# Improvement of association between confidence and accuracy after integration of discrete evidence over time

1  **Zahra Azizi[1], Sajjad Zabbah[2], Azra Jahanitabesh[3], Reza Ebrahimpour[2,4*]**

2  [1] Department of Cognitive Modeling, Institute for Cognitive Science Studies, Tehran, Iran.

3  [2] Institute for Research in Fundamental Sciences, School of Cognitive Sciences, Tehran, Iran.

4  [3] Department of Psychology, University of California, Davis, California, United States.

5  [4] Department of Artificial Intelligence, Faculty of Computer Engineering, Shahid Rejaee Teacher
6  Training University, Tehran, Iran.

7  * Correspondence:

8  Reza Ebrahimpour

9  ebrahimpour@ipm.ir

10  **Keywords: Confidence; metacognition; discrete pieces of evidence; perceptual decision-making;**
11  **pupillometry; ERP.**

12  **Abstract**

13  When making decisions in real-life, we may receive discrete pieces of evidence during a time period.
14  Although subjects are able to integrate information from separate cues to improve their accuracy,
15  confidence formation is controversial. Due to a strong positive relation between accuracy and
16  confidence, we predicted that confidence followed the same characteristics as accuracy and would
17  improve following the integration of information collected from separate cues. We applied a Random-
18  dot-motion discrimination task in which participants had to indicate the predominant direction of dot
19  motions by saccadic eye movement after receiving one or two brief stimuli (i.e., pulse(s)). The interval
20  of two pulses (up to 1s) was selected randomly. Color-coded targets facilitated indicating confidence
21  simultaneously. Using behavioral data, computational models, pupillometry and EEG methodology we
22  show that in double-pulse trials: (i) participants improve their confidence resolution rather than
23  reporting higher confidence comparing with single-pulse trials, (ii) the observed confidence follow
24  neural and pupillometry markers of confidence, unlike in weak and brief single-pulse trials. Overall,
25  our study showed improvement of associations between confidence and accuracy in decision results
26  from the integration of stimulus separated by different temporal gaps.

27  **1    Introduction**

28  Humans and animals can both make choices based on multiple discrete pieces of information. Imagine
29  that a large bus is passing between you and a faraway car as you cross the street. In this situation,
30  simply by collecting discrete pieces of information about the car's position through the windows of the
31  bus, you can decide whether the car is moving toward or away from you. In this scenario, as the number
32  of pieces of information increased, the interpretation of the car's direction would be improved. Indeed,
33  research has shown that the accuracy of decisions can be significantly improved by integrating

34  information from separate cues (Kiani, Churchland, & Shadlen, 2013; Kira, Yang, & Shadlen, 2015;
35  tickle, Tsetsos, Speekenbrink, & Summerfield, 2020; Tohidi-Moghaddam, Zabbah, Olianezhad, &
36  Ebrahimpour, 2019; Waskom & Kiani, 2018). Typically, our decisions are accompanied by feelings
37  that reflect the likelihood that the decision is correct; such a feeling is called confidence (Kiani,
38  Corthell, & Shadlen, 2014). For example, imagine that the scene in the previous scenario is also
39  included a foggy weather. In this case, low visibility may reduce the confidence of your judgments.
40  This diminished confidence per se may lead to change your mind (Fleming, Putten, & Daw, 2018;
41  Resulaj, Kiani, Wolpert, & Shadlen, 2009), impact your behavioral adjustments, and affect how
42  quickly and accurately you make your consecutive decisions (Meyniel, Sigman, & Mainen, 2015; van
43  den Berg, Zylberberg, Kiani, Shadlen, & Wolpert, 2016). Due to the potential effects of confidence on
44  decision-making, in the last few years, considerable progresses had been made in the understanding of
45  the behavioral (Kiani et al., 2014; Zylberberg, Barttfeld, & Sigman, 2012) and the neuronal (Baranski
46  et al., 2017; Gherman & Philiastides, 2015; Kiani & Shadlen, 2009) properties of confidence and its
47  association with perceptual decision-making. However, how confidence is established within a discrete
48  environment is still unclear.

49  According to the leading computational approach in perceptual decision-making (Gold & Shadlen,
50  2007; Shadlen & Kiani, 2013), when the accumulated evidence for one option, called a decision
51  variable (DV), crosses a threshold or a boundary, a decision would be made. In addition, confidence is
52  briefed by the probability that a decision relying on the DV is correct (Kiani et al., 2014; Kiani &
53  Shadlen, 2009; van den Berg et al., 2016; Zylberberg, Fetsch, & Shadlen, 2016). Research has
54  confirmed a strong positive relation between accuracy and confidence (Kiani et al., 2014; Vafaei
55  Shooshtari, Esmaily Sadrabadi, Azizi, & Ebrahimpour, 2019). Moreover, it has been shown that, when
56  we need to decide based on the discrete pieces of evidence, the decision is determined by integrating
57  the DV of all those pieces (Kiani et al., 2013; Waskom & Kiani, 2018) and the accuracy even exceeded
58  expectations predicted by evidence integration models (Kiani et al., 2013). Accordingly, one may
59  suggest that confidence would follow the same characteristics as accuracy and would increase
60  considerably after receiving separate pieces of information.

61  Nevertheless, a large body of evidence (e.g. (Herce Castañón et al., 2019; Zylberberg et al., 2016))
62  determines that human observers do not report their confidence in consistent with their accuracy. From
63  this standpoint, noise can be considered as the key parameter to clarify variations in confidence (Kiani
64  et al., 2014; Zylberberg et al., 2012). For instance, an underestimation of sensory noise in decisions
65  would lead to over and/or under-confidence (De Gardelle & Mamassian, 2015; Herce Castañón et al.,
66  2019; Zylberberg, Roelfsema, & Sigman, 2014) such that observers may ignore evidence in favor of
67  other alternatives (Zylberberg et al., 2012). Moreover, confidence ratings may not only originate from
68  the available sensory evidence (Rahnev & Denison, 2018; Zylberberg et al., 2016). So, the observers
69  may integrate additional evidence into their confidence rating, which was not used for making their
70  decision, allowing them to change their mind after the initiation of a response (Atiya et al., 2020;
71  Resulaj et al., 2009). This suggests that computational description of confidence would be controlled
72  by the attendance of both decision and confidence performance (Balsdon, Wyart, & Mamassian, 2020;
73  Maniscalco & Lau, 2014).

74  To test the hypothetical relation between the accuracy and confidence, in binary decisions, signal
75  detection theory (SDT) can provide a method to characterize how well the observers reporting the
76  confidence ratings by introducing metacognitive sensitivity and efficiency (**Figure 1B**; (Fleming,
77  2017; Maniscalco & Lau, 2012, 2014)). In fact, for years, SDT has provided a simple yet powerful
78  methodology to distinguish between an observer's ability to categorize the stimulus and the behavioral
79  response (Green & Swets, 1966), and to determine confidence resolution.

80 Moreover, levels of confidence can be tracked by behavioral, neural and pupillometry signatures.
81 Higher confidence are accompanied by faster and more accurate decisions (Kiani et al., 2014; van den
82 Berg et al., 2016; Zylberberg et al., 2016). In addition, research on perceptual decision-making has
83 established an EEG potential characterized by a centro-parietal positivity (CPP) as a neural correlate
84 of sensory evidence accumulation (Kelly & O'Connell, 2013; O'connell, Dockree, & Kelly, 2012) and
85 confidence (Boldt, Schiffer, Waszak, & Yeung, 2019; Herding, Ludwig, von Lautz, Spitzer, &
86 Blankenburg, 2019; Tagliabue et al., 2019; Vafaei Shooshtari et al., 2019; Zizlsperger, Sauvigny,
87 Händel, & Haarmeier, 2014). In particular, it has been shown that the CPP, despite the difference, is
88 present for both correct and incorrect decisions (O'connell et al., 2012; Steinemann, O'Connell, &
89 Kelly, 2018) and can reflect not only external evidence but also an internal decision quantity such as
90 decision confidence. In addition, levels of confidence can be tracked by monitoring the pupil. The
91 literature has suggested strong links between pupil dilation and both the decision (Murphy, Boonstra,
92 & Nieuwenhuis, 2016) and confidence (Allen et al., 2016; Lempert, Chen, & Fleming, 2015; Urai,
93 Braun, & Donner, 2017) via pupil-linked dynamics of the noradrenergic system (Laeng, Sirois, &
94 Gredebäck, 2012). For example, pupillometry has provided some evidence that shows a partial
95 dissociation between choice and confidence in decision-making (Balsdon et al., 2020). Considering the
96 potential of response-time, CPP and pupillometry signatures to capture the distinction between choice
97 and confidence in decision-making, they can be considered as informative paradigms to explore the
98 confidence-accuracy association.

99 Accordingly, to bridge the existing gap in the confidence and perceptual decision-making literature,
100 we implemented two separate experiments to explore three questions: First, how participants
101 accumulate discrete evidence to establish confidence judgments. Second, whether the confidence
102 ratings are in accordance with accuracy after integration of discrete evidence. Finally, how implicit
103 markers of confidence —response-time, CPP and pupillometry— change after receiving separated
104 pieces of information. Here, to clarify confidence, we required observers to make a two-alternative
105 decision after viewing either one (single-pulse) or two (double-pulse) motion pulses separated by four
106 various temporal gaps (similar to (Kiani et al., 2013; Tohidi-Moghaddam et al., 2019)). We performed
107 several logistic regression models to measure the impact of stimulus characteristics on confidence.
108 Also, we applied a set of computational models based on SDT to assess how accuracy and confidence
109 varied throughout the experiments. Then, in the second experiment, we used EEG methodology to
110 examine the relation between participants' brain activity and their confidence. We expected a neural
111 indicator of perceptual decision making (CPP) would show amplitude changes between the two levels
112 of confidence. In addition to behavioral data and EEG methodology, participants' pupil response was
113 monitored across both experiments to examine the relation between participants' pupil dilation and
114 their confidence. The findings expose that participant integrated information from pulses, invariant to
115 the temporal gap, to improve the confidence resolution instead of reporting higher confidence.
116 Likewise, in double-pulse trials, behavioral, neural and pupillometry markers of confidence would be
117 distinguishable, entirely unlike in brief and weak single-pulse trials.

118 **2    Materials and Methods**

119 **2.1    Participants**

120 Consistent with methodological considerations in previous studies, a total of 19 observers participated
121 in the two experiments. Six participants (three male; $M_{age}$ = 32.25; $SD_{age}$ = 4.5) attended in our
122 behavioral experiment —Experiment 1— and 13 participants (three males; $M_{age}$ = 31.41; $SD_{age}$ = 5.56)
123 took part in our EEG experiment —Experiment 2. All participants had normal or corrected-to-normal

3

124  vision, and none of them had any history of psychiatric and neurological disorders. Previous studies
125  with the same paradigm in which a large number of trials were presented to a small number of
126  participants (e.g., five participants in(Kiani et al., 2013); six participants in (Kiani et al., 2014); four
127  participants in (van den Berg et al., 2016) and six participants in (Stine, Zylberberg, Ditterich, &
128  Shadlen, 2020)), assume that with extensive training, all participants would reach an acceptable level
129  of performance. As such, a small number of trained participants would perform similar to the
130  performance of a large number of participants. Accordingly, to make participants' performance reach
131  the same criteria and reduce the between-participant variability, all participants received extensive
132  training sessions on the Random-dot-motion discrimination task prior to data collection. Moreover,
133  participants' understanding of the confidence reporting procedure was double-checked prior to the
134  experiments.  In Experiment 1, one participant was excluded due to the difficulty in reporting decision
135  and confidence simultaneously, and another participant decided to leave the experiment shortly after
136  participation. In addition, one participant was excluded from Experiment 2 because of the excessive
137  noise in EEG electrodes crucial to the analysis.

## 2.2  Stimuli

139  We explored the confidence formation in discrete environment with a random-dot-motion (RDM)
140  discrimination task. Participants had to indicate the predominant motion direction of a cloud of moving
141  dots (left or right) presented within a 5° circular aperture at the center of the screen. The dot density
142  was 16.7 dots/degree$^2$/s and the displacement of the coherently moving dots produced an apparent
143  speed of 6 deg/s. The RDM movies were generated by three interleaved sets of dots presented on
144  consecutive video frames. Three video frames later, each dot was redrawn at a location consistent with
145  the direction of motion or at a random location within the stimulus space. More details can be found in
146  previous studies (e.g. (Roitman & Shadlen, 2002)). The experiment code was programmed in
147  MATLAB 2016a (Mathworks Inc., USA) using PsychToolbox (Brainard & Vision, 1997; Kleiner,
148  Brainard, & Pelli, 2007)

## 2.3  Experimental Tasks

150  Participants performed the RDM task in blocks of 200 trials. Each trial started with participants fixating
151  a small red point (diameter 0.3°) at the screen center. After 500 ms, two choice-targets appeared to the
152  left and right of the fixation point (10° eccentricity; **Figure 1A**). Each target was shaped as a gradient
153  rectangle (9° length and 0.5° width). After a variable duration of 200 - 500 ms (truncated exponential
154  distribution), the RDM was presented. Participants had to indicate their choice after receiving one or
155  two pulses of 120ms of motion pulses. The gap interval of double-pulse trials was selected randomly
156  from 0, 120, 360, and 1080ms. On single-pulse trials, motion coherence was randomly selected from
157  these six values: 0%, 3.2%, 6.4%, 12.8%, 25.6%, and 51.2%, whereas, on double-pulse trials, motion
158  coherence of each pulse was randomly chosen from three values: 3.2%, 6.4%, and 12.8%. Both pulses
159  had the same net direction of motion and participants were aware of it. In total, there were 6 single-
160  pulse and 9 × 4 double-pulse trial types. After the offset of one or two motion pulses, a 400 to 1000 ms
161  delay period (truncated exponential) was imposed before the Go signal appeared on the screen. In each
162  trial, participants were required to indicate their response by directing the gaze to one of the targets,
163  the upper extreme of targets representing full decision confidence and the lower extreme representing
164  guessing (**Figure 1A**). To provide the approximate balance within the trials, we constructed a list of
165  all possible conditions of motion coherences and gaps. Then, we shuffled the listed conditions and
166  assigned them randomly to the trials in each block. Participants were instructed to achieve high
167  performance. Distinctive auditory feedback (Beep Tones) was provided for correct and incorrect
168  responses. The type of feedback of 0% coherence trials was selected randomly by a uniform

169  distribution. In Experiment 1, each participant performed the task across multiple blocks on different
170  days (12-20 blocks). Experiment 2 contained the same paradigm as Experiment 1. All variables of
171  stimulus remained constant except, in Experiment 2, the EEG data were also recorded. In Experiment
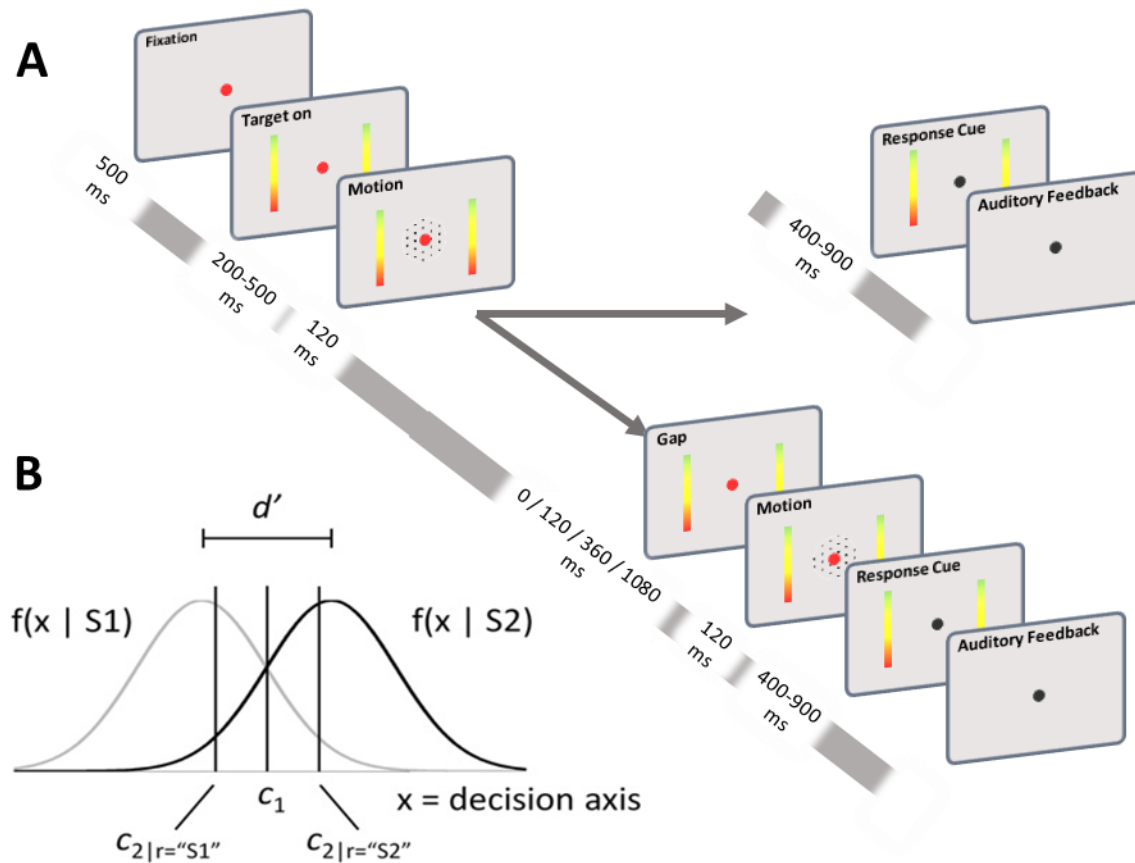172  2, each participant completed a session of 4-5 blocks.

173



**Figure 1. Task paradigm and Signal Detection Theory**. **(A)** Participants had to indicate the predominant direction of motion of moving dots (left or right) by saccadic eye movement to one of the targets after receiving one or two pulse(s) of 120ms stimulus. The intervals between two pulses were selected randomly from 0 to 1080 ms and the direction of both pulses were the same. Color-coded targets enabled participants indicating their confidence simultaneously. **(B)** On each trial, a stimulus generates an internal response $x$ within an observer, who must use $x$ to decide whether the stimulus is $S_1$ or $S_2$, $x$ is drawn from a normal distribution. The distance between these distributions is $d'$, which measures the observer's ability to discriminate $S_1$ from $S_2$. The observer also rates decision confidence on a scale of high and low by comparing $x$ to the additional response specific confidence criteria ($cr_2$ for each option). For details, see Supplementary Appendix 2 and refs (Fleming, 2017; Maniscalco & Lau, 2012, 2014).

174

## 2.4   EEG Recording and pre-processing

176  We used a 32-channel amplifier for the EEG signal recording (eWave, produced by ScienceBeam,
177  http://www.sciencebeam.com/) which provided 1K sample/s of time resolution. EEG was recorded at

178    31 scalp sites (Fp1, Fp2, AF3, AF4, C3, C4, P3, P4, O1, O2, F7, F8, T7, T8, P7, P8, FPz, Fz, Cz, Pz,
179    Oz, POz, FC1, FC2, CP1, CP2, FC5, FC6, CP5, CP6). The EEG signals were referenced to the right
180    mastoid. The recorded data were taken to Matlab (Mathworks Inc., USA) and pre-processed as follows.
181    The signals were filtered using a band-pass filter from 0.1 Hz to 40 Hz (Zizlsperger et al., 2014) for
182    removing high frequency and independent cognitive noises. Then, all trials were inspected, and those
183    containing Electromyography (EMG) or other artifacts were identified and manually removed. The
184    second artifact rejection step included independent components analysis (ICA) using the EEGLAB
185    toolbox (Delorme & Makeig, 2004). To select the removable ICA component, the ADJUST plugin
186    (Mognon, Jovicich, Bruzzone, & Buiatti, 2011) was used.

## 2.5    Pupillometry Recording and pre-processing

188    The eye data were collected using an EyeLink 1000 infrared eye-tracker system (SR Research Ltd.
189    Ontaro, Canada). This device allowed a 1000-Hz sampling rate and was controlled by a dedicated host
190    PC. The system was calibrated and validated before each block by presenting nine targets at the center,
191    edges, and corners of the display monitor. The left eye's data was recorded and passed to the host PC
192    via an Ethernet link during data collection.

193    Missing data and blinks, as detected by the EyeLink software, were padded and interpolated.
194    Additional blinks were spotted using peak detection on the pupil signal's velocity and then linearly
195    interpolated (Mathôt, 2013).

## 2.6    Experimental procedure

197    In this study we employed behavioral, neural, and pupillometry signatures. Participants were given a
198    consent form in which the experiment was described in general terms. After providing written informed
199    consent, in both experiments, participants completed the tasks in a semidark, sound-attenuating room
200    to minimize distraction. All instructions were presented and stimuli were displayed on a CRT monitor
201    (17 inches; PF790; refresh rate, 75 Hz; screen resolution, 800 × 600). A head and chin rest confirmed
202    that the distance between the participants' eyes and the monitor's screen was 57 cm throughout the
203    experiment. Participants were presented demographic questions followed by training sessions and main
204    sessions, respectively. The experimental protocol was approved by the ethics committee of the Iran
205    University of Medical Sciences.

## 2.7    Data Analysis

207    Data analysis was performed using Matlab 2019a (The MathWorks Inc., United States).

### 2.7.1 Quantifying confidence

209    Reported confidence was categorized as high and low. Since the participants were told to choose the
210    upper part of the bar as high confidence and lower part as low confidence, we considered reported
211    confidence higher than midline as high confidence and lower than midline as low confidence
212    respectively. This categorization allowed us to take each confidence report as a binary variable
213    comparable to the choice. Using categorical variables also provided the possibility of comparing the
214    current data with our previous work (Vafaei Shooshtari et al., 2019). However, in addition to the
215    midline, we tested various binary level set methods for categorizing participants' high and low
216    confidence ratings. First, the highest 55% and 45% of each participant's confidence reports were
217    considered high confidence (similar to (Zylberberg, Wolpert, & Shadlen, 2018)). Then, the mean of

This is a provisional file, not the final typeset article

218  each participant's confidence was calculated separately, and the confidence ratings above the mean
219  were considered as high ratings. Using these methods did not significantly alter reported confidence
220  categorization (see **Supplementary Figure 6**).

221  **2.7.2 Behavioral analyses**

222  Except where otherwise specified, we reported behavioral data of the first experiment but all the
223  analyses were repeated for the EEG experiment and if the results were inconsistent, it has been admitted
224  (EEG experiment results were reported in **Supplementary Figures 1**, **2**, **3**, **4** and, **5**).

225  We performed several logistic regression models to measure the impact of stimulus characteristics on
226  binary outcomes after confirming the assumptions of the linear regression were met. For logistic
227  regression models, we used maximum likelihood under a binomial error model (i.e., a GLM) to
228  evaluate the null hypothesis that one or more of the regression coefficients were equal to zero. $P_{high}$
229  was the probability of high confidence, $Logit[P_{high}]$ indicated log $\frac{P_{high}}{1 - P_{high}}$ and $\beta_i$ denoted fitted
230  coefficients. Also, $P_{correct}$ was the probability of correct response and $Logit[P_{correct}]$ indicated log
231  $\frac{P_{correct}}{1 - P_{correct}}$.

232  For single-pulse trials, the probability of a high confidence choice was given by the following:

$$Logit[P_{high}] = \beta_0 + \beta_1 C, \qquad (1)$$

233  where $C$ was motion strengths of the pulse. Likewise, the probability of a correct choice was stated by
234  the logistic regression:

$$Logit[P_{correct}] = \beta_0 + \beta_1 C, \qquad (2)$$

235  To examine whether confidence judgments were associated with more accurate choices, we fitted a
236  logistic regression model to accuracy where the probability of high confidence is given by:

$$Logit[P_{high}] = \beta_0 + \beta_1 A, \qquad (3)$$

237  where $A$ was the accuracy of the response (0 or 1 for incorrect and correct) and our null hypothesis was
238  that the accuracy would not affect reported confidence ($H_0: \beta_1 = 0$). We also used logistic regression
239  to evaluate the effect of interpulse interval on confidence in double-pulse trials:

$$Logit[P_{high}] = \beta_0 + \beta_1 C_1 + \beta_2 C_2 + \beta_3 T + \beta_4 C_1 T + \beta_5 C_2 T, \qquad (4)$$

240  where $C_1$ and $C_2$ were motion strengths of each pulse, and $T$ was the interpulse time interval. For
241  double-pulse trials with equal pulse strength ($C_1 = C_2$), the redundant regression terms ($\beta_2, \beta_4$) were
242  omitted. The null hypothesis was that the interpulse interval would not affect reported confidence
243  ($H_0: \beta_{3-5} = 0$). The similar equation was used to assess relation of accuracy and time interval:

$$Logit[P_{correct}] = \beta_0 + \beta_1 C_1 + \beta_2 C_2 + \beta_3 T + \beta_4 C_1 T + \beta_5 C_2 T, \qquad (5)$$

244  The null hypothesis was that the interpulse interval would not affect performance ($H_0: \beta_{3-5} = 0$). To
245  evaluate the impact of pulse sequence on confidence, the following regression model was fitted:

$$\text{Logit}[P_{high}] = \beta_0 + \beta_1[C_1 + C_2] + \beta_2[C_2 - C_1], \tag{6}$$

246  where $C_1$ and $C_2$ were corresponding motion strengths of each pulse. $\beta_2$ indicated how the confidence
247  varied from trials in which $C_1 > C_2$ to trials with a reversed sequence of motion pulses $C_1 < C_2$. The
248  null hypothesis was that the sequence of motion pulses did not influence the confidence ($H_0: \beta_2 = 0$).

249  To examine the interaction between the two pulses (e.g., a stronger pulse 1 reduced the effect of pulse
250  2), we fitted the following regression model to all double-pulse trials:

$$\text{Logit}[P_{high}] = \beta_0 + \beta_1 C_1 + \beta_2 C_2 + \beta_3 C_1 C_2, \tag{7}$$

251  The null hypothesis was that there was not an interaction between motion strengths of pulses ($H_0: \beta_3 =$
252  $0$). In other words, higher influence of second pulse on confidence was due to higher sensitivity rather
253  than an interaction of motion pulses and $\beta_2 > \beta_1$ confirmed greater sensitivity to the second pulse on
254  the decision.

255  In addition to logistic regression models, to investigate the variation of confidence in double-pulse
256  trials compared to single-pulse trials, we subtracted participants' confidence of double-pulse trials from
257  corresponding confidence in single-pulse trials. For example, the confidence of a sequence of 3.2%,
258  6.4% motion strength trial, subtract separately once from 3.2% and once from 6.4% corresponding
259  confidence in single-pulse trials. The process repeated for the data of each gap too. Moreover, the same
260  method was used to compare accuracy of double-pulse trials and single-pulse trials. To assess the effect
261  of choice accuracy on variation of confidence in double-pulse and single-pulse trials, we fitted the
262  following:

$$S_{Conf} = \beta_0 + \beta_1 A, \tag{8}$$

263  where the $S_{Conf}$ was the subtraction of confidence in double-pulse trials from corresponding single-
264  pulse trials and $A$ was the accuracy of the response (0 or 1 for incorrect and correct). The null hypothesis
265  was the choice accuracy did not affect the variation of $S_{Conf}$ ($H_0: \beta_1 = 0$).

### 266  2.7.2.1 Response-time analysis

267  In the current study, response-time was referred to the time between the cue onset and a participant's
268  response. To evaluate the significance of the effect of response-time on confidence, we fitted the
269  following linear regression model separately in double-pulse and single-pulse trials:

$$\text{Logit}[P_{high}] = \beta_0 + \beta_1 R, \tag{9}$$

270  where $R$ was the response-time of each trial and the null hypothesis was that confidence did not depend
271  on the response-time ($H_0: \beta_1 = 0$). Moreover, to evaluate the relation of delay-time imposed before
272  the cue onset and response-time, we fit a linear regression model as follows:

This is a provisional file, not the final typeset article

$$RT = \beta_0 + \beta_1 D, \tag{10}$$

273  where $D$ was the delay-time. The null hypothesis was that response-time did not rest on the delay-time
274  ($H_0: \beta_1 = 0$).

275  In addition, confidence is tracked by both evidence and response-time (Kiani et al., 2014; van den Berg
276  et al., 2016; Zylberberg et al., 2016), and indeed accuracy is relied on evidence. Furthermore, to study
277  the profile of high and low confidence from behavioral data, an equal number of trials from each
278  participant's trials was selected randomly from single/double-pulse trials. Same procedure repeated
279  100 times, then individual response-time were rank-ordered and binned into four quintiles. Then, the
280  accuracy of high and low confidence trials in each bin was calculated. We expected to see a significant
281  difference between accuracy of each bin grouped by levels of confidence. We only included motion
282  strength of 3.2, 6.4, 12.8 of single-pulse trials (similar to coherence used in double-pulse trials) to
283  control the impact of coherence on response-time.

### 2.7.3 Motion energy analysis

285  Random dot stimulus is stochastic, so the sensory evidence fluctuated within and across trials but
286  around the nominal motion coherence level. To examine the fluctuations in motion during each trial,
287  we filtered the sequence of random by using two pairs of quadrature spatiotemporal filters, as specified
288  in previous studies (Adelson & Bergen, 1985; Kiani, Hanks, & Shadlen, 2008; Zylberberg et al., 2012).
289  Since we aimed to understand the temporal course of choice and confidence, we summed the energies
290  across trials for each pulse in single/double-pulse trials.

291  We used logistic regression to test whether the confidence was more influenced by the second pulse's
292  motion energy than that of the first pulse in double-pulse trials. We tested double-pulse trials with equal
293  motion strength using the following logistic regression model:

$$\text{Logit}[P_{high}] = \beta_0 + \beta_1 C + \beta_2(M_1 + M_2) + \beta_3 M_2, \tag{11}$$

294  where $M_1$ and $M_2$ were the motion energy of each pulse. The null hypothesis was that the second pulse
295  was not more functional ($H_0: \beta_3 = 0$). We tested double-pulse trials with unequal motion strength by
296  modifying the regression model to:

$$\text{Logit}[P_{high}] = \beta_0 + \beta_1 C_1 + \beta_2 C_2 + \beta_3(M_1 + M_2) + \beta_4 M_2, \tag{12}$$

297  and the null hypothesis was ($H_0: \beta_4 = 0$). To evaluate the relation of $P_{high}$ and motion energy in single-
298  pulse trials, we fitted a linear regression model as follows:

$$\text{Logit}[P_{high}] = \beta_0 + \beta_1 C + \beta_2 M, \tag{13}$$

299  where $M$ was the motion energy of the presented motion stimulus and the null hypothesis was that
300  confidence did not depend on the motion energy ($H_0: \beta_2 = 0$).

### 2.7.4 General computational modeling approach

301  We implemented a set of computational models based on signal detection theory to provide a
302  mechanistic explanation of the experimental data. According to SDT, observers set a decision criterion
303  ($cr$) to discriminate between two stimuli (e.g., labeled as $S_1$ and $S_2$). They also set criteria $cr_{2,"S_1"}$ and
304  $cr_{2,"S_2"}$ to determine confidence ratings around the decision criterion $cr$ (**Figure 1B**; for more details,
305  see Supplementary Appendix 2). We computed stimulus sensitivity ($d'$) and measures of metacognitive
306  ability ($Meta\text{–}d'$, $Meta\text{–}d'/d'$). We used code provided by Maniscalco and Lau (Maniscalco & Lau,
307  2012) in which metacognitive sensitivity ($Meta\text{–}d'$) is computed by setting the $d'$ value that would
308  produce the observed confidence. In addition, $Meta\text{–}d'/d'$ were calculated by normalizing
309  $Meta\text{–}d'$ by $d'$ through division. Here, $d'$, $Meta\text{–}d'$ and, $Meta\text{–}d'/d'$ of single-pulse and double-
310  pulse trials were computed separately. In addition, we fitted SDT model with trials simulated by a
311  perfect integrator model (the model is described later). We then addressed the trend of $d'$, $Meta\text{–}d'$
312  and, $Meta\text{–}d'/d'$ of three models for each participant. To support the fact that our findings were not
313  relevant to variation of coherence of single and double-pulse trials, we only included single-pulse trials
314  with motion strength of 3.2, 6.4, 12.8. However, one difference between groups was that they might
315  not be matched for the number of trials: the single-pulse included on average fewer trials for each
316  coherence per participant compared to double-pulse trials. Previous research has suggested that the
317  number of trials could bias measures of metacognitive ability (Fleming, 2017). Therefore, in a control
318  analysis, we created 100 sets of trials randomly from the single/double-pulse trials and from trials
319  simulated by the perfect integrator model. Each set contained the same number of trials for each
320  participant. We then averaged the metacognitive scores obtained from these 100 sets and repeated the
321  comparison procedure (see **Supplementary Figure 6**).

### 2.7.5 Perfect integrator Model

323  To estimate the expected high confidence ($P_{e(high)}$) in double-pulses trials, we assumed that each trial's
324  confidence was achieved based on evidence integrating from both pulses by using a perfect integrator
325  model. In the perfect integrator model, the expected accuracy ($P_{e(correct)}$) for double-pulse trials
326  computed as the following (Kiani et al., 2013):

$$P_{e(correct)} = 1 - \phi(0, e_1 + e_2, \sqrt{2}), \tag{14}$$

328  where $e_1$ and $e_2$ were the pieces of evidence that underlie by $P_1$ and $P_2$ (the probabilities of the correct
329  answer in corresponding single-pulse trials) and were computed as:

$$e_i = \phi^{-1}(P_i, 0, 1), \text{ I} = 1,2, \tag{15}$$

330  Where $\phi^{-1}$ was inverse $\phi$, which represented the cumulative Gaussian distribution (Kiani et al., 2013).

331  To predict the confidence of double-pulse trials by this model, after calculating $cr$ and $d'$ (see
332  Supplementary Appendix 2), $cr$ was shifted to zero and $d'$ was normalized. Then, confidence Hit Rate
333  and False Alarm Rate were calculated based on confidence performance from corresponding single-
334  pulse trials (similar to Eq.14, 15). Accordingly, high confidence probability (for both correct response
335  or incorrect response) would be predicted by the perfect integrator model. Besides, the model
336  parameters, including confidence criteria along with $Meta\text{–}d'$ were computed.

### 2.7.6 Confidence optimized model

In the confidence optimized model, we optimized the confidence criteria in the perfect integrator model by providing each participant's confidence performance computed of double-pulse trials. The purpose of this simulation was to understand why the perfect integrator model was not able to predict confidence well.

### 2.7.7 Model evaluation

We evaluated the models qualitatively (i.e., parameter recovery exercises) and quantitatively (i.e., maximum likelihood estimation).

In the qualitative method, based on the calculated parameters of the model, the probability of choosing high confidence for all combinations of motion strength for each participant were calculated (see Supplementary Appendix 2). We compared the expected high confidence predicted by models to the observed confidence in double-pulse trials using regression, as follows:

$$P_{high} = \beta_1 P_{e(high)} + \beta_0, \tag{16}$$

where $P_{e(high)}$ was the expected probability of high confidence. We regressed predicted vs. observed $P_{high}$ and compare slope ($\beta_1$) against the 1:1 line in each model. In this linear regression, we expected the predicted values to be close to the actual values.

In addition, to compare models quantitatively, an equal number of trials from each subject's trials selected randomly and then each model fitted to the selected data. This procedure repeated for 100 times, then the computed MLEs of each model was averaged.

### 2.7.8 Confidence suboptimality

The optimal decision-making is disrupted by several sources of suboptimality (Balsdon et al., 2020). In SDT, an added noise, $\xi_n$, represents a potential loss of information between sensory decision information and metacognitive information, such as confidence rating. This noise has a Gaussian distribution with zero mean, and standard deviation $\sigma$ (Maniscalco & Lau, 2014). The parameter $\sigma$ determines how much noisier the metacognitive variable is than the decision variable (Maniscalco & Lau, 2014).

$$\xi_n = N(0, \sigma)) \tag{17}$$

This noise is correlated to metacognitive efficiency ($Meta\text{--}d'/d'$) (Maniscalco & Lau, 2014). To consider this suboptimality, we simulated trials using the same parameter values resulted from the perfect integrator model except this noise was increased.

### 2.7.9 EEG analysis

The EEG analysis focused on a neural marker of perceptual decision-making linked with stimulus preparation and stimulus processing. The component we focused on was the centro-parietal positivity (CPP) which possibly identical to the classic P300 component (Herding et al., 2019; Twomey, Murphy, Kelly, & O'connell, 2015). The CPP is associated with the sampling of available evidence in perceptual decisions and confidence rating at time period of 200-500 ms after stimulus onset (Herding et al., 2019; Rausch, Zehetleitner, Steinhauser, & Maier, 2020; Vafaei Shooshtari et al., 2019; Zizlsperger et al.,

11

372 2014) or at the time of the response (Boldt et al., 2019). Here, CPP amplitude was measured as the
373 mean amplitude in a time-window ranging from 200 ms to 500 ms after stimulus onset in an electrode
374 cluster containing the electrodes CP1, CP2, Cz, and Pz (Boldt et al., 2019; Herding et al., 2019; Rausch
375 et al., 2020; Tagliabue et al., 2019; Twomey, Kelly, & O'Connell, 2016; Vafaei Shooshtari et al., 2019;
376 Zizlsperger et al., 2014). We epoched the EEG responses were aligned with respect to the stimulus
377 onset, from 200 ms pre-stimulus to 500 ms post-stimulus of each pulse. Then, these epochs were
378 baselined to a window -100 ms to stimulus-locked to prevent differences in the visual response to the
379 stimulus affecting the baseline. The ERP signals were examined for levels of confidence separately in
380 double-pulse trials and single-pulse trials. We analyzed correct trials of each coherence level distinctly
381 to support the fact that our findings were not relevant to participants' performance and motion pulse
382 strength. Also, in double-pulse trials, we tested double-pulse trials with non-zero gaps and equal motion
383 strength pulses.

### 2.7.10 Pupillometry analysis

385 Previous work showed that pupil dilation after choice and before feedback reflected decision
386 uncertainty (Colizoli, De Gee, Urai, & Donner, 2018; Urai et al., 2017). Accordingly, as the confidence
387 is uncertainty complement (Hebart, Schriever, Donner, & Haynes, 2014; Kepecs & Mainen, 2012), to
388 study the confidence profile, the method was implemented here. The mean baseline-corrected pupil
389 signal throughout 200 ms before feedback was calculated as our single-trial measure of pupil response.
390 We epoched trials and baselined each trial by subtracting the mean pupil diameter 50 ms before the
391 response. We included all trials of both experiments in the analyses reported in this paper.

392 According to the temporal low-pass characteristics of the slow peripheral pupil apparatus (Hoeks &
393 Ellenbroek, 1993), trial-to-trial variations in response-time can impact trial-to-trial pupil responses,
394 even in the absence of amplitude variations in the underlying neural responses (Urai et al., 2017). To
395 isolate trial-to-trial variations in the amplitude (not duration) of the underlying neural responses, we
396 removed components explained by response-time via linear regression:

$$y' = y - (y^T R)R, \qquad (18)$$

397 where $y$ was the original vector of pupil responses, $R$ was the vector of the corresponding response-
398 time (log-transformed and normalized to a unit vector), and $T$ indicated matrix transpose.
399 Consequently, after removing the variance explained by trial-by-trial response-time, the residual $y'$
400 reflected pupil responses. This residual pupil response was used for analyses reported in this study. To
401 evaluate the relation of confidence and pupil response, we fit a linear regression model as follows:

$$\text{Logit}[P_{high}] = \beta_0 + \beta_1 P, \qquad (19)$$

402 where the $P$ was pupil response in each trial. The null hypothesis was that confidence did not change
403 with the pupil response ($H_0: \beta_1 = 0$). To control the impact of coherence on pupil response, we only
404 included motion strength of 3.2, 6.4, 12.8 of single-pulse trials.

### 2.7.11 General statistical analysis

406 We used repeated-measures two-tailed $t$-tests. As suggested, we considered small (d = .2), medium (d
407 = .5), and large (d = .8) effect sizes for this assessment (see (Cohen, 1970)) and the statistical
408 significance for $t$-tests was set to a probability from data ≥ .90.

This is a provisional file, not the final typeset article

409  Moreover, to test our hypotheses, a series of regression analyses were run after confirming the
410  assumptions of the linear regression are met. Effect sizes were reported and as suggested, here, we
411  considered small ($f^2 = .02$), medium ($f^2 = .15$), and large ($f^2 = .35$) effect sizes (see (Cohen, 1970)) at
412  the alpha level of 5%.

413  For tests of pupil response signals and ERPs between two levels of confidence, statistical inferences
414  were performed using $t$-tests at each time-point (at a statistical threshold of $p < .05$).

## 3    Results

416  We tested our predictions in two studies that applied the same paradigm (**Figure 1A**). The first study
417  used behavioral measures and pupillometry analyses, whereas for the second experiment, we recorded
418  EEG signals as well. Participants decided about the direction of the RDM motion based on brief motion
419  pulses. The task design contained different conditions which allowed us to compare participants'
420  behavior in (i) double-pulse vs single-pulse, (ii) different coherence of motion stimulus, and (iii) four
421  distinct gaps intervals.

### 3.1    Behavioral results

423  We used the single-pulse trials to benchmark the effect of coherence on choice accuracy and
424  confidence. As shown in **Figure 2A**, for single-pulse trials, participants were more confident for high
425  coherence stimuli (**Figure 2A**; Eq.1; $\beta_1 = .06$, $p < .001$, 95% CI = [.04, .08], $f^2 = .23$), ranged from .96
426  for 51.2 coherence to .43 for 3.2 coherence. Also, accuracy improved with motion strength reached
427  from .56 for 3.2% to .99 for 51.2% (**Figure 2B**, black line; Eq.2; $\beta_1 = .10$, $p < .001$, 95% CI = [.08,
428  .12], $f^2 = .36$). They also had better performance whenever they had reported higher confidence
429  comparing to lower confidence (**Figure 2B**, red and green; Eq.3; $\beta_1 = 1.1$, $p < .001$, 95% CI = [.88,
430  1.31], $f^2 = .09$). Moreover, in double-pulse trials, the accuracy improved with motion strength (**Figure
431  2C**, black dots) and participants were more accurate while reporting higher confidence (**Figure 2C**,
432  green dots). Along with accuracy (**Figure 2D**; Eq.5; $p > .1$, (Kiani et al., 2013; Tohidi-Moghaddam et
433  al., 2019), see **Supplementary Figure 1A** for Experiment 2 data), the confidence was largely
434  unaffected by interpulse interval in both double-pulse trials with equal pulse strength (**Figure 2E**; Eq.4;
435  $p > .1$; **Supplementary Table 2**, results of individual participants) and those with unequal pulse
436  strength (**Figure 2E**; Eq.4; $p > .1$; **Supplementary Table 2**; see **Supplementary Figure 1B**  for
437  Experiment 2 data). The two pulses separated by up to 1 s supported a level of confidence that was
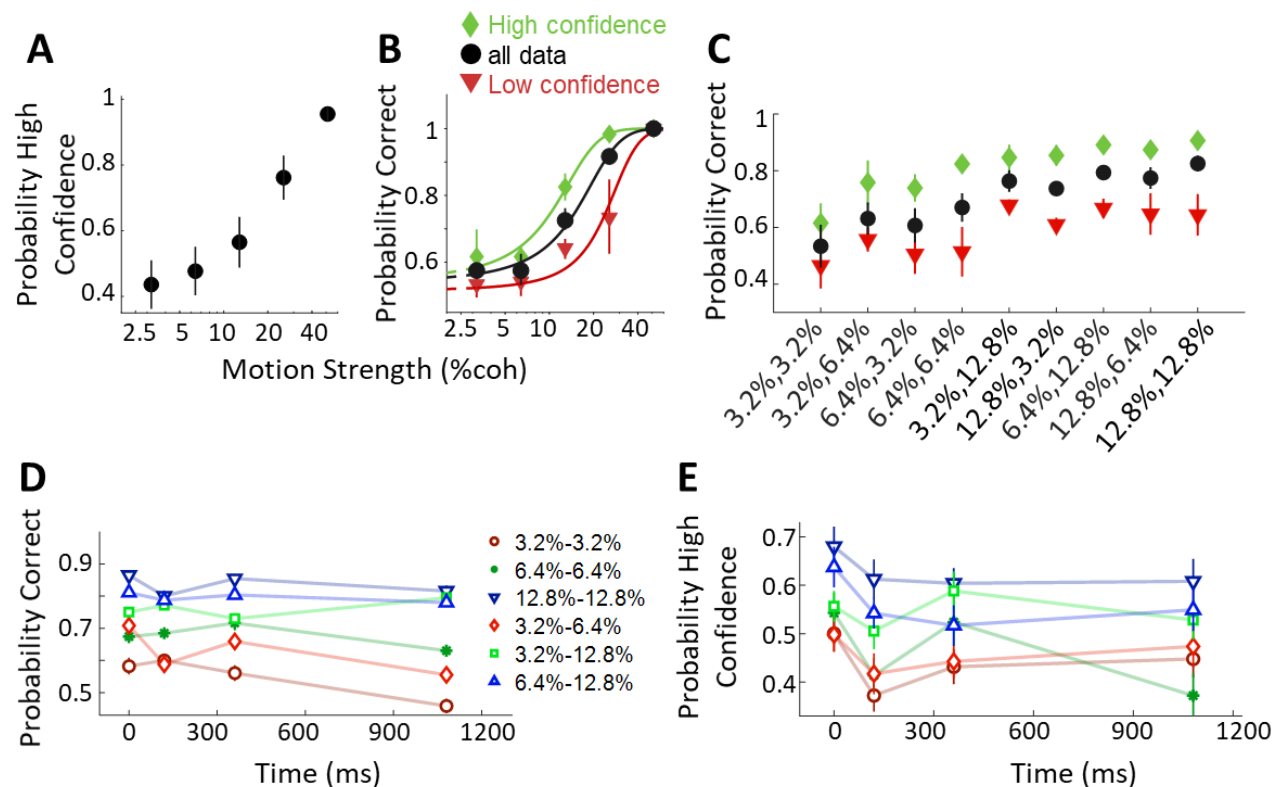438  indistinguishable from a pair of pulses separated by no gap.

439

**Figure 2. Interplay between confidence, accuracy, and coherence in single/double-pulse trials, and interpulse interval in double-pulse trials. (A)** Probability of high confidence as a function of motion coherence. **(B) (C)** Accuracy in single-pulse trials and double-pulse trials in all trials (black), split by high (green) and low (red) confidence decisions. In (B) curves are model fits. **(D)** Choice accuracy for double-pulse trials grouping in all possible interval conditions. **(E)** Confidence of double-pulse trials was calculated by pooling data across all time intervals. In (D) and (E) each data point reports pooled data from indicated sequence pulse and its reverse order (e.g., 12.8– 3.2% and 3.2 –12.8%).

440

441 Direct comparison between single-pulse and double-pulse trials, along with previous studies (Kiani et
442 al., 2013; Tohidi-Moghaddam et al., 2019), showed that participants' accuracy significantly differed
443 (t(11490) = -3.09, *p* < .05, 95% CI = [-.08, -.02], Cohen's d = .11). However, in double-pulse trials
444 participants were not more confident comparing to single-pulse trials (t(11490) = 1.35, *p* = .18, 95%
445 CI = [-.01, .06], Cohen's d = -.05).

446 Although, the order of the pulses affected accuracy (**Figure 3A**, (Kiani et al., 2013; Tohidi-
447 Moghaddam et al., 2019)), participants were not more confident in double-pulse trials with unequal
448 pulse strength where the stronger motion appeared in a second order (**Figure 3B**; Eq. 6; $\beta_2$ = .01, p =
449 .08, 95% CI = [.00, .02] ], $f^2$ = .02; see **Supplementary Figure 2B** for Experiment 2 data). Also, the
450 increased confidence was not because of an interaction of motion pulses (Eq. 7; $\beta_3$ = -.01, p = .13, 95%
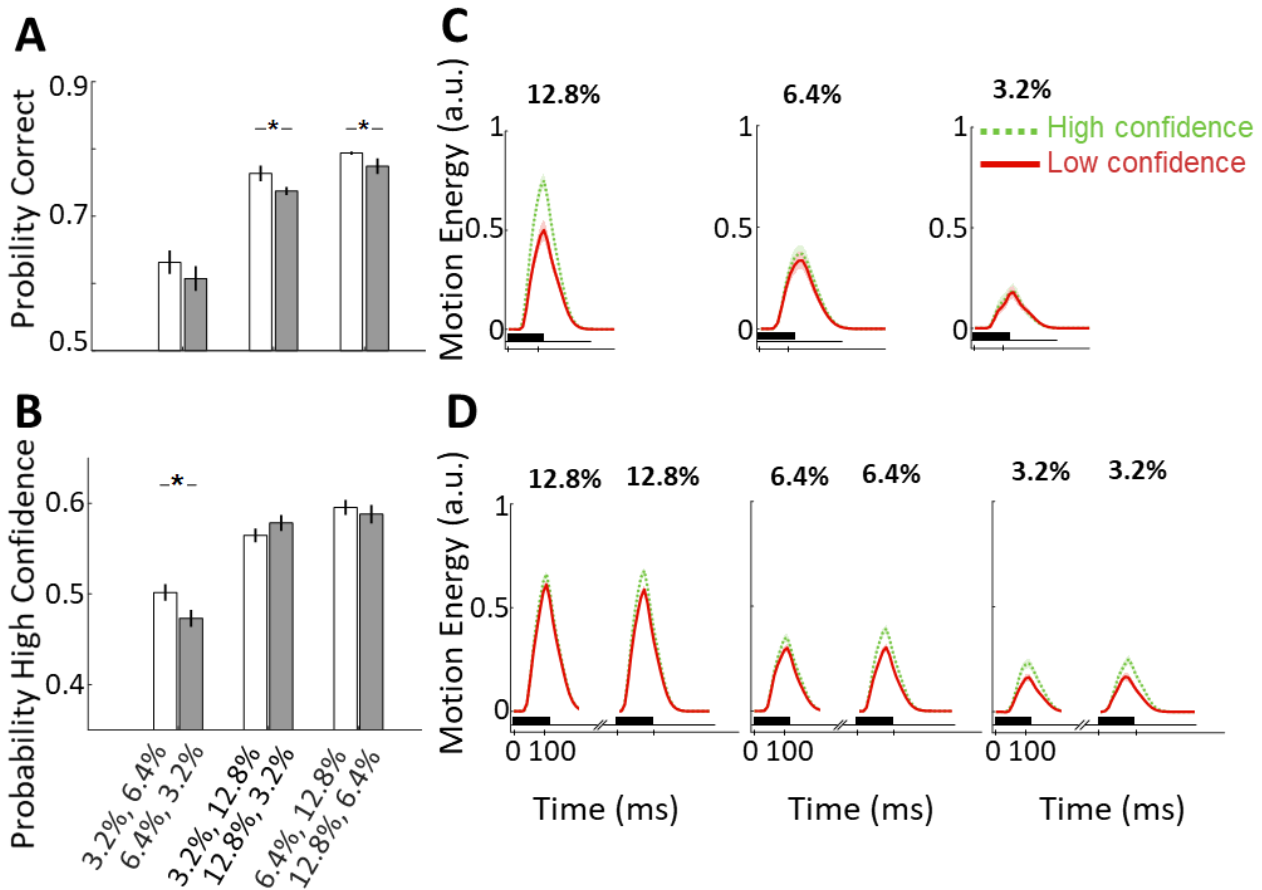451 CI = [-.02, .00], $f^2$ = .03).

**Figure 3. Choice confidence was not depended of the sequence of motion pulses (A)** The weak–strong pulse sequence contributed higher accuracy than the strong–weak sequence. **(B)** The weak–strong pulse sequence did not contribute higher confidence than the strong–weak sequence. In all panels, data are represented as group mean ± SEM. (*p<0.05) **(C)** In single-pulse trials, low and high confidence cannot be determined by motion energy profiles in weaker pulses **(D)** The second pulse had slightly more impact on confidence. Data were pooled for all nonzero interpulse intervals. Only correct trials with equal pulse strength are included. In (C) and (D), the shaded region around the mean indicates SEM. The black horizontal bars show the duration of the stimulus display. The units of motion energy are arbitrary and the same for all motion strengths.

452

### 3.2 Motion energy results

To yield a precise estimate of the decision-relevant sensory evidence accommodated in the stochastic stimuli, we employed motion energy filtering to the random dot motion stimuli. **Figure 3D** displays the average motion energy in double-pulse trials when the strength of pulses was the same. Accordingly, the difference of the motion energy profiles for high and low confidence responses was slightly larger for the second pulse than the first pulse. A logistic regression confirmed the influence of trial-to-trial fluctuations of motion energy on confidence (Eq.11; $\beta_2 = .10$, $p = .001$, 95% CI = [.04, .16], $f^2 = .19$). Also, there was slightly larger impact of motion energy of the second pulse with equal pulse strength (Eq.11; $\beta_3 = .11$, $p = .04$, 95% CI = [.07, .15], $f^2 = .13$). On the contrary, the impact of motion energy of the second pulse was not significant (Eq.12, $\beta_4 = .10$, $p = .06$, 95% CI = [.06, .14], $f^2$

463  = .08). Consequently, motion energy analysis could not provide independent confirmation of
464  asymmetric effect of both pulses for confidence.

465  As well, in single-pulse trials, the difference of the motion energy profiles for high and low confidence
466  with stronger pulse strength (12.8%, 6.4%) was significant (**Figure 3C**; Eq.13; $\beta_2$ = .41, $p$ = 2.25 × 10$^{-5}$
467  , 95% CI = [.23, .59], $f^2$ = .24). However, the difference in weak motion pulse was insignificant
468  (**Figure 3D**; Eq.13; $\beta_2$ = .17, $p$ = .44, 95% CI = [-.26, .60], $f^2$ = .10). Thus, motion energy analysis thus
469  suggests that when the pulses' motion strengths are weak, the subjects decide about their confidence
470  almost randomly.

## 3.3    The Interplay between confidence in single vs double-pulse trials

472  To address accuracy and confidence variation in double-pulse from single-pulse trials, we consider
473  $P_{correct}$ or $P_{high}$ of each coherence (3.2%, 6.4% and 12.8%) in single-pulse as baseline and measure
474  the $P_{correct}$ or $P_{high}$ variation of any corresponding sequence in double-pulse trials. As we expected
475  in all combinations of three coherence as the baseline, $P_{correct}$ improved (**Figure 4A**). Additionally,
476  when considering all the trials, in all combinations of three coherence as the baseline, $P_{high}$ increased
477  when the other pulse was a strong pulse (12.8%) (**Figure 4B** and **Figure 4C** for correct trials). On the
478  contrary, $P_{high}$ decreased or not changed considerably whenever the other pulse was a weak motion
479  strength (3.2%, 6.4%). Interestingly, in incorrect trials, the confidence decreased comparing to single-
480  pulse for all the coherence and conditions (**Figure 4D**). These data did not correlate with the interval
481  duration (**Figure 4A**, **B**, **C**, **D**).

This is a provisional file, not the final typeset article

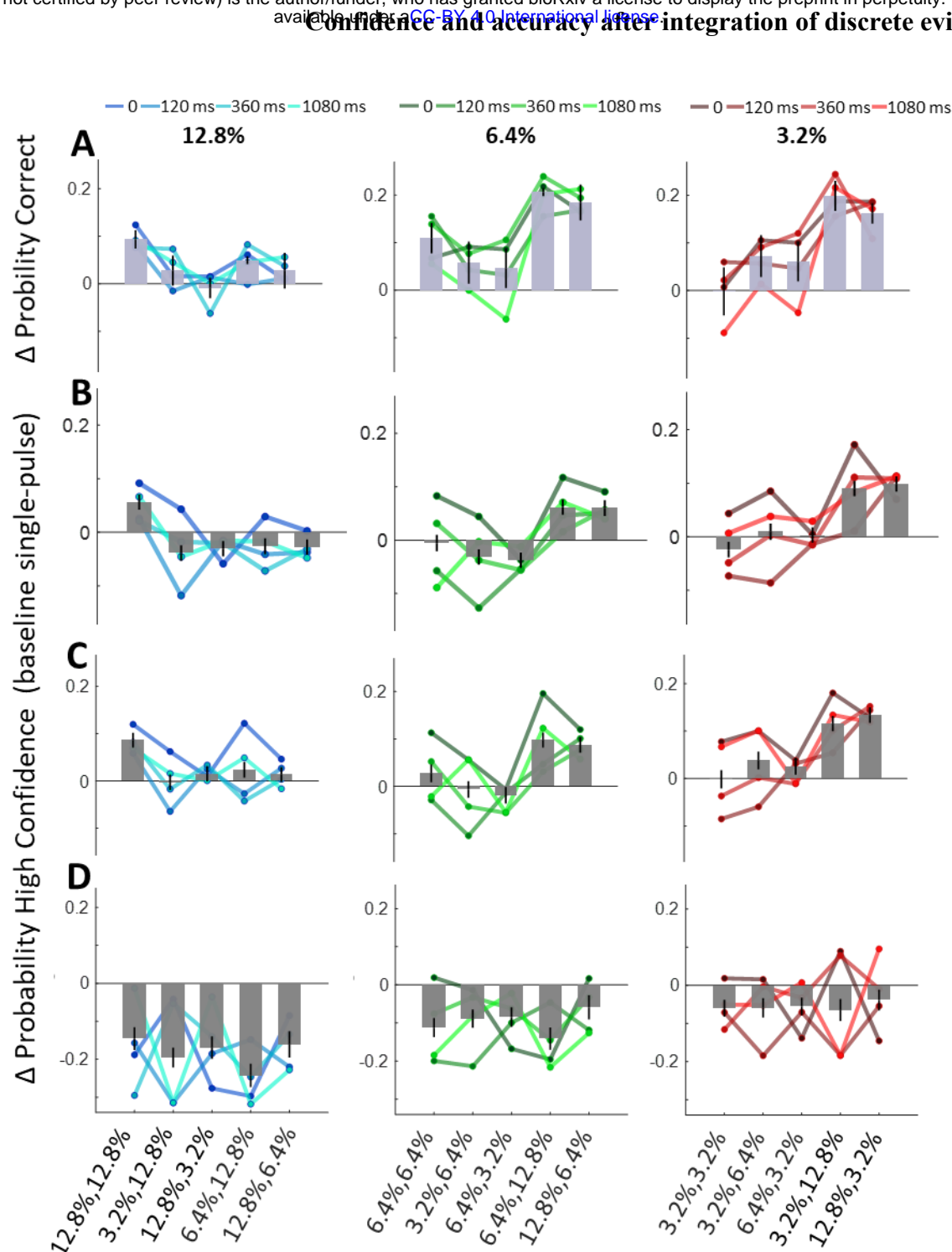Confidence and accuracy after integration of discrete evidence



**Figure 4. Variation of accuracy or confidence in double-pulse trials baselined by corresponding coherence (3.2%, 6.4% and 12.8% for each column). (A)** Considering all the trials, the accuracy improved in almost all pulses combination. **(B)** Considering all the trials, the confidence improved in combination with stronger pulses while the confidence in sequence with a weaker pulse either decreased or remained constant. **(C)** In correct-choice trials, the increasing effect of stronger pulses is more significant and the confidence even slightly improved in combination with weaker pulses comparing to corresponding baseline. **(D)** Interestingly, in incorrect trials, the confidence decreased in every condition. The colored line representing matching data for each of four possible gaps. The data are represented as group mean ± SEM.

482

483 In other words, the participants reported lower confidence in double-pulse trials compared to single-
484 pulse trials for incorrect choices but reported higher confidence for correct choices (**Figure 4** Eq.8, $\beta_1$
485 = .15, $p$ < .001, 95% CI = [.13, .17], $f^2$ = .29; **Supplementary Figure 3** for Experiment 2;
486 **Supplementary Table** 1, results of individual participants). This data is in line with the fact that the
487 good metacognitive sensitivity will provide higher confidence for correct responses, and lower for
488 incorrect ones.

489 ### 3.4 Computational models

490 The accuracy in double-pulse trials surpasses the expectation measured by the perfect integrator
491 (**Figure 5A**; (Kiani et al., 2013)). Considering the strong positive relation of accuracy and confidence
492 (Kiani et al., 2014; Vafaei Shooshtari et al., 2019), we expected the observed confidence would exceed
493 the predicted confidence (Eq.16) calculated by the perfect integrator model (Eq.14), but it did not
494 (**Figure 5B**).

495 According to SDT models, $d'$ is stimulus sensitivity and has relation to task performance. As the $d'$ of
496 the perfect integrator model was calculated based on single-pulse trials performance, if participants'
497 performance in single-pulse trials failed, their performance prediction missed the double-pulse trials
498 (**Figure 5C**).

499 $Meta\text{-}d'$ in all participants increased in double-pulse trials but perfect integrator model failed to imitate
500 the increasing (**Figure 5C**). We also computed metacognitive efficiency ($Meta\text{-}d'/d'$), as another
501 index of the ability to discriminate between correct and incorrect trials. Here, $Meta\text{-}d'/d'$ in all
502 participants missed to track their $Meta\text{-}d'/d'$ in double-pulse trials. Altogether, the perfect integrator
503 was incapable of employment observed metacognitive ability in double-pulse trials. The same
504 modeling procedure of data from EEG experiment has provided similar results (**Supplementary
505 Figure 4**).

506 As a control investigation, we examined whether the differences in estimated metacognitive ability
507 between models could result from the different number of trials. We averaged the metacognitive scores
508 obtained from equal numbers of samples, and found very similar results. Thus, the difference in the
509 estimated metacognitive efficiency cannot be explained by the difference in the number of trials
510 between the single-pulse, double-pulse, and perfect integrator models (**Supplementary Figure 7**).
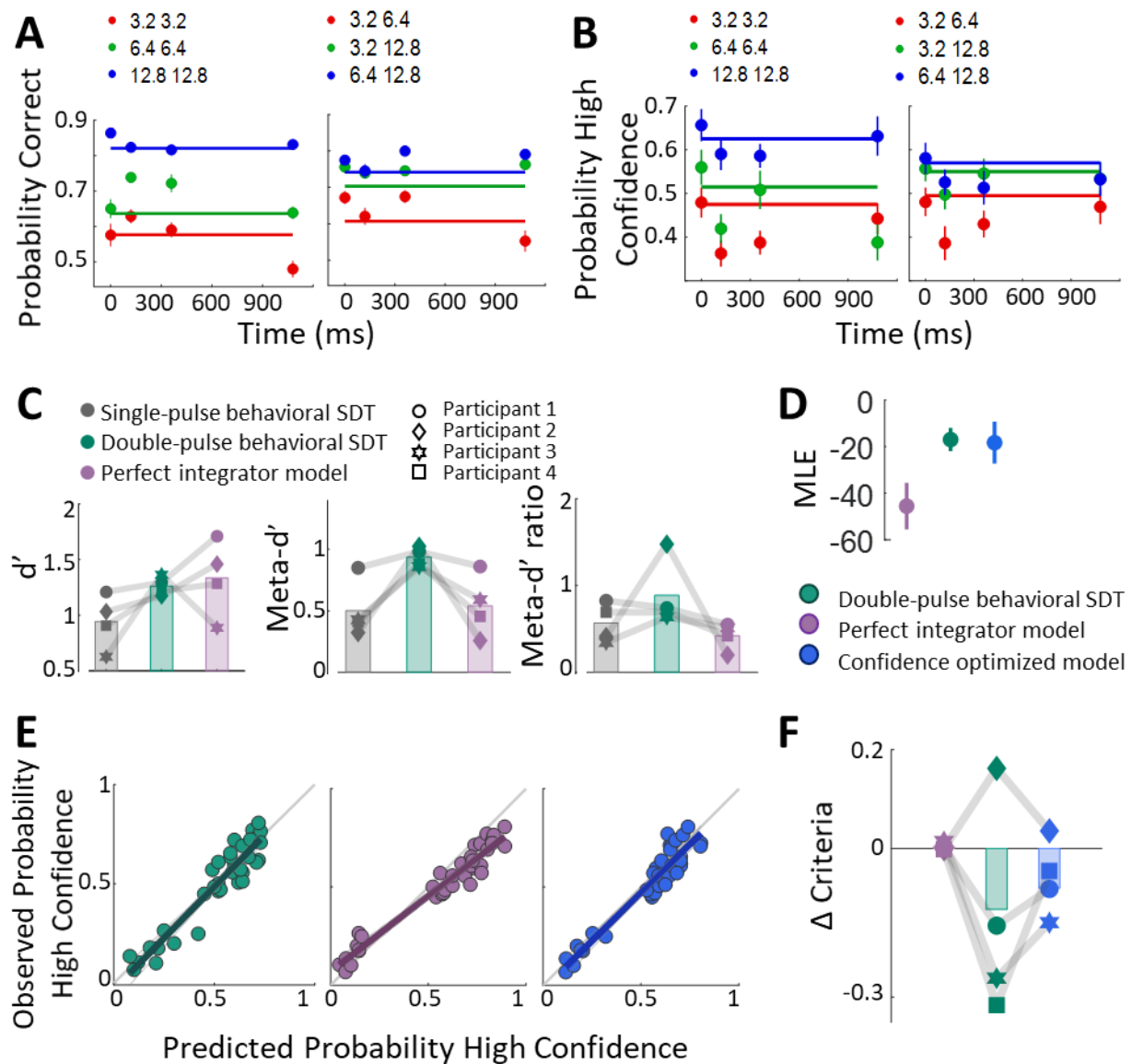
511

This is a provisional file, not the final typeset article

**Confidence and accuracy after integration of discrete evidence**



**Figure 5. Comparison of the models and human behavior. (A)** Accuracy in double-pulse trials. Horizontal lines show accuracy prediction by the perfect integrator model. **(B)** Confidence in double-pulse trials. Horizontal lines show confidence prediction by the perfect integrator model. In (A) and (B) Each data point represents pooled data from the pulse sequence indicated by the legend and its reverse order. **(C)** Stimulus sensitivity ($d'$), metacognitive sensitivity ($Meta-d'$), metacognitive efficiency ($Meta-d'/d'$) estimated for single-pulse trials, double-pulse trials and perfect integrator models for each participant. **(D)** Model comparison suggests strong evidence in favor of the confidence optimized model over the perfect integrator **(E)** Relation of predicted confidence and observed data. SDT model fitted to double-pulse trials (green), the perfect integrator model (purple), and optimized model (blue). Colored lines indicate best-fitting slope of a linear regression analysis. Each data point represents pooled data from different sequence of pulses of each participant. **(F)** Variation of confidence criteria comparing to single-pulse trials in perfect integrator vs optimized model. For panels A, B and, D, data are represented as group mean ± SEM.

512

19

513 As in the perfect integrator model, the $\beta_1$ (slope in Eq.16) differed from 1:1 line and confidence
514 prediction failed to account for behavioral data (**Figure 5E**; Eq.16, $\beta_1 = .77$, $p < .001$, 95% CI = [.71,
515 .83], $f^2 = 15.66$), we introduced a model in which the metacognitive sensitivity (*Meta–d'*) calculated
516 in the perfect integrator model was optimized. The stimulus parameter ($d'$) remained constant whereas
517 the placement of confidence criteria was optimized to fit best to observed data. So, the predicted $P_{high}$
518 improved intensely (**Figure 5E**, Eq.16, $\beta_1 = .98$, $p < .001$, 95% CI = [.88, 1.08], $f^2 = 10.11$).
519 Additionally, we take the confidence criteria of the single-pulse model as the baseline and measure the
520 variation of criteria of the perfect integrator and the optimized model. This variation in the optimized
521 model has changed comparing to the perfect integrator (**Figure 5F**, and **Supplementary Figure 5** for
522 EEG experiment). Failure to predict the proper change in confidence criteria in the perfect integrator
523 model was the factor that made the model unable to estimate the confidence from single-pulse trials.

524 In addition, to consider the suboptimality in confidence reporting, we simulated data using the perfect
525 integrator model's parameters while setting higher confidence noise (Eq.15). The predicted $P_{high}$ from
526 this simulation improved (Eq.16, $\beta_1 = .97$, $p < .001$, 95% CI = [.83, 1.07], $f^2 = 9.00$). Consequently, the
527 perfect integrator model simply highlighted accumulating decision evidence and ignored the effect of
528 confidence noise.

### 529 3.4.1 Models' evaluation

530 We conducted parameter recovery simulations to evaluate models fitted to single/double-pulse trials.
531 We regressed predicted vs. observed $P_{high}$ confidence for each coherence of each participant. In single-
532 pulse trials, linear regression indicated that there was a significant effect between the predicted and
533 observed $P_{high}$, (Eq.16, $\beta_1 = 1.04$, $p < .001$, 95% CI = [.90, 1.18], $f^2 = 8.09$). In double-pulse trials,
534 regression coefficient was statistically significant and close to 1:1 line (**Figure 5E**; Eq.16, $\beta_1 = 1.03$, $p$
535 $< .001$, 95% CI = [.91, 1.15], $f^2 = 11.05$) meaning predicted $P_{high}$ by classic SDT also explained a
536 significant proportion of variance in the observed $P_{high}$.

537 A quantitative model comparison unsurprisingly favored the optimized model (mean MLE = -18.31)
538 and the SDT behavioral model (mean MLE = -16.98) over the perfect integrator model (mean MLE =
539 -45.53) (**Figure 5D**).

540 In summary, comparing between the models, both quantitatively (**Figure 5E**) and qualitatively (**Figure
541 5D**) in doube-pulse trials, also showed that the confidence optimized model has a better prediction in
542 estimating confidence. Accordingly, these investigations indicated: (i) participants integrated the
543 decision evidence perfectly but to report their confidence, their confidence resolution improved rather
544 than reporting higher confidence, (ii) the inability to predict the proper change in confidence criteria in
545 the perfect integrator model was the factor that made the model unable to estimate the confidence from
546 single-pulse trials, (iii) the confidence noise was changed after receiving the second pulse in double-
547 pulse trials.

### 548 3.5 Response-time analysis

549 Response-time had a significant effect on confidence in double-pulse (Eq.9, $\beta_1 = .09$, $p = .04$, 95% CI
550 = [0.01, 0.17], $f^2 = .10$) but not in single-pulse trials (Eq.9, $\beta_1 = - .02$, $p = .90$, 95% CI = [-1.78, 1.74],
551 $f^2 = .00$). Moreover, the confidence profile as a function of response-time was significant in double-
552 pulse trials (**Figure 6A; Table 1**) but not in single-pulse trials (**Figure 6B; Table 1**).

553

This is a provisional file, not the final typeset article

554 *Table 1. Result from t-tests to compare confidence profile in single/double-pulse trials in each response-time bin.*

| Trial type | Double-pulse | | | | Single-pulse | | | |
|---|---|---|---|---|---|---|---|---|
| RT bin | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 |
| *tstat* | 8.89 | 6.77 | 9.20 | 8.10 | 1.00 | 2.71 | 1.46 | 1.69 |
| *df* | 2652 | 2652 | 2652 | 2652 | 2150 | 2150 | 2150 | 2152 |
| *p* | $.11 \times 10^{18}$ | $.16 \times 10^{11}$ | $.71 \times 10^{20}$ | $.77 \times 10^{16}$ | .31 | .007 | .14 | .09 |
| 95% CI | [-.17, -.11] | [-.15, -.08] | [-.19, -.12] | [-.19, -.12] | [-.18, .06] | [-.30, -.05] | [-.22, .03] | [-.24, .02] |
| Cohen's *d* | .35 | .26 | .36 | .32 | .14 | .37 | .20 | .23 |

555

556 Additionally, in our double-pulse trials, participants decided faster than single-pulse trials in all interval
557 durations (**Figure 6C**). We regress the delay-time before cue onset (0.4 to 1 s truncated exponential)
558 and response-time in both single-pulse and double-pulse trials to examine the effect of imposed delay
559 time on response-time. The effect was small in both single-pulse (Eq.10, $\beta_1 = -.01 \times 10^{-5}$, $p = .004$,
560 95% CI = [0.00, 0.00], $f^2 = .00$) and double-pulse trials (Eq.10, $\beta_1 = -.0003 \times 10^{-6}$, $p < .001$, 95% CI =
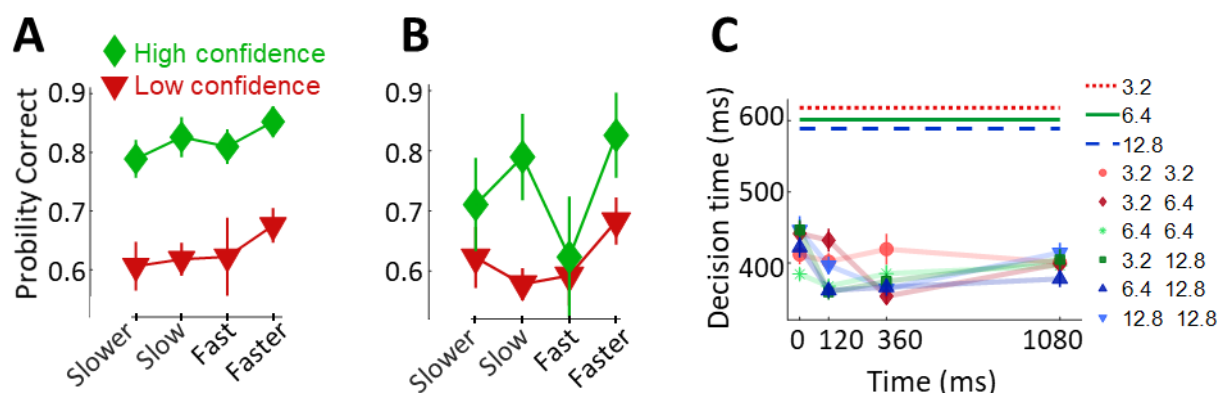561 [0.00, 0.00], $f^2 = .01$).

562



**Figure 6. Response-time profiles in single and double-pulse trials. (A) (B)** Accuracy as a function of response-time split by high (green) and low (red) confidence in double-pulse trials (A) and single-pulse trials (B). **(C)** Response-time of all coherence combination clustered by gap interval in double-pulse trials (dots) comparing to single-pulse trials (lines). Data are represented as group mean ± SEM.

563

564 **3.6 EEG Analysis**

565 We derived the ERPs of averaged signals for two levels of confidence to verify whether there was a
566 significant difference in the centro-parietal ERPs across confidence levels. **Figure 7** exhibit ERPs and
567 scalp topographies for confidence levels time-locked to the stimulus onset in low and high confidence
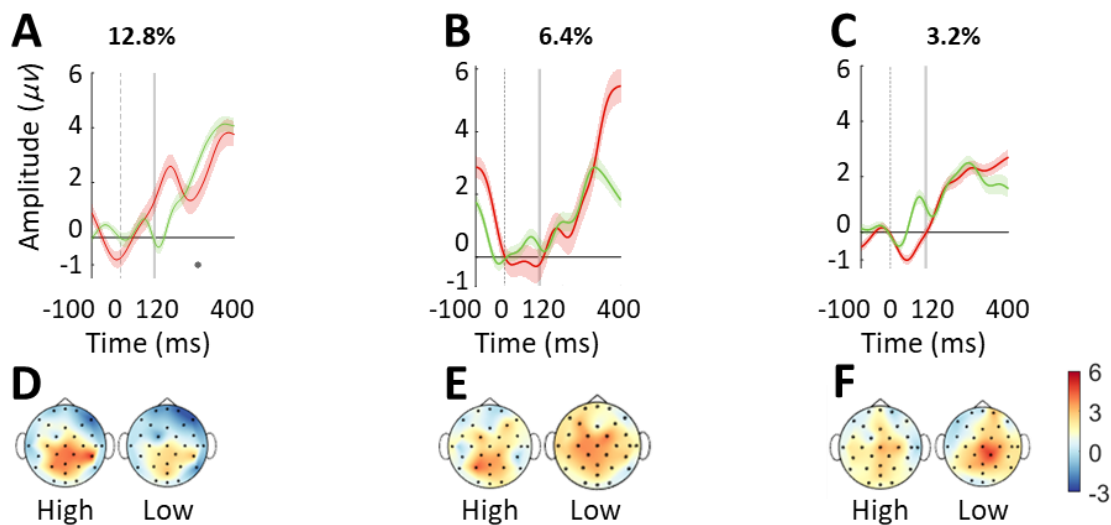568 in single-pulse trials.

21

**Figure 7. ERPs and scalp topographies in single-pulse trials. (A) (B) (C)** ERPs of correct single-pulse trials shows an insignificant difference in weaker motion strength in high and low confidence level trials. **(D) (E) (F)** Scalp topographies in two levels of confidence (the mean amplitude in a time-window ranging from 200 ms to 500 ms after stimulus onset). The shading region around the mean indicates SEM. * indicate $p<.05$ from a *t*-test, of the difference between the two-time.

569

570    **Figure 8** exhibit ERPs and scalp topographies for confidence levels time-locked to the stimulus onset
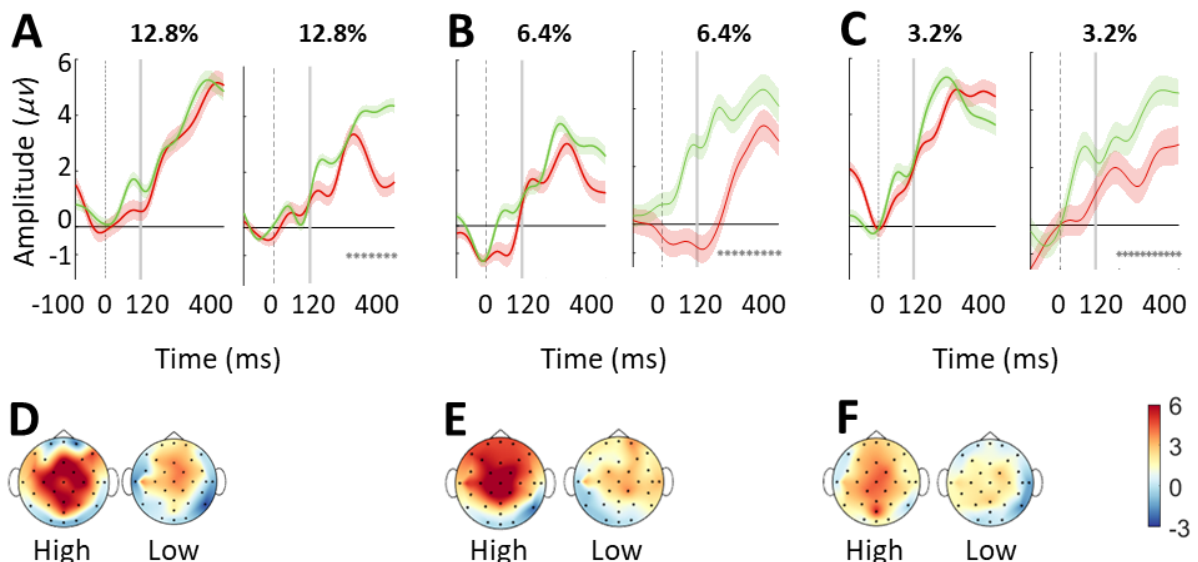571    in low and high confidence in double-pulse trials.

572



**Figure 8.  ERPs and scalp topographies in double-pulse trials. (A) (B) (C)** ERPs in the two levels of confidence are distinct after the stimulus onset. **(D) (E) (F)** Scalp topographies in two levels of confidence (the mean amplitude in a time-window ranging from 200 ms to 500 ms after stimulus

onset of second pulse). The shading region around the mean indicates SEM. * indicate p<.05 from a *t*-test, of the difference between the two-time.

573

574 Interestingly the effect of different confidence profiles in centro-parietal was considerable in double-
575 pulse trials (**Figure 8**) but not in single-pulse trials (**Figure 7**).

576 **3.7 Pupil responses**

577 We took the mean baseline-corrected pupil signal during 200 ms before feedback delivery as our
578 measure of pupil response. In line with previous work (Urai et al., 2017) pupil responses reflect
579 decision confidence in our double-pulse trials (**Figure 9A**; Eq.19, $\beta_1 = -.95$, $p < .001$, 95% CI = [-1.10,
580 -.79] , $f^2 = .31$) while in single-pulse trials the confidence profile is not significant (**Figure 9B**; Eq.19,
581 $\beta_1 = -.42$, $p = .21$, 95% CI = [-1.09, -.24], $f^2 = .22$).
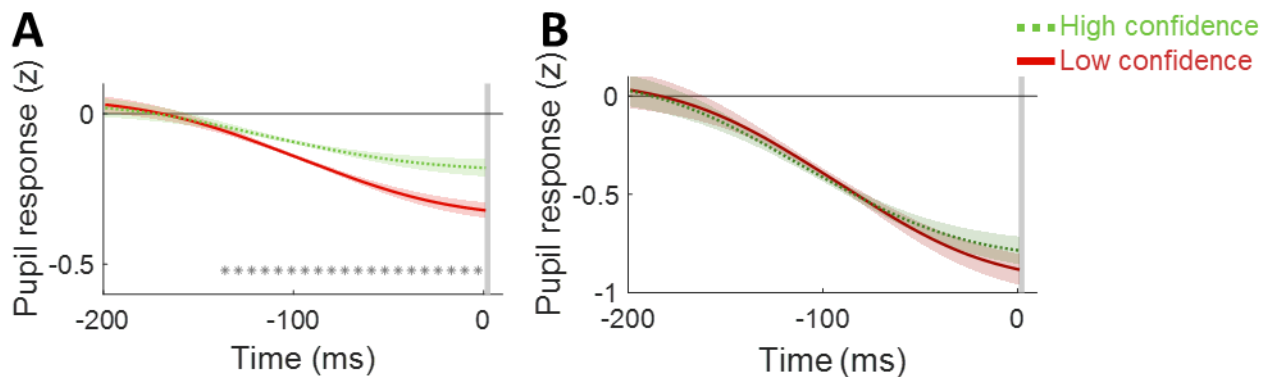
582



**Figure 9. Standardized pupil response across time-window aligned to the feedback. (A) (B)** Standardized pupil response, high confidence trials (green) vs low confidence trials (red) in double-pulse trials (A) and single-pulse trials (B). The shading region around the mean indicates SEM. * indicate p < .05 from a *t*-test.

583

584 **4 Discussion**

585 The current study was designed to clarify the confidence of decisions in more real-world contexts
586 where the evidence arrives separately. Using an experimental design, we examined how human
587 subjects combined the pieces of information to form their decision and confidence and how the two
588 are related to each other. We performed two experiments with either single or double pulses of RDM
589 stimuli. To this end, we investigated behavioral modeling, EEG responses and pupillometry. In
590 summary, the results across experiments showed that participants used both pulses to decide about their
591 confidence. Also, while their confidence was largely invariant to the gap interval, confidence scoring
592 was not noticeably enhanced in double-pulse trials compared to single-pulse trials. Instead, participants
593 reported their confidence with higher resolution and their metacognitive sensitivity improved.
594 Furthermore, using RT, EEG and pupillometry analysis, we could considerably track the confidence
595 profiles in double-pulse trials, unlike in single-pulse trials.

## 4.1 Behavioral and motion energy findings

Remarkably, unlike accuracy, confidence ratings in double-pulse trials have not increased significantly comparing to single-pulse trials. We hypothesize that participants mainly trust on the evidence of one of the pulses and ignore the other one. The trusted pulse can either be the first or second pulse; it also can simply be the stronger pulse. However, the effect of sequence and interaction of pulses on confidence was examined and no effect was observed. Moreover, motion stimulus fluctuations are known to influence the choice (Kiani et al., 2008; Resulaj et al., 2009) and confidence (Van Den Berg et al., 2016; Zylberberg et al., 2012), so they can inform us about the parts of the stimulus that bear more intensely on the choice and confidence (Kiani et al., 2013, 2008; Nienborg & Cumming, 2009). The motion energy analysis could not confirm the asymmetric influences of the two pulses for confidence. However, the motion energy analysis does provide independent confirmation of the unequal influences of the two pulses for choice (Kiani et al., 2013). Since, pervious research suggests that participants obtained more information from a second pulse (Kiani et al., 2013; Tohidi-Moghaddam et al., 2019), we hypothesized, in line with a large body of evidence (De Gardelle & Mamassian, 2015; Herce Castañón et al., 2019; Rahnev & Denison, 2018; Zylberberg et al., 2016, 2014), here observers do not make their decisions exactly in accordance with confidence rating.

Moreover, once comparing confidence in double-pulse trials grouped by accuracy, we show that the participants had lower confidence in double-pulse trials than single-pulse trials for incorrect choices but higher confidence for correct choices. In other words, compared with single-pulse trials, in double-pulse trials, participants adjusted their confidence by enhancing their confidence resolution or metacognitive sensitivity.

Typically, confidence facilitates evidence accumulation and drives a confirmation bias in perceptual decision-making (Rollwage et al., 2020). Likewise, we suggest that an extra brief and weak evidence can validate confidence and improve metacognitive sensitivity.

## 4.2 Computational modeling findings

To understand the nature of the differences in participants' metacognitive sensitivity in double-pulse vs single-pulse trials, we compared corresponding estimated metacognitive parameters. Likewise, we included the expected parameters that would be achieved in double-pulse trials under the assumption of perfect integration. Accordingly, we computed $Meta\text{-}d'/d'$ as a measure of 'metacognitive efficiency'. In the case of $Meta\text{-}d' = d'$, the observer is metacognitively 'ideal'. Indeed, all the information available for the decision would be used to report the confidence. Yet, in many cases, we might find that $Meta\text{-}d' < d'$, along with some degree of noise or suboptimality (Fleming & Lau, 2014; Maniscalco & Lau, 2012). Conversely, we may find that $Meta\text{-}d' > d'$, if subjects are able to draw on additional information such as hunches (Rausch & Zehetleitner, 2016), further processing of stimulus information (Charles, Van Opstal, Marti, & Dehaene, 2013) or extra prior knowledge on the task. In the model fitted to double-pulse trials, $Meta\text{-}d'/d'$ was around .8 and near to ideal for almost all participants. However, as in single-pulse trials, it varies considerably between participants, the value could not be adjusted in perfect integrator model similar to the behavioral model.

Previously, the better-than-expected performance in double-pulse trials was explained by underperformance in single-pulse trials (Kiani et al., 2013). Here, metacognitive sensitivity in double-pulse trials surpasses the value predicted by the perfect integrator model (**Figure 5C** and **Supplementary Figure 4**). This effect can be followed in all of our participant (except one of

638 participants from EEG experiment) and can be explained by low confidence resolution in single-pulse
639 trials.

640 Metacognitive noise is the noise that affects confidence estimates but not perceptual decisions (De
641 Martino, Fleming, Garrett, & Dolan, 2013; Jang, Wallsten, & Huber, 2012; Maniscalco & Lau, 2016;
642 Mueller & Weidemann, 2008; Rahnev, Nee, Riddle, Larson, & D'Esposito, 2016; Shekhar & Rahnev,
643 2018; Van den Berg, Yoo, & Ma, 2017). A recent work categorized sources of metacognitive
644 inefficiency (Shekhar & Rahnev, 2020). Accordingly, metacognitive noise is a superordinate term for
645 all noise sources that impact the confidence formation process (Shekhar & Rahnev, 2020, 2021)
646 ranging from systematic to nonsystematic input and computation. Nevertheless, the exact source of
647 metacognitive noise remains unclear (Shekhar & Rahnev, 2020). This noise can be tracked in our
648 perfect integrator model, which was capable of accumulating decision evidence perfectly but could not
649 predict confidence formation in our task. We suggest that the perfect integrator model was unable to
650 adjust to confidence criteria when predicting confidence in double-pulse trials. However, an improved
651 SDT capable of addressing metacognitive noise might be able to empower the employed perfect
652 integrator model. Furthermore, SDT is not the only available model to implement a perfect integrator
653 model. Previous studies suggested attractor models as a candidate model to implement the perfect
654 integrator model (Kiani et al., 2013; Waskom & Kiani, 2018). Attractor models are a group of networks
655 that formed a bridge between cognitive theory and biological data which exploits inhibition to achieve
656 a competition among alternatives (Wang, 2002; Wong & Huk, 2008). Although these models can
657 integrate momentary evidence to establish a decision, they have specific failure behaviors that would
658 be apparent when the sources of evidence are separated by gaps in time (Kiani et al., 2013). Besides,
659 when the stimulus is very short, mostly, none of the attractors could be reached and, the network would
660 revert back to the resting state after the stimulus offset (Wang, 2002). Therefore, the choice would be
661 assigned randomly. However, our experiments' data represent a noteworthy performance in single-
662 pulse trials, which does not support this expectation. Consequently, to implement a perfect integrator
663 model by implementing an attractor model, a mechanism for simulating a very short stimulus might be
664 considered. Moreover, our behavioral assays highlighted different relationships between confidence
665 and accuracy in the different conditions of the task. So, a dedicated neural module with a plausible
666 circuit of confidence might be a better option to implement a perfect integrator model. Recently, multi-
667 layer recurrent network models has been developed to account for decision confidence mechanisms
668 (Atiya, Rañó, Prasad, & Wong-Lin, 2019; Paz, Insabato, Zylberberg, Deco, & Sigman, 2016). These
669 models consist of multiple layers of neural integrators and in line with neural evidence of decision
670 confidence (Kepecs, Uchida, Zariwala, & Mainen, 2008; Murphy, Robertson, Harty, & O'Connell,
671 2015), they are suggested to justify the observed behavior.

672 Furthermore, perceptual decisions are often modeled using ideal observers (e.g., SDT). However, a
673 source of suboptimal behavior in decision-making is 'lapse' (Gold & Ding, 2013; Pisupati,
674 Chartarifsky-Lynn, Khanal, & Churchland, 2021). Lapses are an additional constant rate of errors
675 independent of the evidence strength (Gold & Ding, 2013; Pisupati et al., 2021). Lapse rate has been
676 shown to increase with higher perceptual uncertainty (Pisupati et al., 2021) and would be accounted
677 by fitting extra parameter to psychometrics models. Accordingly, as the perfect integrator model was
678 based on SDT, ignoring lapse in the single-pulse trials might lead to mis-estimation of decision
679 parameters in double-pulse trials. Consequently, further models including the lapse parameters
680 (Pisupati et al., 2021), may improve the perfect integrator model's predictivity.

681 ### 4.3 Implicit confidence markers

682 Although research suggests faster decisions accompanied by higher confidence (Kiani et al., 2014;
683 Vafaei Shooshtari et al., 2019; van den Berg et al., 2016; Zylberberg et al., 2016), our results do not
684 show such an association in the presence of a brief piece of evidence. Moreover, our participants decide
685 much faster in double-pulse trials comparing to single-pulse trials. We hypothesized that the decrease
686 of response-time in double-pulse trials would be reflected with higher internal confidence. However,
687 another hypothesis of this variation pointed to the extra time duration in double-pulse trials, which can
688 be used to increase readiness to decide. We regress the delay-time before response cue onset and
689 response-time in both single-pulse and double-pulse trials to explore the hypothesis. If the variation of
690 response-time was primarily dependent on extra delay time, the delay time should have had a
691 considerable effect on response-time, especially in our 120ms single-pulse trials when the stimulus
692 duration was concise and the delay time varied. Nevertheless, the effects in both double-pulse and
693 single-pulse trials are weak. Accordingly, the hypothesis that faster decision reflect higher confidence
694 in double-pulse trials is supported. In addition, the confidence profile as a function of response-time
695 was significant in double-pulse trials unlike in single-pulse trials.

696 Our findings furthermore suggest that reported confidence might not follow confidence marker in EEG
697 response. We focused on the CPP —a neural correlate of perceptual processing believed to reflect
698 evidence accumulation and correlated to confidence (Boldt et al., 2019; Herding et al., 2019; Rausch
699 et al., 2020; Vafaei Shooshtari et al., 2019; Zizlsperger et al., 2014). However, our findings suggest
700 that in the presence of a brief and weak stimulus, entirely unlike in double-pulse trials, CPP amplitudes
701 show no significant variation in high and low level of confidence. As confidence in single and double-
702 pulse trials did not vary significantly, we suggest that variation of CCP amplitude share more
703 commonalities with implicit confidence measure rather than explicit confidence measures like ratings.
704 Moreover, we propose that pupil response relation to confidence rating varies as the task condition
705 changes; when participants access brief and weak stimuli, no association detected, unlike in the
706 presence of a pair of separated stimuli. Our current observations are not easily reconciled with existing
707 theoretical accounts of the impact of the confidence level on pupil response (Allen et al., 2016; Lempert
708 et al., 2015; Urai et al., 2017).

709 To sum up, when participants access brief and mainly weak stimuli, the confidence ratings are not
710 reliable and confidence profile could not be tracked from response-time, pupil and EEG response. In
711 other words, implicit confidence markers, in some case, might be incapable of following the conscious
712 confidence rating. This is in line with innovative findings abstracting implicit confidence measures
713 from explicit confidence measures (Logan & Crump, 2010).

714 **4.4 Limitations and future directions**

715 To the best of our knowledge, how evidence accumulation processes improve the accuracy confidence
716 association was not addressed using the combination of behavioral, neural, and pupillometry signatures
717 before. Obviously, our results were grounded in assumptions of integration strategy in decision-
718 making. However, this insight has recently been reconsidered (Carland, Marcos, Thura, & Cisek, 2016;
719 Stine et al., 2020). Participants' decisions might be better explained by an urgency-gate model (Evans,
720 Hawkins, Boehm, Wagenmakers, & Brown, 2017; Thura, Beauregard-Racine, Fradet, & Cisek, 2012)
721 rather than an integration strategy such as perfect integrator. A participant's strategy could be
722 something between no integration and perfect integration or in a completely different space of models
723 (Stine et al., 2020) and might be change depending on task paradigm or even subject's internal state
724 (Evans & Hawkins, 2019; Najafi & Churchland, 2018; Tsetsos, Gao, McClelland, & Usher, 2012).
725 Consequently, further models to discuss the decision strategy in the presence of separated pulses could
726 guide future works. In addition, future experiments could develop computational approaches and

This is a provisional file, not the final typeset article

727 attempt to implement other scenarios in a discrete environment to study choice and confidence
728 formation and examine the involved processes.

729  In addition, although the vast number of trials for each participant allowed us to do a robust subject-
730 wise analysis and our EEG study replicated the same behavioral and modeling data, the small number
731 of participants we used prevents us from making general claims. Future research might capitalize on
732 our paradigm to provide a situation in which confidence remains persistent but metacognitive
733 sensitivity improved. In this way, future research continues studying the neural basis of metacognitive
734 ability and consciousness in addition to previous works (Feuerriegel, Blom, & Hogendoorn, 2021;
735 Fleming & Dolan, 2012).

## 5   Conclusion

737 To sum, the present study sheds new light on confidence formation, especially in perceptual decision-
738 making when a pair of visual cues separated by diverse temporal gaps. Our data suggest that
739 accumulated evidence from both pulses shapes confidence but not in line with accuracy. Moreover, we
740 showed that the classic perfect integrator model merely highlighted evidence accumulation which
741 predict the choice and ignored the effect the metacognitive noise that affects confidence. Finally,
742 integrating evidence from two separated pieces of information makes the confidence profiles in RT,
743 EEG and pupil responses show up, unlike the situation in which participants have to decide based on
744 a brief and weak pulse of information.

## 6   Conflict of Interest

746 The authors declare no conflict of interest.

## 7   Author Contributions

748 ZA: conceptualization, data acquisition, analysis, visualization, writing - original draft, writing - review
749 and editing; SZ: conceptualization, supervision, writing - review and editing; AJ: supervision, writing
750 - review and editing; RE: conceptualization, supervision, writing - review and editing.

## 8   Acknowledgments

## 9   Ethics

758 The ethics committee of the Iran University of Medical Sciences (protocol #IR.IUMS.REC1399648)
759 approved the experimental protocol, and subjects gave written informed consent.

## 10   Supplementary Material

761 The Supplementary Material for this article can be found online at:

## 11   Data Availability Statement

The datasets generated for this study are available on request to the corresponding author.

## 12   References

Adelson, E. H., & Bergen, J. R. (1985). Spatiotemporal energy models for the perception of motion. *Josa A*, *2*(2), 284–299.

Allen, M., Frank, D., Schwarzkopf, D. S., Fardo, F., Winston, J. S., Hauser, T. U., & Rees, G. (2016). Unexpected arousal modulates the influence of sensory noise on confidence. *Elife*, *5*, e18103.

Atiya, N. A. A., Rañó, I., Prasad, G., & Wong-Lin, K. F. (2019). A neural circuit model of decision uncertainty and change-of-mind. *Nature Communications*, *10*(1). https://doi.org/10.1038/s41467-019-10316-8

Atiya, N. A. A., Zgonnikov, A., O'Hora, D., Schoemann, M., Scherbaum, S., & Wong-Lin, K. (2020). Changes-of-mind in the absence of new post-decision evidence. *PLoS Computational Biology*, *16*(2), e1007149.

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

Baranski, J. V, Petrusic, W. M., Peters, M. A. K., Thesen, T., Ko, Y. D., Maniscalco, B., … Dolan, R. J. (2017). The relationship between perceptual decision variables and confidence in the human brain. *Journal of Neuroscience*, *32*(18), 412–428.

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

Brainard, D. H., & Vision, S. (1997). The psychophysics toolbox. *Spatial Vision*, *10*, 433–436.

Carland, M. A., Marcos, E., Thura, D., & Cisek, P. (2016). Evidence against perfect integration of sensory information during perceptual decision making. *Physiology*, *115*(1), 1–13.

Charles, L., Van Opstal, F., Marti, S., & Dehaene, S. (2013). Distinct brain mechanisms for conscious versus subliminal error detection. *Neuroimage*, *73*, 80–94.

Cohen, J. (1970). Approximate power and sample size determination for common one-sample and two-sample hypothesis tests. *Educational and Psychological Measurement*, *30*(4), 811–831.

Colizoli, O., De Gee, J. W., Urai, A. E., & Donner, T. H. (2018). Task-evoked pupil responses reflect internal belief states. *Scientific Reports*, *8*(1), 1–13.

De Gardelle, V., & Mamassian, P. (2015). Weighting mean and variability during confidence judgments. *PLoS ONE*, *10*(3). https://doi.org/10.1371/journal.pone.0120870

De Martino, B., Fleming, S. M., Garrett, N., & Dolan, R. J. (2013). Confidence in value-based choice. *Nature Neuroscience*, *16*(1), 105.

This is a provisional file, not the final typeset article

795  Delorme, A., & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG
796      dynamics including independent component analysis. *Journal of Neuroscience Methods*, *134*(1),
797      9–21.

798  Evans, N. J., & Hawkins, G. E. (2019). When humans behave like monkeys: Feedback delays and
799      extensive practice increase the efficiency of speeded decisions. *Cognition*, *184*, 11–18.

800  Evans, N. J., Hawkins, G. E., Boehm, U., Wagenmakers, E.-J., & Brown, S. D. (2017). The
801      computations that support simple decision-making: A comparison between the diffusion and
802      urgency-gating models. *Scientific Reports*, *7*(1), 1–13.

803  Feuerriegel, D., Blom, T., & Hogendoorn, H. (2021). Predictive activation of sensory representations
804      as a source of evidence in perceptual decision-making. *Cortex*, *136*, 140–146.

805  Fleming, S. M. (2017). HMeta-d: hierarchical Bayesian estimation of metacognitive efficiency from
806      confidence ratings. *Neuroscience of Consciousness*, *2017*(1), nix007.

807  Fleming, S. M., & Dolan, R. J. (2012). The neural basis of metacognitive ability. *Philosophical
808      Transactions of the Royal Society B: Biological Sciences*, *367*(1594), 1338–1349.

809  Fleming, S. M., & Lau, H. C. (2014). How to measure metacognition. *Frontiers in Human
810      Neuroscience*, *8*, 443.

811  Fleming, S. M., Putten, E. J., & Daw, N. D. (2018). Neural mediators of changes of mind about
812      perceptual decisions. *Nature Neuroscience*, 1.

813  Gherman, S., & Philiastides, M. G. (2015). Neural representations of confidence emerge from the
814      process of decision formation during perceptual choices. *Neuroimage*, *106*, 134–143.

815  Gold, J. I., & Ding, L. (2013). How mechanisms of perceptual decision-making affect the psychometric
816      function. *Progress in Neurobiology*, *103*, 98–114.

817  Gold, J. I., & Shadlen, M. N. (2007). The neural basis of decision making. *Annual Review of
818      Neuroscience*, *30*.

819  Green, D. M., & Swets, J. A. (1966). *Signal detection theory and psychophysics* (Vol. 1). Wiley New
820      York.

821  Hebart, M. N., Schriever, Y., Donner, T. H., & Haynes, J.-D. (2014). The relationship between
822      perceptual decision variables and confidence in the human brain. *Cerebral Cortex*, *26*(1), 118–
823      130.

824  Herce Castañón, S., Moran, R., Ding, J., Egner, T., Bang, D., & Summerfield, C. (2019). Human noise
825      blindness drives suboptimal cognitive inference. *Nature Communications*, *10*(1), 1–11.
826      https://doi.org/10.1038/s41467-019-09330-7

827  Herding, J., Ludwig, S., von Lautz, A., Spitzer, B., & Blankenburg, F. (2019). Centro-parietal EEG
828      potentials index subjective evidence and confidence during perceptual decision making.
829      *NeuroImage*, *201*, 116011.

830  Hoeks, B., & Ellenbroek, B. A. (1993). A neural basis for a quantitative pupillary model. *Journal of*

831    *Psychophysiology*, *7*, 315.

832    Jang, Y., Wallsten, T. S., & Huber, D. E. (2012). A stochastic detection and retrieval model for the
833        study of metacognition. *Psychological Review*, *119*(1), 186.

834    Kelly, S. P., & O'Connell, R. G. (2013). Internal and external influences on the rate of sensory evidence
835        accumulation in the human brain. *Journal of Neuroscience*, *33*(50), 19434–19441.

836    Kepecs, A., & Mainen, Z. F. (2012). A computational framework for the study of confidence in humans
837        and animals. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*,
838        *367*(1594), 1322–1337.

839    Kepecs, A., Uchida, N., Zariwala, H. A., & Mainen, Z. F. (2008). Neural correlates, computation and
840        behavioural impact of decision confidence. *Nature*, *455*(7210), 227.

841    Kiani, R., Churchland, A. K., & Shadlen, M. N. (2013). Integration of direction cues is invariant to the
842        temporal gap between them. *Journal of Neuroscience*, *33*(42), 16483–16489.

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

845    Kiani, R., Hanks, T. D., & Shadlen, M. N. (2008). Bounded integration in parietal cortex underlies
846        decisions even when viewing duration is dictated by the environment. *Journal of Neuroscience*,
847        *28*(12), 3017–3029.

848    Kiani, R., & Shadlen, M. N. (2009). Representation of confidence associated with a decision by
849        neurons    in    the    parietal    cortex.    *Science*,    *324*(5928),    759–764.
850        https://doi.org/10.1126/science.1169405

851    Kira, S., Yang, T., & Shadlen, M. N. (2015). A neural implementation of Wald's sequential probability
852        ratio test. *Neuron*, *85*(4), 861–873.

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

854    Laeng, B., Sirois, S., & Gredebäck, G. (2012). Pupillometry: A window to the preconscious?
855        *Perspectives on Psychological Science*, *7*(1), 18–27.

856    Lempert, K. M., Chen, Y. L., & Fleming, S. M. (2015). Relating pupil dilation and metacognitive
857        confidence during auditory decision-making. *PLoS One*, *10*(5), e0126588.

858    Logan, G. D., & Crump, M. J. C. (2010). Cognitive illusions of authorship reveal hierarchical error
859        detection in skilled typists. *Science*, *330*(6004), 683–686.

860    Maniscalco, B., & Lau, H. (2012). A signal detection theoretic approach for estimating metacognitive
861        sensitivity from confidence ratings. *Consciousness and Cognition*, *21*(1), 422–430.

862    Maniscalco, B., & Lau, H. (2014). Signal detection theory analysis of type 1 and type 2 data: meta-d′,
863        response-specific meta-d′, and the unequal variance SDT model. In *The cognitive neuroscience of*
864        *metacognition* (pp. 25–66). Springer.

865    Maniscalco, B., & Lau, H. (2016). The signal processing architecture underlying subjective reports of

This is a provisional file, not the final typeset article

866    sensory awareness. *Neuroscience of Consciousness*, *2016*(1).

867    Mathôt, S. (2013). A simple way to reconstruct pupil size during eye blinks. *Retrieved From*, *10*, m9.

868    Meyniel, F., Sigman, M., & Mainen, Z. F. (2015). Confidence as Bayesian probability: From neural
869        origins to behavior. *Neuron*, *88*(1), 78–92.

870    Mognon, A., Jovicich, J., Bruzzone, L., & Buiatti, M. (2011). ADJUST: An automatic EEG artifact
871        detector based on the joint use of spatial and temporal features. *Psychophysiology*, *48*(2), 229–
872        240.

873    Mueller, S. T., & Weidemann, C. T. (2008). Decision noise: An explanation for observed violations of
874        signal detection theory. *Psychonomic Bulletin & Review*, *15*(3), 465–494.

875    Murphy, P. R., Boonstra, E., & Nieuwenhuis, S. (2016). Global gain modulation generates time-
876        dependent urgency during perceptual choice in humans. *Nature Communications*, *7*, 13526.

877    Murphy, P. R., Robertson, I. H., Harty, S., & O'Connell, R. G. (2015). Neural evidence accumulation
878        persists after choice to inform metacognitive judgments. *Elife*, *4*, e11946.

879    Najafi, F., & Churchland, A. K. (2018). Perceptual Decision-Making: A Field in the Midst of a
880        Transformation. *Neuron*, *100*(2), 453–462.

881    Nienborg, H., & Cumming, B. G. (2009). Decision-related activity in sensory neurons reflects more
882        than a neuron's causal effect. *Nature*, *459*(7243), 89–92.

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

885    Paz, L., Insabato, A., Zylberberg, A., Deco, G., & Sigman, M. (2016). Confidence through consensus:
886        a neural mechanism for uncertainty monitoring. *Scientific Reports*, *6*, 21830.

887    Pisupati, S., Chartarifsky-Lynn, L., Khanal, A., & Churchland, A. K. (2021). Lapses in perceptual
888        decisions reflect exploration. *Elife*, *10*, e55490.

889    Rahnev, D., & Denison, R. N. (2018). Suboptimality in perceptual decision making. *Behavioral and
890        Brain Sciences*, *41*.

891    Rahnev, D., Nee, D. E., Riddle, J., Larson, A. S., & D'Esposito, M. (2016). Causal evidence for frontal
892        cortex organization for perceptual decision making. *Proceedings of the National Academy of
893        Sciences*, *113*(21), 6059–6064.

894    Rausch, M., & Zehetleitner, M. (2016). Visibility is not equivalent to confidence in a low contrast
895        orientation discrimination task. *Frontiers in Psychology*, *7*, 591.

896    Rausch, M., Zehetleitner, M., Steinhauser, M., & Maier, M. E. (2020). Cognitive modelling reveals
897        distinct electrophysiological markers of decision confidence and error monitoring. *NeuroImage*,
898        *218*, 116963.

899    Resulaj, A., Kiani, R., Wolpert, D. M., & Shadlen, M. N. (2009). Changes of mind in decision-making.
900        *Nature*, *461*(7261), 263.

901   Roitman, J. D., & Shadlen, M. N. (2002). Response of neurons in the lateral intraparietal area during a
902        combined visual discrimination reaction time task. *Journal of Neuroscience*, *22*(21), 9475–9489.

903   Rollwage, M., Loosen, A., Hauser, T. U., Moran, R., Dolan, R. J., & Fleming, S. M. (2020). Confidence
904        drives a neural confirmation bias. *Nature Communications*, *11*(1), 1–11.

905   Shadlen, M. N., & Kiani, R. (2013). Decision making as a window on cognition. *Neuron*, *80*(3), 791–
906        806.

907   Shekhar, M., & Rahnev, D. (2018). Distinguishing the roles of dorsolateral and anterior PFC in visual
908        metacognition. *Journal of Neuroscience*, *38*(22), 5078–5087.

909   Shekhar, M., & Rahnev, D. (2020). Sources of Metacognitive Inefficiency. *Trends in Cognitive*
910        *Sciences*.

911   Shekhar, M., & Rahnev, D. (2021). The nature of metacognitive inefficiency in perceptual decision
912        making. *Psychological Review*, *128*(1), 45.

913   Steinemann, N. A., O'Connell, R. G., & Kelly, S. P. (2018). Decisions are expedited through multiple
914        neural adjustments spanning the sensorimotor hierarchy. *Nature Communications*, *9*(1), 1–13.

915   Stine, G. M., Zylberberg, A., Ditterich, J., & Shadlen, M. N. (2020). Differentiating between
916        integration and non-integration strategies in perceptual decision making. *Elife*, *9*, e55365.

917   Tagliabue, C. F., Veniero, D., Benwell, C. S. Y., Cecere, R., Savazzi, S., & Thut, G. (2019). The EEG
918        signature of sensory evidence accumulation during decision formation closely tracks subjective
919        perceptual experience. *Scientific Reports*, *9*(1), 1–12.

920   Thura, D., Beauregard-Racine, J., Fradet, C.-W., & Cisek, P. (2012). Decision making by urgency
921        gating: theory and experimental support. *Journal of Neurophysiology*, *108*(11), 2912–2930.

922   tickle, hannah, Tsetsos, K., Speekenbrink, M., & Summerfield, C. (2020). Optional Stopping in a
923        Heteroscedastic World. *PsyArXiv*. https://doi.org/10.31234/OSF.IO/T7DN2

924   Tohidi-Moghaddam, M., Zabbah, S., Olianezhad, F., & Ebrahimpour, R. (2019). Sequence-dependent
925        sensitivity explains the accuracy of decisions when cues are separated with a gap. *Attention,*
926        *Perception, and Psychophysics*, *81*(8), 2745–2754. https://doi.org/10.3758/s13414-019-01810-8

927   Tsetsos, K., Gao, J., McClelland, J. L., & Usher, M. (2012). Using time-varying evidence to test models
928        of decision dynamics: bounded diffusion vs. the leaky competing accumulator model. *Frontiers*
929        *in Neuroscience*, *6*, 79.

930   Twomey, D. M., Kelly, S. P., & O'Connell, R. G. (2016). Abstract and effector-selective decision
931        signals exhibit qualitatively distinct dynamics before delayed perceptual reports. *Journal of*
932        *Neuroscience*, *36*(28), 7346–7352.

933   Twomey, D. M., Murphy, P. R., Kelly, S. P., & O'connell, R. G. (2015). The classic P300 encodes a
934        build-to-threshold decision variable. *European Journal of Neuroscience*, *42*(1), 1636–1643.

935   Urai, A. E., Braun, A., & Donner, T. H. (2017). Pupil-linked arousal is driven by decision uncertainty
936        and alters serial choice bias. *Nature Communications*, *8*(1), 1–11.

This is a provisional file, not the final typeset article

937    Vafaei Shooshtari, S., Esmaily Sadrabadi, J., Azizi, Z., & Ebrahimpour, R. (2019). Confidence
938        Representation of Perceptual Decision by EEG and Eye Data in a Random Dot Motion Task.
939        *Neuroscience*, *406*. https://doi.org/10.1016/j.neuroscience.2019.03.031

940    Van Den Berg, R., Anandalingam, K., Zylberberg, A., Kiani, R., Shadlen, M. N., & Wolpert, D. M.
941        (2016). A common mechanism underlies changes of mind about decisions and confidence. *ELife*,
942        *5*(FEBRUARY2016), 1–21. https://doi.org/10.7554/eLife.12192

943    Van den Berg, R., Yoo, A. H., & Ma, W. J. (2017). Fechner's law in metacognition: A quantitative
944        model of visual working memory confidence. *Psychological Review*, *124*(2), 197.

945    van den Berg, R., Zylberberg, A., Kiani, R., Shadlen, M. N., & Wolpert, D. M. (2016). Confidence Is
946        the Bridge between Multi-stage Decisions. *Current Biology*, *26*(23), 3157–3168.
947        https://doi.org/10.1016/j.cub.2016.10.021

948    Wang, X.-J. (2002). Probabilistic decision making by slow reverberation in cortical circuits. *Neuron*,
949        *36*(5), 955–968.

950    Waskom, M. L., & Kiani, R. (2018). Decision Making through Integration of Sensory Evidence at
951        Prolonged Timescales. *Current Biology*, *28*(23), 3850–3856.e9.
952        https://doi.org/10.1016/j.cub.2018.10.021

953    Wong, K.-F., & Huk, A. C. (2008). Temporal dynamics underlying perceptual decision making:
954        Insights from the interplay between an attractor model and parietal neurophysiology. *Frontiers in*
955        *Neuroscience*, *2*, 28.

956    Zizlsperger, L., Sauvigny, T., Händel, B., & Haarmeier, T. (2014). Cortical representations of
957        confidence in a visual perceptual decision. *Nature Communications*, *5*, 3940.

958    Zylberberg, A., Barttfeld, P., & Sigman, M. (2012). The construction of confidence in a perceptual
959        decision. *Frontiers in Integrative Neuroscience*, *6*, 79.

960    Zylberberg, A., Fetsch, C. R., & Shadlen, M. N. (2016). The influence of evidence volatility on choice,
961        reaction time and confidence in a perceptual decision. *ELife*, *5*(OCTOBER2016), e17688.
962        https://doi.org/10.7554/eLife.17688

963    Zylberberg, A., Roelfsema, P. R., & Sigman, M. (2014). Variance misperception explains illusions of
964        confidence in simple perceptual decisions. *Consciousness and Cognition*, *27*(1), 246–253.
965        https://doi.org/10.1016/j.concog.2014.05.012

966    *Zylberberg, A., Wolpert, D. M., & Shadlen, M. N. (2018). Counterfactual reasoning underlies the*
967        *learning of priors in decision making.* Neuron*, 99(5), 1083–1097.*

968      Supplementary Material

969   **13      Supplementary Appendix 1: Figures and Tables**

970   **13.1  Supplementary Tables**

971   **Supplementary Table 1.** Subtraction of confidence in double-pulse from single-pulse trials was
972   significantly affected by choice accuracy.

| Participant | $\beta_1$ |
|---|---|
| $P_1$ | 0.17 ($p < 0.01$)<br>CI = [.15, .19] |
| $P_2$ | $0.21 \pm 0.01$ ($p < 0.01$)<br>CI = [.19, .23] |
| $P_3$ | $0.10 \pm 0.01$ ($p < 0.01$)<br>CI = [.08, .12] |
| $P_4$ | $0.12 \pm 0.01$ ($p < 0.01$)<br>CI = [.10, .14] |

973   Each row shows the coefficients of Eq.10 of manuscript, their related $p$ values and a 95% confidence
974   interval.

This is a provisional file, not the final typeset article

**Supplementary Table 2.** Performance was largely unaffected by interpulse interval for double-pulse trials with equal pulse strength and with unequal pulse strength.

| Participant | Equal strength | | Unequal strength | | |
|---|---|---|---|---|---|
| | $\beta_3$ | $\beta_4$ | $\beta_3$ | $\beta_4$ | $\beta_5$ |
| **P$_1$** | $-0.45 \pm 1.07$ ($p = 0.67$) | $-0.01 \pm 0.01$ ($p = 0.91$) | $0.39 \pm 0.84$ ($p = 0.63$) | $-0.10 \pm 0.09$ ($p = 0.29$) | $-0.01 \pm 0.01$ ($p = 0.77$) |
| | CI = [-2.55, 1.65] | 95% CI = [-.03, .01] | CI = [-1.26, 2.03] | CI = [-.28, .08] | CI = [-.03, .01] |
| **P$_2$** | $-2.58 \pm 1.50$ ($p = 0.09$) | $-0.01 \pm 0.08$ ($p = 0.93$) | $0.78 \pm 1.23$ ($p = 0.52$) | $0.01 \pm 0.15$ ($p = 0.95$) | $0.02 \pm 0.07$ ($p = 0.77$) |
| | CI = [-5.52, .36] | CI = [-.17, .15] | CI = [-1.63, 3.19] | CI = [-.28, .30] | CI = [-.11, .17] |
| **P$_3$** | $0.48 \pm 1.08$ ($p = 0.66$) | $0.07 \pm 0.11$ ($p = 0.52$) | $2.18 \pm 0.95$ ($p = 0.03$) | $-0.19 \pm 0.10$ ($p = 0.06$) | $0.14 \pm 0.10$ ($p = 0.17$) |
| | CI = [-1.64, 2.60] | CI = [-.15, .28] | CI = [.32, 4.04] | CI = [-.39, .01] | CI = [-.06, .34] |
| **P$_4$** | $0.74 \pm 0.97$ ($p = 0.44$) | $-0.12 \pm 0.08$ ($p = 0.15$) | $-0.57 \pm 0.75$ ($p = 0.45$) | $0.02 \pm 0.08$ ($p = 0.81$) | $-0.07 \pm 0.07$ ($p = 0.33$) |
| | CI = [-1.16, 2.64] | CI = [-.28, .04] | CI = [-2.04, .90] | CI = [-.14, .18] | CI = [-.21, .07] |

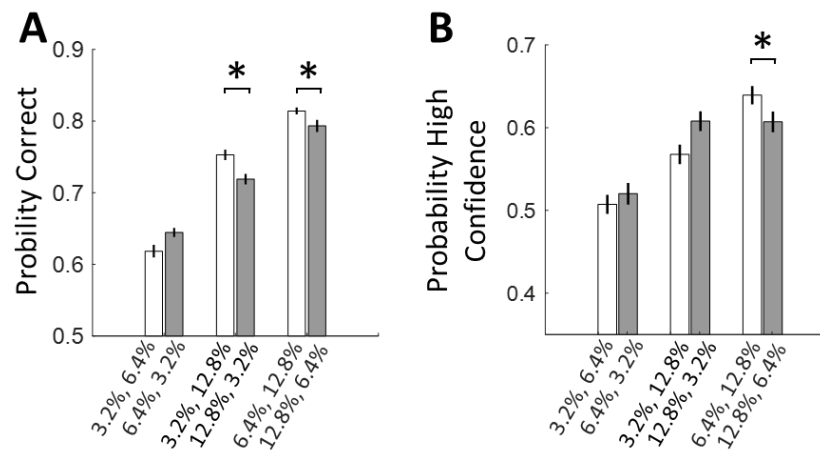Each row shows the coefficients of Eq.4 and 5, their related $p$ values and a 95% confidence interval of $\beta_i$.

35

980  **Supplementary Table 3.** Pairwise comparisons across models (1: single-pulse trials, 2: double-pulse
981  trials, 3: perfect integrator) for SDT parameters.

| | $d'$ | | | $Meta\text{-}d'$ | | | $Meta\text{-}d'/d'$ | | |
|---|---|---|---|---|---|---|---|---|---|
| | 1 vs 2 | 1 vs 3 | 2 vs 3 | 1 vs 2 | 1 vs 3 | 2 vs 3 | 1 vs 2 | 1 vs 3 | 2 vs 3 |
| *tstat* | 3.79 | 4.05 | 1.25 | 1.98 | 0.21 | 2.48 | 1.74 | 0.12 | 3.58 |
| *df* | 22 | 22 | 22 | 22 | 22 | 22 | 22 | 22 | 22 |
| *pValue* | 0.001 | $0.52 \times 10^{-5}$ | 0.22 | 0.05 | 0.83 | 0.02 | 0.09 | 0.90 | 0.002 |
| **95% CI** | [.13, .44] | [.20, 0.62] | [-.08, 0.32] | [-.02, 1.12] | [-.66, 0.82] | [.10, 1.15] | [-.17, 1.98] | [-1.09, 1.23] | [.35, 1.32] |
| **Cohen's *d*** | 1.55 | 1.65 | - 0.51 | 0.81 | 0.09 | 1.01 | 0.71 | 0.05 | 1.46 |

982

This is a provisional file, not the final typeset article

983    **13.2  Supplementary Figures**



**Supplementary Figure 10. Choice confidence was not depended on the sequence of motion pulses. (A)** The weak–strong pulse sequence contributed higher accuracy than the strong–weak sequence. **(B)** The weak–strong pulse sequence did not contribute higher confidence comparing to the strong–weak sequence. In all panels, data are represented as group mean ± SEM. (*$p < 0.05$)
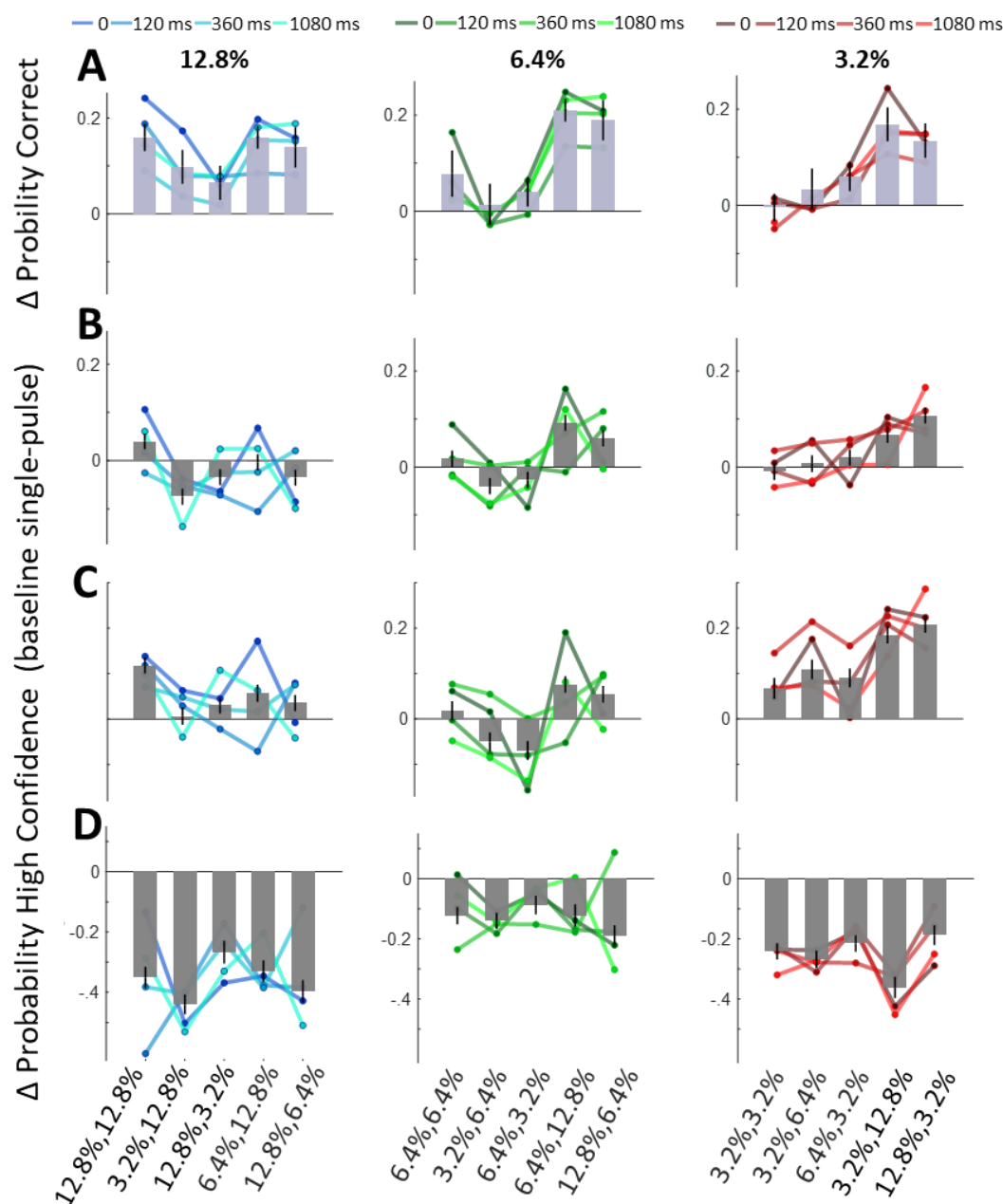
984

985

**Confidence and accuracy after integration of discrete evidence**



**Supplementary Figure 11. Interplay between confidence/accuracy and interpulse interval in double-pulse trials. (A)** Choice accuracy for double-pulse trials grouping in all possible interval conditions. **(B)** Confidence of double-pulse trials was calculated by pooling data across all time intervals. In (A) and (B) each data point addresses pooled data from indicated sequence pulse and its reverse order (e.g., 12.8– 3.2% and 3.2 –12.8%).
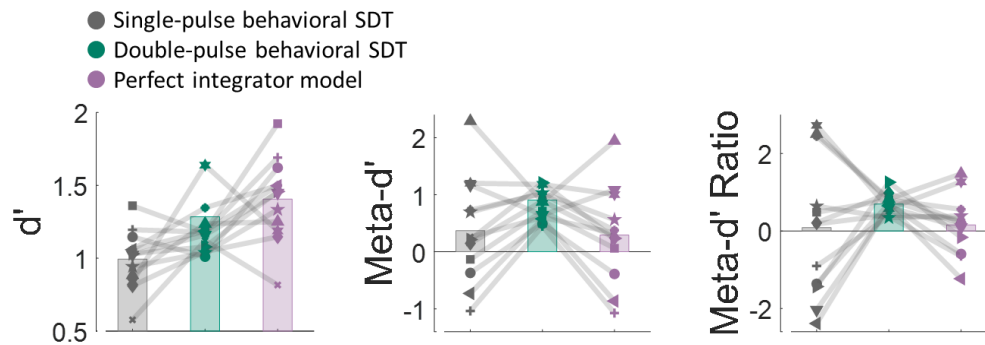
986

987

This is a provisional file, not the final typeset article

**Supplementary Figure 12. Variation of accuracy or confidence in double-pulse trials baselined by corresponding coherence (3.2%, 6.4% and 12.8% for each column). (A)** Considering all the trials, accuracy improved in combination with almost all pulses comparing to the baseline. **(B)** Considering all the trials, confidence improved in combination with stronger pulses while the confidence in sequence with a weaker pulse either decreased or remained constant. **(C)** In correct-choice trials, the increasing effect of stronger pulses is more significant and the confidence even slightly improved in combination with weaker pulses comparing to corresponding baseline. **(D)** Interestingly, in incorrect trials, the confidence decreased in every condition. The colored line representing matching data for each of four possible gaps. In bar graph, the data are represented as group mean ± SEM.
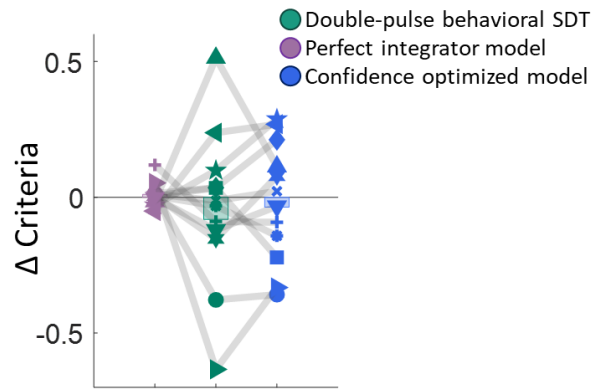
988

**Supplementary Figure 4. Comparison of models and human behavior.** Stimulus sensitivity ($d'$), metacognitive sensitivity ($Meta- d'$) and, metacognitive efficiency ($Meta- d'/d'$) estimated for single-pulse trials, double-pulse trials and the perfect integrator models.

989

990    A univariate ANOVA showed that $d'$ between models fit to double/single-pulse trials and the perfect
991    integrator model significantly differed (F(2,33) = 9.99; $p = 0.41 \times 10^{-4}$). Also, a univariate ANOVA
992    showed that $Meta- d'$ between models fit to double/single-pulse trials and the perfect integrator model
993    partialy differed (F(2,33) = 1.04; $p$ = 0.09). We also computed metacognitive efficiency
994    ($Meta- d'/d'$), A univariate ANOVA revealed a significant difference on all three models (F(2,33) =
995    2.50; $p$ = 0.10).), We also applied the *t*-test as a post hoc procedure to compare all pairs of $d'$, $Meta- d'$,
996    $Meta- d'/d'$ from three models (Supplementary **Supplementary** **Table 3.** Pairwise comparisons
997    across models (1: single-pulse trials, 2: double-pulse trials, 3: perfect integrator) for SDT
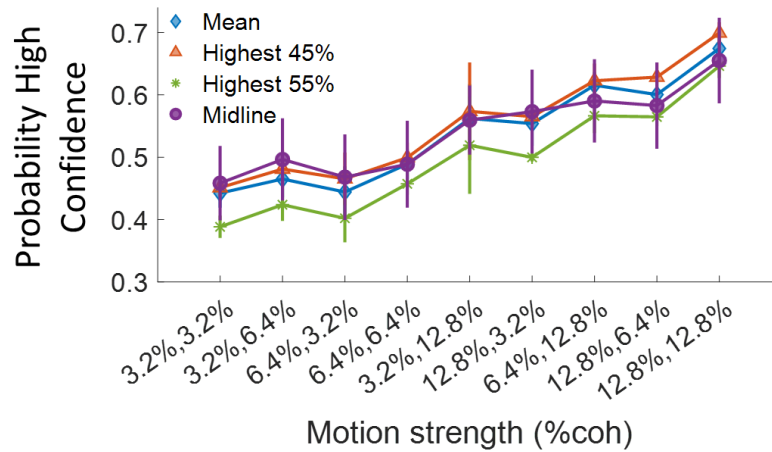998    parameters.Table 3).

999

**Confidence and accuracy after integration of discrete evidence**



**Supplementary Figure 5.** Variation of confidence criteria comparing to single-pulse trials in perfect integrator vs double-pulse trials and optimized model.
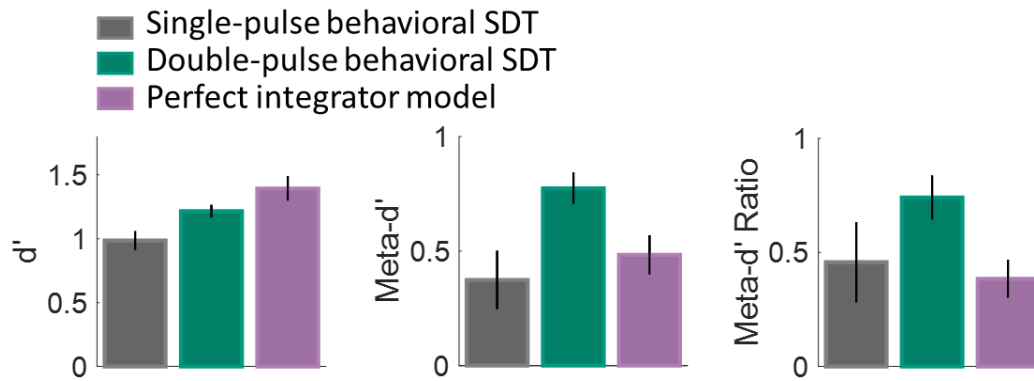
1000

1001

**Supplementary Figure 6.** A univariant Anova showed that confidence categorized by four different approaches in double-pulse trials not significantly differed (F(3,140) = 9.99; $p$ = 0.34). Three paired-samples $t$-tests between our confidence categorization with other methods showed no difference (all $ps$ > 0.36).

1002

This is a provisional file, not the final typeset article

**Supplementary Figure 7. Comparison of models and human behavior considering the same numbers of trials.** (a) Stimulus sensitivity ($d'$), metacognitive sensitivity ($Meta\text{-}d'$), metacognitive efficiency ($Meta\text{-}d'/d'$) estimated for single-pulse trials, double-pulse trials and the perfect integrator models.

1003

1004 We compared $d'$, $Meta\text{-}d'$ and, $Meta\text{-}d'/d'$ of fitted models to single/double-pulse trials and
1005 simulated data by perfect integrator model, following up with three Dunn pair tests. A Kruskal-Wallis
1006 test showed that $d'$ between models fit to double/single-pulse trials and the perfect integrator model
1007 not significantly differed ($H(3) = 3.23$; $p = 0.20$). We also applied the Dunn test as a post hoc procedure
1008 to compare all pairs of $d'$ from three models. No $d'$ in models significantly differed from others (all *ps*
1009 > 0.21).

1010 Also, a Kruskal-Wallis test showed that $Meta\text{-}d'$ between models fit to double/single-pulse trials and
1011 the perfect integrator model significantly differed ($H(3) = 6.96$; $p = 0.03$). Post-hoc Dunn were used to
1012 compare all pairs of $Meta\text{-}d'$ from three models. The difference of $Meta\text{-}d'$ of single-pulse trials
1013 and double-pulse was significant ($p = 0.03$, $CI = [-12.58, -0.41]$). However, the difference of $Meta\text{-}d'$
1014 was insignificant for single-pulse trials and perfect integrator model ($p = 0.87$, $CI = [-7.83, 4.33]$) and
1015 for double-pulse trials and perfect integrator model ($p = 0.17$, $CI = [4.75, 10.83]$).

1016 We also computed metacognitive efficiency ($Meta\text{-}d'/d'$), A Kruskal-Wallis test revealed a
1017 significant difference on all three models ($H(3) = 7.42$, $p = 0.02$), metacognitive efficiency in double-
1018 pulse and perfect integrator differed significantly ($p = 0.04$, $CI = [0.16\ 12.33]$) while in double-pulse
1019 and single-pulse, it partially differed ($p = 0.07$, $CI = [-11.83, 0.33]$). The difference of single-pulse and
1020 perfect integrator was not significant ($p = 0.99$, $CI = [-5.58\ 6.58]$).

1021

43

## 14  Supplementary Appendix 2: Signal detection theory models

In the binary decision, the observer must discriminate between stimuli labeled $S_2$ or labeled $S_1$. Each stimulus presentation generates a value on an internal decision axis (Figure 1b), corresponding to the evidence in favor of $S_1$ or $S_2$. Evidence generated by each stimulus class is normally distributed across the decision axis, and the distance between these distributions in standard deviation units ($d'$) measures how well the observer can discriminate $S_1$ from $S_2$. The observer sets a decision criterion $cr$, such that all signals exceeding $cr$ are labeled $S_2$ and all those failing to exceed $cr$ are labeled $S_1$. The observer also sets criteria $cr_{2,"S1"}$ and $cr_{2,"S2"}$ to determine confidence ratings around the decision criterion $cr$. These two thresholds must be well-ordered so that $cr_{2,"S1"} < cr < cr_{2,"S2"}$ (Figure 1b). When a $S_2$ response is made, a confident $S_2$ response requires the evidence also to have surpassed the $cr_{2,"S2"}$ threshold [1]. Consider only trials where the observer responds $S_2$, which means the decision axis exceeding $cr$. Then the $S_2$ distribution corresponds to the distribution of evidence for correct responses (i.e., $S_2$ stimuli classified as $S_2$), and the $S_1$ distribution corresponds to the distribution of evidence for incorrect responses (i.e., $S_1$ stimuli classified as $S_2$).

### 14.1  Confidence *Hit Rate* and *False Alarm Rate*

Sweeping the $cr_{2,"S2"}$ criterion across the decision axis generates different values for confidence false alarm rate ($Prob(conf = "h" \mid stim \neq resp)$) and confidence hit rate ($Prob(conf = "h" \mid stim = resp)$). A summary of the observer's confidence performance is provided by hit rate (Hit Rate2) and false alarm rate (False Alarm Rate2)[1]:

$$\text{Hit Rate2} = \text{Prob}(conf = "h" \mid stim = resp) = \frac{n(\text{high conf correct})}{n(\text{correct})},$$

$$\text{False Alarm Rate2} = \text{Prob}(conf = "h" \mid stim \neq resp) = \frac{n(\text{high conf incorrect})}{n(\text{incorrect})},$$

(1)

where $n(cond)$ denotes a count of the total number of trials satisfying the condition $cond$.

### 14.2  Decision *Hit Rate* and *False Alarm Rate*

In the SDT model, the decision hit rate (Hit Rate1) and the decision false alarm rate (False Alarm Rate1) are also calculated as follows:

$$\text{Hit Rate1} = \frac{n(resp = Si, \ stim = Si)}{n(stim = Si)}, i = 1,2$$

$$\text{False Alarm Rate1} = \frac{n(resp = Si, \ stim = Sj)}{n(stim = Sj)}, i = 1, j = 2 \text{ or } i = 2, j = 1$$

(2)

This is a provisional file, not the final typeset article

1048  where $i$ and $j$ represent the stimulus classification. After calculating the Hit Rate1
1049  and False Alarm Rate1 of each participant, $d'$ and $cr$ are calculated as follows for each participant:

$$d' = \phi^{-1}(\text{Hit Rate1}, 0, 1) - \phi^{-1}(\text{False Alarm Rate1}, 0, 1)$$

$$(3)$$

$$cr = -0.5 * [\phi^{-1}(\text{Hit Rate1}, 0, 1) + \phi^{-1}(\text{False Alarm Rate1}, 0, 1)]$$

1050

1051  here, $\phi^{-1}$ is the inverse of a function that represents a normal cumulative distribution and is calculated
1052  as follows:

$$\phi(s, \mu, \sigma) = \int_{-\infty}^{0} N(v, \mu, \sigma)dv, \qquad (4)$$

1053

1054  where $N(v, \mu, \sigma)$ is a Normal distribution with mean ($\mu$) and standard deviation ($\sigma$). After the above
1055  calculations, to simplify the next calculations, we may consider the value of $cr$ as zero point and move
1056  the distribution diagrams related to each option on the axis of the evidence.

1057  By setting $d'$, $cr$ and two criteria $cr_{2,"S1"}$ and $cr_{2,"S2"}$ (Figure 1B), the probabilities of each confidence
1058  rating conditional on a given stimulus and response (Hit Rate2 and False Alarm Rate2) can be
1059  calculated theoretically according to the following equations:

1060

$$\text{Prob(conf} = "h"|\text{stim} = S1, \text{resp} = "S1") = \text{HitRate2}_{"S1"} = \frac{\phi\left(cr_{2,"S1"}, -\frac{d'}{2}\right)}{\phi\left(cr, -\frac{d'}{2}\right)}$$

$$\text{Prob(conf} = "h"|\text{stim} = S2, \text{resp} = "S1") = \text{FalseAlarmRate2}_{"S1"} = \frac{\phi\left(cr_{2,"S1"}, \frac{d'}{2}\right)}{\phi\left(cr, \frac{d'}{2}\right)} \qquad (5)$$

$$\text{Prob(conf} = "h"|\text{stim} = S1, \text{resp} = "S2") = \text{HitRate2}_{"S2"} = \frac{1 - \phi\left(cr_{2,"S2"}, \frac{d'}{2}\right)}{1 - \phi\left(cr, \frac{d'}{2}\right)}$$

$$\text{Prob}(\text{conf} = "h"|\text{stim} = S2, \text{resp} = "S2") = \text{FalseAlarmRate2}_{"S2"} = \frac{1 - \phi\left(\text{cr}_{2,"S2"}, -\frac{d'}{2}\right)}{1 - \phi\left(\text{cr}, -\frac{d'}{2}\right)}$$

1061

1062 In the SDT model, there are different methods for adjusting the model with the data obtained from the
1063 experiments. In the method we used, $d'$ and $cr$ were calculated from the participants' performance (Eq.
1064 3). Then, using maximum likelihood estimation (MLE) and Eq. 1 and 5 and by altering the value of
1065 the confidence criteria while holding $d'$ and $cr$ constant, a set of (Hit Rate2, False Alarm Rate2) pairs
1066 ranging between $(0, 0)$ and $(1, 1)$ were generated. Moreover, $Meta-d'$ was found by fitting the decision
1067 SDT model to response-specific confidence.

1068 **15 Supplementary references**

1069 [1] B. Maniscalco and H. Lau, "Signal detection theory analysis of type 1 and type 2 data: meta-
1070 d′, response-specific meta-d′, and the unequal variance SDT model," in *The cognitive neuroscience of*
1071 *metacognition*, Springer, 2014, pp. 25–66.

1072

This is a provisional file, not the final typeset article