex, Orbitofrontal cortex (combining medial and lateral partitions),

444 rostral Middle Frontal cortex, and initially the Anterior Cingulate Cortex (combining rostral and caudal

445 partitions; **Figure S7a**). The results of the Anterior Cingulate Cortex were similar to the neighbouring

446 Orbitofrontal region, so we decided not to present this in the manuscript for simplicity. We show the results

447 in **Figure S7b**, for left and right hemispheres separately (statistical analyses were performed on the
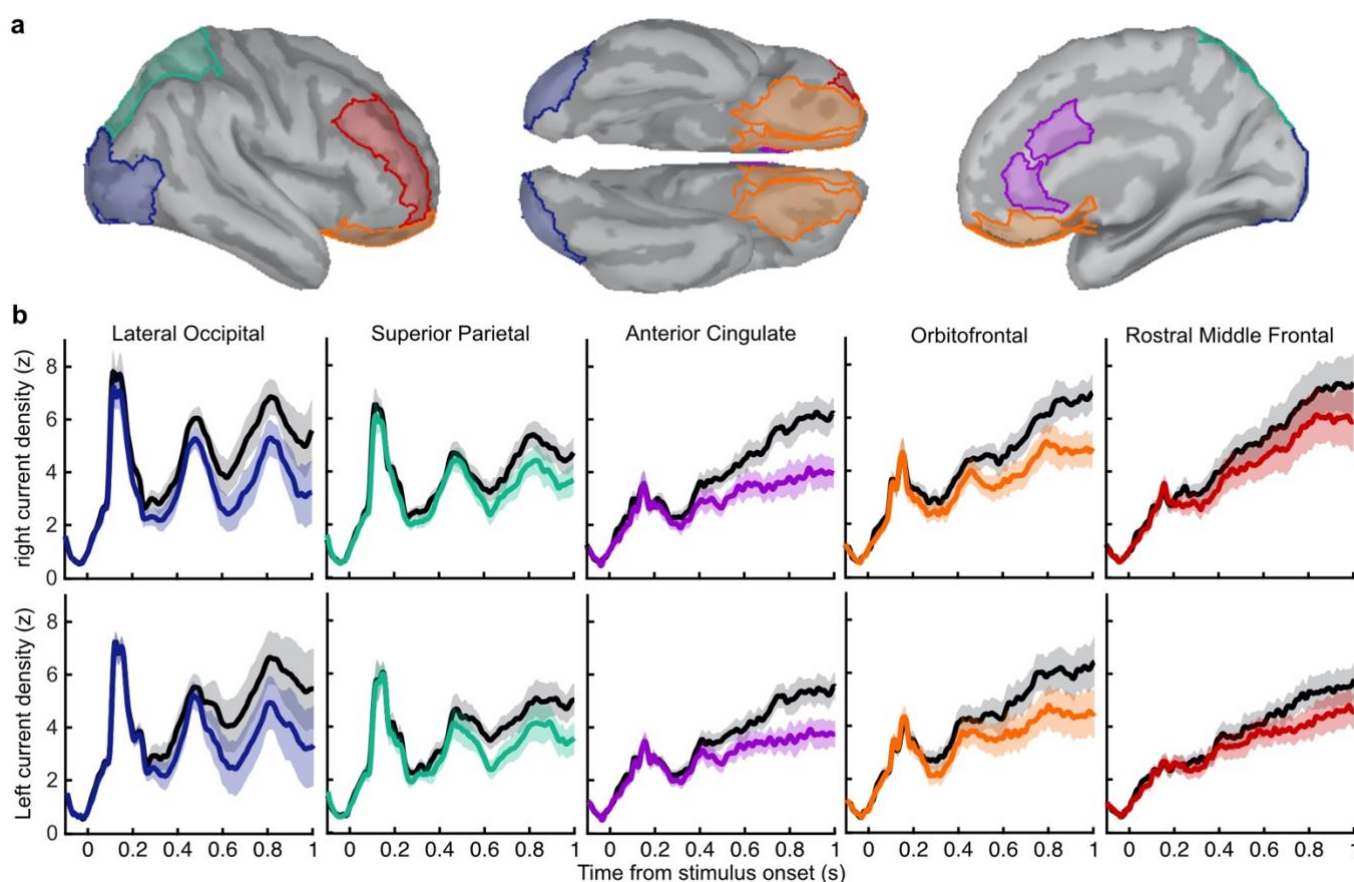
448 average).



450 ***Figure S7. Regions of interest and corresponding current density. a)*** *Regions of interest based on*

451    *Mindboggle coordinates. **b)** Average normalised rectified current density in the right (top) and left (bottom)*

452    *hemispheres. Noise Min epochs are shown coloured, Noise Max in black, with shaded regions showing the 95%*

453    *within-subject confidence interval.*

# References

Balsdon, T., Wyart, V., & Mamassian, P. Confidence controls perceptual evidence accumulation. *Nature Communications*, 2020; **11**(1), 1-11

Klein, A., Ghosh, S. S., Bao, F. S., Giard, J., Häme, Y., Stavsky, E., ... & Keshavan, A. Mindboggling morphometry of human brains. *PLoS Computational Biology,* 2017; **13**(2), e1005350.

Mamassian, P., & de Gardelle, V. Modelling perceptual confidence and the confidence forced-choice paradigm. *Under Review.*

Maniscalco, B., & Lau, H. The signal processing architecture underlying subjective reports of sensory awareness. *Neuroscience of Consciousness*, 2016; **1**.

Neill, D. B. Fast Bayesian scan statistics for multivariate event detection and visualization. *Statistics in Medicine*, 2011; **30**(5), 455-469.

Neill, D. B. Bayesian Scan Statistics. In: Glaz J., Koutras M. (eds) *Handbook of Scan Statistics.* 2019*;* Springer, New York, NY.

Oostenveld, R., Fries, P., Maris, E., & Schoffelen, J. FieldTrip: open source software for advanced analysis of MEG, EEG, and invasive electrophysiological data. *Computational Intelligence and Neuroscience*, 2011.

Salvador, A., Arnal, L. H., Vinckier, F., Domenech, P., Gaillard, R., & Wyart, V. Premature commitment to uncertain beliefs during human NMDA receptor hypofunction. bioRxiv 2020