1

# Dynamic influences on static measures of metacognition

3

**Kobe Desender[1,2,3], Luc Vermeylen[3], & Tom Verguts[3]**

1. Brain and Cognition, KU Leuven, Belgium

2. Department of Neurophysiology and Pathophysiology, University Medical Center Hamburg-Eppendorf, Germany

3. Department of Experimental Psychology, Ghent University, Belgium

9

10

**Corresponding author:**

Dr. Kobe Desender

Brain and Cognition

KU Leuven

Tiensestraat 102, 3000 Leuven

Belgium

E-mail: Kobe.Desender@KULeuven.be

18                                    **Abstract (245 words)**

19            Humans differ in their capability to judge the accuracy of their own choices via confidence

20     judgments. Signal detection theory has been used to quantify the extent to which confidence tracks

21     accuracy via M-ratio, often referred to as metacognitive efficiency. This measure, however, is static in

22     that it does not consider the dynamics of decision making. This could be problematic because humans

23     may shift their level of response caution to alter the tradeoff between speed and accuracy. Such shifts

24     could induce unaccounted-for sources of variation in the assessment of metacognition. Instead, evidence

25     accumulation frameworks consider decision making, including the computation of confidence, as a

26     dynamic process unfolding over time. We draw on evidence accumulation frameworks to examine the

27     influence of response caution on metacognition. Simulation results demonstrate that response caution has

28     an influence on M-ratio. We then tested and confirmed that this was also the case in human participants

29     who were explicitly instructed to either focus on speed or accuracy. We next demonstrated that this

30     association between M-ratio and response caution was also present in an experiment without any

31     reference towards speed. The latter finding was replicated in an independent dataset. In contrast, when

32     data were analyzed with a novel dynamic measure of metacognition, which we refer to as v-ratio, in all of

33     the three studies there was no effect of speed-accuracy tradeoff. These findings have important

34     implications for research on metacognition, such as the question about domain-generality, individual

35     differences in metacognition and its neural correlates.

36

## Introduction

37

38      When asked to explicitly report how sure they are about their decisions, humans often claim high

39 confidence for correct and low confidence for incorrect decisions. This capacity to evaluate the accuracy

40 of decisions is often referred to as metacognitive accuracy. Although metacognitive accuracy about

41 perceptual decisions is generally high [1], it varies significantly between participants [2] and between

42 conditions [3]. Such differences in metacognitive accuracy have important real-life consequences, as they

43 relate, for example, to political extremism [4] and psychiatric symptoms [5].

44      A debated question is how to quantify metacognitive accuracy. One prominent issue why one

45 cannot simply calculate the correlation between confidence and accuracy [6] is that this confounds task

46 accuracy with metacognitive accuracy; i.e. it is much easier to detect one's own mistakes in an easy task

47 than in a hard task. Different solutions have been proposed in the literature, such as using coefficients

48 from a logistic mixed-model [7], type 2 ROC curves [2], and meta-$d'$ [8,9]. Rather than providing an in-depth

49 discussion and comparison of these different measures, we here focus on one of these static approaches,

50 namely the meta-$d'$ framework, the state-of-the-art measure of metacognitive accuracy [10]. The meta-$d'$

51 approach is embedded within signal detection theory, and quantifies the accuracy with which confidence

52 ratings discriminate between correct and incorrect responses (*meta-d'*) while controlling for first-order

53 task performance ($d'$). Because both measures are on the same scale, one can calculate the ratio between

54 both, meta-$d'$/$d'$, also called M-ratio, often referred to as metacognitive *efficiency*. When M-ratio is 1, all

55 available first-order information is used in the (second-order) confidence judgment. When M-ratio is

56 smaller than 1, metacognitive sensitivity is suboptimal, meaning that not all available information from

57 the first-order response is used in the metacognitive judgment (Fleming & Lau, 2014). This measure has

58 been used to address a variety of issues, such as whether metacognition is a domain-general capacity

59 [3,11,12], the neural correlates of metacognition [13–16], how bilinguals differ from monolinguals [17], and how

60 individual differences in metacognitive accuracy correlate with various constructs [4,5].
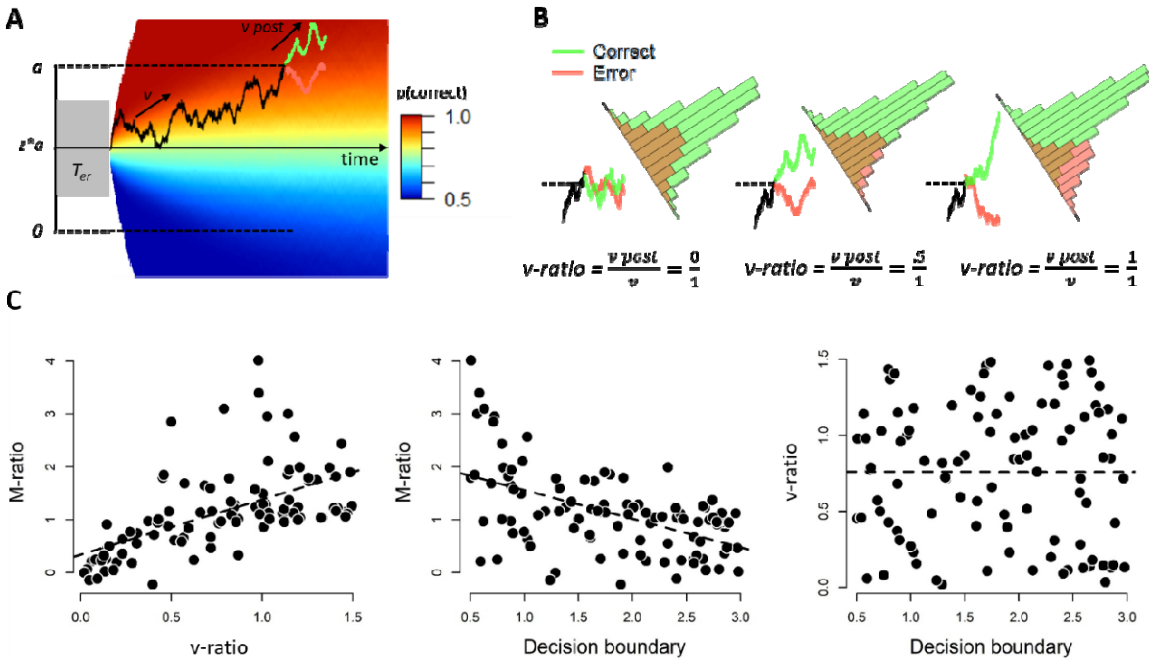
61      An important limitation is that the meta-$d'$ framework (just like the other static approaches cited

62 above), does not consider dynamic aspects of decision making. Put simply, this measure depends on end-

63 of-trial confidence and accuracy, but not on the response process governing the choice and its resulting

64 reaction time. It is well known, however, that choice accuracy depends on response caution; i.e. accuracy

65 decreases when instructing participants to be fast rather than to be correct. The fact that static approaches

66 of metacognition do not consider response caution is problematic because it confounds ability with

67 caution: when focusing on speed rather than accuracy, one will produce many errors due to premature

68 responding, and those errors are much easier to detect compared to errors resulting from low signal

69  quality [18]. Importantly, detecting "premature" errors does not imply "good metacognition" per se, but

70  instead simply depends on one's level of response caution.

71  To account for dynamic influences on metacognition, we propose to instead quantify

72  metacognitive accuracy in a dynamic probabilistic framework [19,20]. Sequential sampling models explain

73  human decision making as a dynamic process of evidence accumulation [21–23]. Specifically, decisions are

74  conceptualized as resulting from the accumulation of noisy evidence towards one of two decision

75  boundaries. The first boundary that is reached, triggers its associated decision. The height of the decision

76  boundary controls the response caution with which a decision is taken [24,25]. When lowering the boundary,

77  decisions will be faster but less accurate; when increasing the boundary, decisions will be slower but

78  more accurate. The prototypical dynamic sampling model is the drift diffusion model (DDM). In this

79  model, confidence can be quantified as the probability of being correct, given evidence, decision time,

80  and the decision that was made [26–28]. The relation between these three variables is represented by the heat

81  map in Figure 1A. It captures the typical finding that trials with strong evidence are more likely to be

82  correct than trials with weak evidence; and that trials with short RTs are more likely to be correct than

83  trials with long RTs. As mentioned, the process of evidence accumulation terminates at the first boundary

84  crossing. Formally, at that time the probability that the choice was correct can be quantified as

85  $p(correct|e_t, t, X)$, where $e_t$ is the level of evidence at time $t$, $t$ is the timing of boundary crossing and $X$ is

86  the choice made [26,28,29]. In typical experiments, however, confidence judgments are provided separately in

87  time (at time $t + s$, i.e., in a separate judgment after the choice), allowing evidence to further accumulate

88  after boundary crossing. As a consequence, confidence should then be quantified as $p(correct|e_{t+s}, t+s, X)$,

89  [19,20,30].

90  Within this formulation, good metacognitive accuracy can be considered as the ability to

91  distinguish corrects versus errors based on $p(correct|e_{t+s}, t+s, X)$. Critically, the difference in the quantity

92  $p(correct|e_{t+s}, t+s, X)$ for corrects versus errors, directly depends on the strength of post-decision

93  accumulation. Thus, we can use post-decision drift rate as a dynamic measure of metacognitive accuracy.

94  For comparison with the M-ratio framework, we quantified v-ratio as the ratio between post-decision drift

95  rate and drift rate. Figure 1B shows post-decision accumulation for three scenarios with varying levels of

96  v-ratio. As can be seen, if v-ratio is zero (left panel), additional evidence meanders adrift for both corrects

97  and errors, and the model does not detect its own errors, i.e., representing a case of poor metacognitive

98  accuracy. If however, v-ratio equals 1 (i.e., post-decision drift rate and drift rate are the same), additional

99  evidence confirms most of the correct choices (i.e., leading to high confidence) and disconfirms most of

100 the error choices (i.e., leading to low confidence), i.e., representing good metacognitive accuracy. We

101    thus propose that v-ratio can be used as a dynamic measure of metacognitive accuracy. In the following,

102    we will shed light on the role of variation in response caution on both M-ratio and v-ratio.



103

104    *Figure 1. Quantifying metacognitive accuracy within an evidence accumulation framework. A.*

105    *Noisy sensory evidence accumulates over time, until the integrated evidence reaches one of two decision*

106    *boundaries (a or 0). After the decision boundary is reached, evidence continues to accumulate. The heat*

107    *map shows the probability of being correct conditional on time, evidence, and the choice made (the*

108    *choice corresponding to the upper boundary, in this example). Confidence is quantified as just this*

109    *probability. B. Histograms of model-predicted confidence for different levels of v-ratio (reflecting the*

110    *ratio between post-decision drift rate and drift rate). Higher levels of v-ratio are associated with better*

111    *dissociating corrects from errors. C. Simulations from this dynamic evidence accumulation model show*

112    *that M-ratio captures variation in v-ratio (r = .58; left panel), and critically, that M-ratio is also related*

113    *to the differences in decision boundary (r = -.52; middle panel). By design, decision boundary and v-ratio*

114    *are unrelated to each other (r ~ 0; right panel).*

115                                              **Results**

116           **Model simulations reveal a link between response caution and M-ratio**

117           We simulated data from a drift diffusion model with additional post-decisional evidence

118    accumulation (see Figure 1A). Decision confidence was quantified as the probability of being correct

119    given evidence, time and choice [26,30,31]. We simulated data for 100 agents with 500 observations each; for

120    each agent, a different random value was selected for drift rate, non-decision time, decision boundary and

121    post-decision drift rate (see Methods). We then used these data to compute M-ratio. As explained before,

122    v-ratio was computed as the ratio between post-decision drift rate and drift rate. The results of our

123    simulation study showed that, first, there was a clear positive relation between M-ratio and v-ratio, $r(98)$

124    $= .58$, $p < .001$, reflecting that M-ratio captures individual variation in metacognition (Figure 1C, left

125    panel). However, we also observed a strongly negative relation between M-ratio and decision boundary,

126    $r(98) = -.52$, $p < .001$ (Figure 1C, central panel). This shows that M-ratio is highly dependent on the

127    speed-accuracy tradeoff that one adopts. This occurs because lowering the decision boundary increases

128    the probability of "fast errors" (i.e. due to noise), which are very likely to generate conflicting evidence in

129    the post-decisional period (i.e. to be detected as an error). Finally, by design there was no relation

130    between v-ratio and decision boundary, $r(98) = .006$, $p = .95$ (Figure 1C, right panel). The full correlation
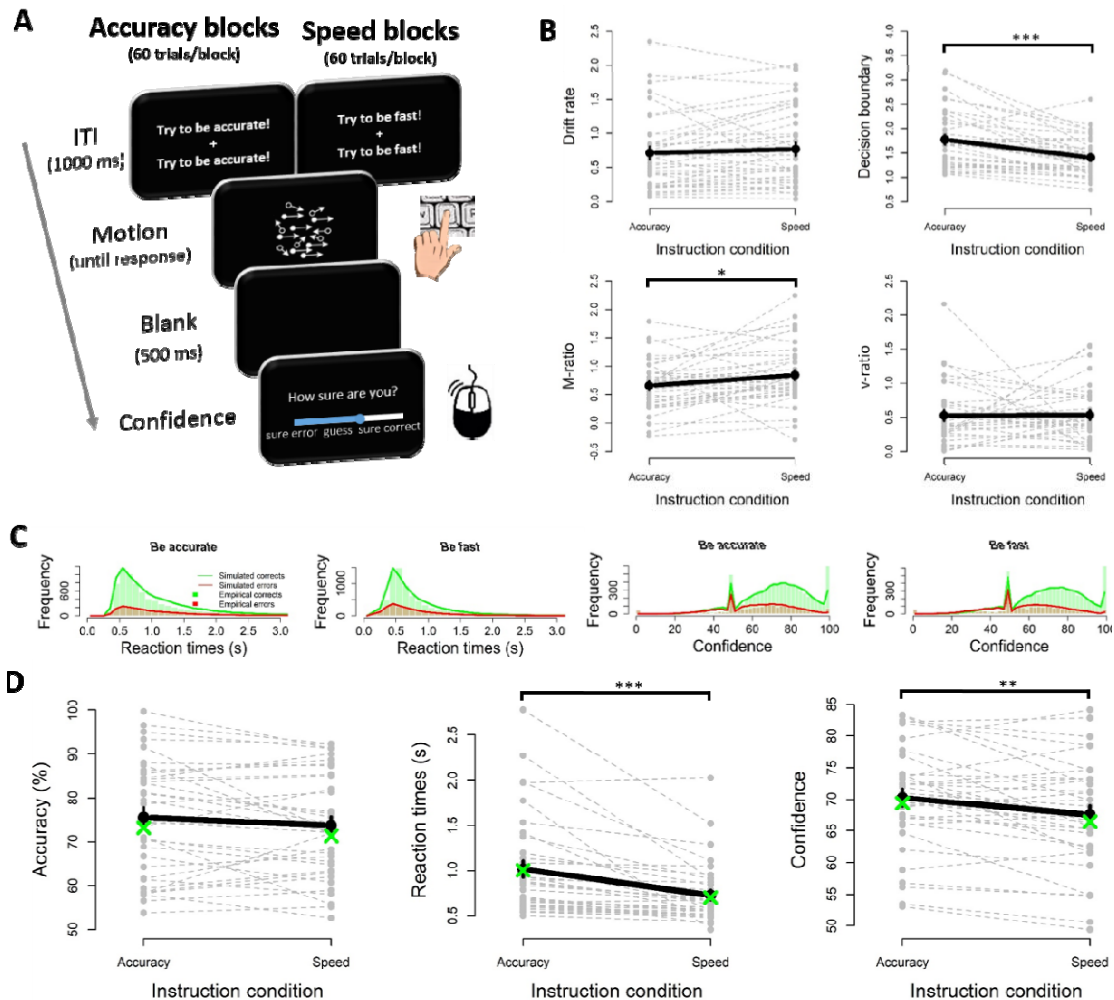
131    matrix is shown in Table 1.

|  | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1. Drift rate | - | | | | |
| 2. Non-decision time | .05 | - | | | |
| 3. Decision boundary | -.05 | .09 | - | | |
| 4. V-ratio | .01 | .01 | .00 | - | |
| 5. M-ratio | .03 | -.08 | -.52[***] | .58[***] | - |

132

133    *Table 1. Correlation table of the parameters from the model simulation. Note: ***<.001*

134    **Experiment 1: Explicit speed-accuracy instructions affect static but not dynamic measures**

135    **of confidence**

136         Next, we tested these model predictions in an experiment with human participants. We recruited

137    36 human participants who performed a task that has been widely used in the study of evidence

138    accumulation models: discrimination of the net motion direction in dynamic random dot displays [21].

139    Participants were asked to decide whether a subset of dots was moving coherently towards the left or the

140    right side of the screen (See Figure 2A). The percentage of dots that coherently moved towards the left or

141    right side of the screen (controlling decision difficulty) was held constant throughout the experiment at

142    20%. After their choice, and a blank screen, participants indicated their level of confidence using a

143    continuous slider. Critically, in each block, participants either received the instruction to focus on

144    accuracy ("try to decide as accurate as possible"), or to focus on speed ("try to decide as fast as

145    possible"). Consistent with the instructions, participants were faster in the speed condition than in the

146    accuracy condition, $M_{speed} = 727$ms versus $M_{accuracy} = 1014$ms, $t(35) = 4.47$, $p < .001$, and numerically

147     more accurate in the accuracy condition than in the speed condition, $M_{accurate} = 75.6\%$ vs $M_{speed} = 73.8\%$,

148     $t(35) = 1.63$, $p = .111$. Participants were also more confident in the accuracy condition than in the speed

149     condition, $M_{accuracy} = 70$ versus $M_{speed} = 67$, $t(35) = 3.57$, $p = .001$ (See Figure 2D).



150

151             *Figure 2. The influence of speed-accuracy instructions on metacognitive accuracy*

152     *(**Experiment 1**). **A.** Sequence of events in the experimental task. Participants decided whether the*

153     *majority of dots were moving left or right, by pressing "E" or "T" with their left hand. After a short*

154     *blank, they then indicated their level of confidence on a continuous scale. Depending on the block,*

155     *instructions during the ITI were either to focus on accuracy or to focus on speed. **B.** Fitted parameters of*

156     *a drift diffusion model with additional post-decision accumulation. Fitted decision boundaries were lower*

157     *in the speed vs accuracy condition, whereas drift rates did not differ. Critically, M-ratio was higher in the*

158     *speed vs accuracy condition whereas v-ratio did not differ between both instruction conditions. **C.***

159     *Distribution of reaction times and confidence for empirical data (bars) and model fits (lines), separately*

160    *for corrects (green) and errors (red). **D**. Participants were faster, less accurate and less confident when*

161    *instructed to focus on speed rather than on accuracy. Note: grey lines show individual data points; black*

162    *lines show averages; green dots show model fits; error bars reflect SEM; \*\*\*p<.001, \*\*p<.01, \*p<.05.*
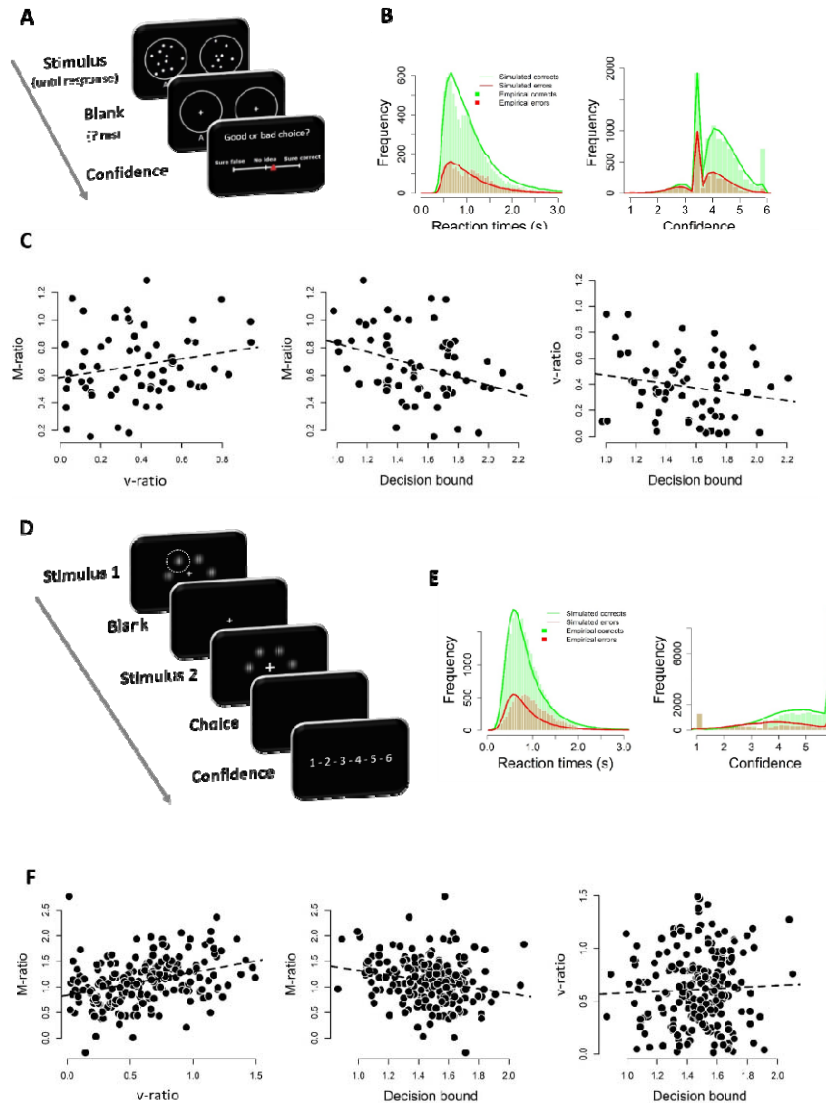
163

164        To shed further light on the underlying cognitive processes, we fitted these data using the

165    evidence accumulation model described in Figure 1A. The basic architecture of our model was a DDM, in

166    which noisy perceptual evidence accumulates over time until a decision boundary is reached. Afterwards,

167    evidence continued to accumulate for a specified amount of time [19]. In addition to drift rate, decision

168    boundary and non-decision time, our model featured a free parameter controlling the strength of the post-

169    decision evidence accumulation (v-ratio, reflecting the ratio between post-decision drift rate and drift rate)

170    and two further parameters controlling the mapping from *p*(*correct*) onto the confidence scale (see

171    Methods). Generally, our model fitted the data well, as it captured the distributional properties of both

172    reaction times and decision confidence (see Figure 2C). As a first sanity check, we confirmed that

173    decision boundaries were indeed different between the two instruction conditions, $M_{speed} = 1.40$ versus

174    $M_{accuracy} = 1.77$, $t(35) = 4.60$, $p < .001$, suggesting that participants changed their decision boundaries as

175    instructed. Also non-decision time tended to be a bit shorter in the speed condition compared to the

176    accuracy condition, $M_{speed} = 309$ms versus $M_{accuracy} = 390$ms, $t(35) = 3.19$, $p = .003$. Drift rates did not

177    differ between both instruction conditions, $p = .368$. There was a small but significant difference between

178    the two instruction conditions in the two additional parameters controlling the idiosyncratic mapping

179    between *p*(*correct*) and the confidence scale, reflecting that in the accuracy condition confidence

180    judgments were slightly higher, $t(35) = 2.506$, $p = .017$, and less variable, $t(35) = 2.206$, $p = .034$,

181    compared to the speed condition.

182        We next focused on metacognitive accuracy in both conditions (see Figure 2B). In line with the

183    model simulations, our data showed that M-ratio was significantly affected by the speed-accuracy tradeoff

184    instructions, $M_{speed} = 0.84$ versus $M_{accuracy} = 0.66$, $t(35) = 2.26$, $p = .030$. Moreover, apart from these

185    between-condition differences we also observed significant correlations between M-ratio and decision

186    boundary both in the accuracy condition, $r(34) = -.36$, $p = .030$, and in the speed condition, $r(34) = -.53$, $p$

187    $< .001$. Consistent with the notion that metacognitive accuracy should not be affected by differences in

188    decision boundary, v-ratio did not differ between both instruction conditions, $p = .938$.

189        **Experiment 2: Spontaneous differences in response caution relate to static but not dynamic**

190    **measures of metacognitive accuracy**

191        Although Experiment 1 provides direct evidence that changes in decision boundary affect M-

192    ratio, it remains unclear to what extent this is also an issue in experiments without speed stress. Notably,

193    in many metacognition experiments, participants do not receive the instruction to respond as fast as

194    possible. Nevertheless, it remains possible that participants implicitly decide on a certain level of response

195    caution. For example, a participant who is eager to finish the experiment quickly might adopt a lower

196    decision boundary compared to a participant who is determined to perform the experiment as accurate as

197    possible, thus leading to a natural across-subject variation in decision boundaries. To examine this

198    possibility, in Experiment 2 we analyzed data from an experiment in which participants (N = 63) did not

199    receive any specific instructions concerning speed or accuracy. Participants decided which of two boxes

200    contained more dots, and afterwards indicated their level of confidence on a continuous scale (see Figure

201    3A). The same evidence accumulation model as before was used to fit these data, and again this model

202    captured both reaction times and decision confidence distributions (Figure 3B). Consistent with our

203    model simulations, model fits showed a positive correlation between M-ratio and v-ratio, $r(61) = .21$, $p =$

204    .092, although this correlation was not statistically significant (Figure 3C). However, we again observed

205    that M-ratio correlated negatively with the fitted decision boundary, $r(61) = -.34$, $p = .006$, whereas v-

206    ratio did not, $r(61) = -.19$, $p = .129$.

207

*Figure 3. The influence of spontaneous variations in speed-accuracy tradeoff on metacognitive*

*accuracy. **A.** Sequence of events in Experiment 2. On each trial participants decided which of the two*

*circles contained more dots. Afterwards, they indicated their level of confidence on a continuous scale.*

*Note that participants did not receive any instructions concerning speed or accuracy. **B.** Distribution of*

*reaction times and confidence for Experiment 2, using the same conventions as in Figure 2. **C.** The data*

*of Experiment 2 showed a non-significant positive relation between M-ratio and v-ratio (r=.21).*

*Critically, only M-ratio correlated negatively with decision boundary (r=-.34) whereas this relation was*

*not significant for v-ratio (r= -.19). **D.** Sequence of events in Experiment 3. On each trial, participants*

*decided in which temporal interval (first or second) one of the Gabor patches had a higher contrast. After*

*this choice, participants indicated confidence on a continuous scale. **E.** Distribution of reaction times and*

*confidence for Experiment 3, using the same conventions as in Figure 2. **F.** The data of Experiment 3*

219    *showed a significant positive relation between M-ratio and v-ratio (r=.38) and a significant negative*

220    *correlation between M-ratio and decision boundary (r=-.18) but not between v-ratio and decision*

221    *boundary (r= -.04).*

222

223    **Experiment 3: Replication in an independent dataset**

224    To assess the robustness of our findings, in Experiment 3 we aimed to replicate our analysis in an

225    independent dataset with high experimental power. To achieve this, we searched the confidence database

226    [32] for studies with high power (N > 100) in which a 2CRT task was performed with separate confidence

227    ratings given on a continuous scale. Moreover, because our fitting procedure was not designed for

228    multiple levels of difficulty, we focused on studies with a single difficulty level. We identified one study

229    that satisfied all these constraint (Figure 3D; Prieto, Reyes & Silva, *under review*). Their task was highly

230    similar to the one reported above, but their high experimental power (N=204) assured a very sensitive

231    analysis of our claims. Consistent with the previous analysis, model fits on this independent dataset

232    showed a positive and statistically significant correlation between M-ratio and v-ratio, $r(202) = .38$, $p <$

233    .001, suggesting that both variables capture shared variance reflecting metacognitive accuracy (see Figure

234    3F). We again observed that M-ratio correlated negatively with the fitted decision boundary, $r(202) = -$

235    $.18$, $p = .009$, whereas no relation with decision bound was found for v-ratio, $r(202) = .04$, $p = .535$.

**Discussion**

236

237        Metacognitive accuracy is a quickly emerging field in recent years. Crucial to its study is a

238    method to objectively quantify the extent to which participants are able to detect their own mistakes,

239    regardless of decision strategy. We here report that a commonly used *static* measure of metacognitive

240    accuracy (M-ratio) highly depends on the decision boundary – reflecting decision strategy – that is set for

241    decision making. This was the case in simulation results, in an experiment explicitly manipulating the

242    tradeoff between speed and accuracy, and in two datasets in which participants received no instructions

243    concerning speed or accuracy. We propose an alternative, *dynamic,* measure of metacognitive accuracy

244    (v-ratio) that does not depend on decision boundary.

245

246        **Caution is warranted with static measures of metacognition**

247        The most important consequence of the current findings is that researchers should be cautious

248    when interpreting static measures of metacognitive accuracy, such as M-ratio. In the following, we will

249    discuss several examples where our finding might have important implications. In the last decade there

250    has been quite some work investigating to what extent the metacognitive evaluation of choices is a

251    domain-general process or not. These studies often require participants to perform different kinds of tasks,

252    and then examine correlations in accuracy and in metacognitive accuracy between these tasks [3,11–14,33]. For

253    example, Mazancieux and colleagues [11] asked participants to perform an episodic memory task, a

254    semantic memory task, a visual perception task and a working memory task. In each task, participants

255    rated their level of confidence after a decision. The results showed that whereas correlations between

256    accuracy on these different tasks were limited, there was substantial covariance in metacognitive accuracy

257    across these domains. Because in this study participants received no time limit to respond, it remains

258    unclear whether this finding can be interpreted as evidence for a domain-general metacognitive monitor,

259    or instead a domain-general response caution which caused these measures to correlate. Another popular

260    area of investigation has been to unravel the neural signatures supporting metacognitive accuracy [13,14,34–

261    36]. For example, McCurdy et al. observed that both visual and memory metacognitive accuracy correlated

262    with precuneus volume, potentially pointing towards a role of precuneus in both types of metacognition.

263    It remains unclear, however, to what extent differences in response caution might be responsible for this

264    association. Although differences in response caution are usually found to be related to pre-SMA and

265    anterior cingulate [24,25], there is some suggestive evidence linking precuneus to response caution [37].

266    Therefore, it is important that future studies on neural correlates of metacognition rule out the possibility

267    that their findings are caused by response caution. Finally, our study has important consequences for

268    investigations into differences in metacognitive accuracy between specific, e.g. clinical, groups. For

269    example, Folke and colleagues [17] reported that M-ratio was reduced in a group of bilinguals compared to

270     a matched group of monolinguals. Interestingly, they also observed that on average bilinguals had shorter

271     reaction times than monolinguals, but this effect was unrelated to the group difference in M-ratio.

272     Because these authors did not formally model their data using evidence accumulation models, however, it

273     remains unclear whether this RT difference results from a difference in boundary, and if so to what extent

274     this explains the difference in M-ratio between both groups that was observed. In a similar vein,

275     individual differences in M-ratio have been linked to psychiatric symptom dimensions, and more

276     specifically to a symptom dimension related to depression and anxiety [5]. At the same time, it is also

277     known that individual differences in response caution are related to a personality trait known as *need for*

278     *closure* [38]. Given that need for closure is, in turn, related to anxiety and depression [39], it remains a

279     possibility that M-ratio is only indirectly related to these psychiatric symptoms via response caution.

280

281     **The potential of dynamic measures of metacognition**

282     In order to control for potential influences of response caution on measures of metacognitive

283     accuracy, one approach could be to estimate the decision boundary and examine whether the relation

284     between metacognitive accuracy and the variable of interest remains when controlling for decision

285     boundary (e.g., using mediation analysis). However, a more direct approach would be to estimate

286     metacognitive accuracy in a dynamic framework, thus taking into account differences in response caution.

287     In the current work, we proposed v-ratio (reflecting the ratio between post-decision drift rate and drift

288     rate) as such a dynamic measure of metacognitive accuracy (following the observation that post-decision

289     drift rate indexes how accurate confidence judgments are[19,20]). In both simulations and empirical data, we

290     observed a positive relation between v-ratio and M-ratio, suggesting they capture shared variance.

291     Critically, v-ratio was not correlated with decision boundary, suggesting it is not affected by differences

292     in response caution. Thus, our dynamic measure of metacognition holds promise as a novel approach to

293     quantify metacognitive accuracy while taking into account the dynamics of decision making.

294     In our approach we allowed the drift rate and the post-decision drift rate to dissociate. This

295     proposal is in line with the view of metacognition as a second-order process whereby dissociations

296     between confidence and accuracy might arise because of noise or bias at each level [40–42]. However, when

297     formulating post-decision drift rate as a continuation of evidence accumulation, it remains underspecified

298     which evidence the post-decision accumulation process is exactly based on. It has been suggested that

299     participants can accumulate evidence that was still in the processing pipeline (e.g. in a sensory buffer)

300     even after a choice was made [30,43]. However, it is not very likely that this is the only explanation,

301     particularly in tasks without much speed stress. One other likely possibility, is that during the post-

302     decision process, participants resample the stimulus from short-term memory [44]. Because memory is

303     subject to decay, dissociations between the post-decision drift rate and the drift rate can arise. Other

304     sources of discrepancy might be contradictory information quickly dissipating from memory [45] which

305     should lower metacognitive accuracy, or better assessment of encoding strength with more time [46] which

306     should increase metacognitive accuracy.

307        To sum up, we provided evidence from simulations and empirical data that a common static

308     measure of metacognition, M-ratio, is confounded with response caution. We proposed an alternative

309     measure of metacognition based on a dynamic framework, v-ratio, which is insensitive to variations in

310     caution, and may thus be suitable to study how metacognitive accuracy varies across subjects and

311     conditions.

312 **Methods**

313 **Computational model**

314 **Simulations**

315 Data were simulated for 100 observers with 500 trials each. For each simulated observer, we

316 randomly selected a value for the drift rate (uniform distribution between 0 and 2.5), for the decision

317 boundary (uniform distribution between .5 and 3), for the non-decision time (uniform distribution

318 between .2 and .6) and for the v-ratio (uniform distribution between 0 and 1.5; see below for details). To

319 estimate meta-*d'*, data is needed for both of the possible stimuli (i.e., to estimate bias); therefore, for half

320 of the trials we multiplied the drift rate by -1. Finally, we fixed the values for starting point ($z = .5$),

321 within-trial noise ($\sigma = 1$) and post-decision processing time (1s).

322 **Fitting procedure**

323 We coded an extension of the drift diffusion model (DDM) that simultaneously fitted choices,

324 reaction times and decision confidence. The standard DDM is a popular variant of sequential sampling

325 models of two-choice tasks. We used a random walk approximation, coded in the rcpp R package to

326 increase speed [47], in which we assumed that noisy sensory evidence started at $z*a$; 0 and a are the lower

327 and upper boundaries, respectively, and $z$ quantifies bias in the starting point ($z = .5$ means no bias). At

328 each time interval $\tau$ a displacement $\Delta$ in the integrated evidence occurred according to the formula shown

329 in equation (1):

$$\Delta = v * \tau + \sigma * \sqrt{\tau} * \mathcal{N}(0,1)$$
(1)

330 Evidence accumulation strength is controlled by $v$, representing the drift rate, and within-trial

331 variability, $\sigma$, was fixed to 1. The random walk process continued until the accumulated evidence crossed

332 either 0 or $a$. After boundary crossing, the evidence continued to accumulate for a duration depending on

333 the participant-specific median confidence reaction time. Importantly, consistent with the signal detection

334 theoretical notion that primary and secondary evidence can dissociate, we allowed for dissociations

335 between the drift rate governing the choice and the post-decision drift rate. For compatibility with the M-

336 ratio framework, we quantified metacognitive accuracy as the ratio between post-decision drift rate and

337 drift rate, as shown in equation (2):

$$v\text{-}ratio = \frac{post-decision\ drift\ rate}{drift\ rate}$$
(2)

338 As a consequence, when v-ratio = 1, this implies that post-decision drift and drift are the same.

339 When v-ratio = .5, the magnitude of the post-decision drift rate is half the magnitude of the drift rate. To

340　calculate decision confidence, we first quantified for each trial the probability of being correct given

341　evidence, time, and choice. The heat map representing $p(correct|e, t, X)$ is shown in Figure 1A, and was

342　constructed by means of 300.000 random walks without absorbing bounds, with drift rates sampled from

343　a uniform distribution between zero and ten. This assured sufficient data points across the relevant part of

344　the heat map. Subsequently, the average accuracy was calculated for each (response time, evidence,

345　choice) combination, based on all trials that had a data point for that (response time, evidence, choice)

346　combination. Smoothing was achieved by aggregating over evidence windows of .01 and $\tau$ windows of 3.

347　Next, to take into account idiosyncratic mappings of $p(correct|e, t, X)$ onto the confidence scale used in

348　the experiment, we added two extra free parameters that controlled the mean (M) and the width (SD) of

349　confidence estimates, as shown in equation (3):

$$confidence = \frac{p(correct|e_{t+s}, t + s, X) + M}{SD} \qquad (3)$$

350　　　Although empirical confidence distributions appeared approximately normally distributed, there

351　was an over-representation of confidence values at the boundaries (1 and 100 in Experiment 1; 1 and 6 in

352　Experiments 2 and 3) and in the middle of the scale (50 in Experiment 1, 3.5 in Experiment 2). Most

353　likely, this resulted from the use of verbal labels placed at exactly these values. To account for frequency

354　peaks at the endpoints of the scale, we relabeled predicted confidence values that exceeded the endpoints

355　of the scale as the corresponding endpoint (e.g., in Experiment 1 a predicted confidence value of 120 was

356　relabeled as 100), which naturally accounted for the frequency peaks at the endpoints. To account for

357　peaks in the center of the scale, we assumed that confidence ratings around the center were pulled towards

358　the center value. Specifically, we relabeled $P$% of trials around the midpoint as the midpoint (e.g., in

359　Experiment 1, $P = 10$% implies that 10% of the data closest to 50 were (re)labeled as 50). Note that $P$ was

360　not a free parameter, but instead its value was taken to be the participant-specific proportion based on the

361　empirical data. Note that the main conclusions reported in this manuscript concerning the relation

362　between M-ratio, decision boundary and post-decision drift rate, remain the same in a reduced model

363　without $P$, and also in a reduced model without $P$, $M$ and $SD$. Because these reduced models did not

364　capture confidence distributions very well though, we here report only the findings of the full model.

365　　　To estimate these 6 parameters ($v, a, Ter, v$-ratio, $M,$ and $SD$) based on choices, reaction times

366　and decision confidence, we implemented quantile optimization. Specifically, we computed the

367　proportion of trials in quantiles .1, .3, .5, .7, and .9, for both reaction times and confidence; separately for

368　corrects and errors (maintaining the probability mass of corrects and errors, respectively). We then used

369　differential evolution optimization, as implemented in the DEoptim R package [48], to estimate these 6

370　parameters by minimizing the chi square error function shown in equation (4):

$$x^2 = \sum \frac{(oRT_i - pRT_i)^2}{pRT_i} + \sum \frac{(oCJ_i - pCJ_i)^2}{pCJ_i}$$

(4)

371    with $oRT_i$ and $pRT_i$ corresponding to the proportion of observed/predicted responses in quantile $i$,

372    separately calculated for corrects and errors both reaction times, and $oCJ_i$ and $pCJ_i$ reflecting their

373    counterparts for confidence judgments. We set $\tau$ to 1e-2. Model fitting was done separately for each

374    participant. Note that in Experiment 3 there was no clear peak in the middle of the scale so $P$ was fixed to

375    0 in that experiment.

376    **Parameter recovery**

377        To assure that our model was able to recover the parameters, we here report parameter recovery.

378    In order to assess parameter recovery with a sensible set of parameter combinations, we used the fitted

379    parameters of Experiment 1 (N = 36), simulated data from these parameters with a varying number of

380    trials, and then tested whether our model could recover these initial parameters. As a sanity check, we

381    first simulated a large number of trials (25000 trials per participant), which as expected provided excellent

382    recovery for all six parameters, $r$s > .97. We then repeated this process with only 200 trials per

383    participants, which was the trial count in Experiment 2 (note that Experiment 1 and 3 both had higher trial

384    counts). Recovery for v-ratio was still very good, $r$ = .85, whereas it remained excellent for all other

385    parameters, $r$s > .98.

386

387    **Experiment 1**

388    **Participants**

389        Forty healthy participants (18 males) took part in Experiment 1 in return for course credit (mean

390    age = 19.82, between 18 and 30). All reported normal or corrected-to-normal vision. Two participants

391    were excluded because they required more than 10 practice blocks in one of the training blocks (see

392    below) and two participants were excluded because their accuracy, averaged per block and then compared

393    against chance level using a one-sample t-test, was not significantly above chance level. The final sample

394    thus comprised thirty-six participants. All participants provided their informed consent and all procedures

395    adhered to the general ethical protocol of the ethics committee of the Faculty of Psychology and

396    Educational Sciences of Ghent University.

397    **Stimuli and apparatus**

398        The data for Experiment 1 were collected in an online study, due to COVID-19. Participants were

399    allowed to take part in the experiment only when they made us of an external mouse. Choices were

400    provided with the keyboard, and decision confidence was indicated with the mouse. Stimuli in

401    Experiment 1 consisted of 50 randomly moving white dots (radius: 2 pixels) drawn in a circular aperture

402    on a black background centered on the fixation point. Dots disappeared and reappeared every 5 frames.

403    The speed of dot movement (number of pixel lengths the dot will move in each frame) was a function of

404    the screen resolution (screen width in pixels / 650).

405    **Task procedure**

406        Each trial started with the presentation of a fixation cross for 1000 ms. Above and below this

407    fixation cross specific instructions were provided concerning the required strategy. In accuracy blocks the

408    instruction was to respond as accurately as possible; in speed blocks the instruction was to respond as fast

409    as possible. The order of this block-wise manipulation was counterbalanced across participants. Next,

410    randomly moving dots were shown on the screen until a response was made or the response deadline was

411    reached (max 5000 ms). On each trial, 20% of the dots coherently moved towards the left or the right side

412    of the screen, with an equal number of leftward and rightward movement trials in each block. Participants

413    were instructed to decide whether the majority of dots was moving towards the left or the right side of the

414    screen, by pressing "E" or "T", respectively, with their left hand. After their response, a blank screen was

415    shown for 500 ms, followed by the presentation of a continuous confidence scale. Below the scale the

416    labels "Sure error", "guess", and "sure correct" were shown, arranged outer left, centrally and outer right,

417    respectively. After clicking the confidence scale, participants had to click a centrally presented

418    "Continue" button (below the confidence scale) that ensured that the position of the mouse was central

419    and the same on each trial.

420        The main part of Experiment 1 consisted of 10 blocks of 60 trials, half of which were from the

421    accuracy instruction condition and half from the speed instruction condition. The experiment started with

422    24 practice trials during which participants only discriminated random dot motion at 50% coherence, no

423    confidence judgments were asked. This block was repeated until participants achieved 85% accuracy

424    (mean = 2 blocks). Next, participants completed again 24 practice trials with the only difference that now

425    the coherence was decreased to 20% (mean = 1.05 blocks). When participants achieved 60% accuracy,

426    they then performed a final training block of 24 trials during which they practiced both dot discrimination

427    and indicated their level of confidence (mean = 1.05 blocks).

428    **Experiment 2**

429    Full experimental details are described in Drescher et al. [49]. On each trial participants were

430    presented with two white circles (5.1° diameter) on a black background, horizontally next to each other

431    with a distance of 17.8° between the midpoints. Fixation crosses were shown for 1s in each circle,

432    followed by dots clouds in each circle for 700ms. The dots had a diameter of 0.4°. Dot positions in the

433    boxes, as well as the position of the box containing more dots were randomly selected on each trial.

434    Participants indicated which circle contained more dots by pressing "S" or "L" on a keyboard. Then, the

435    question "correct or false?" appeared on the screen, with a continuous confidence rating bar, with the

436    labels "Sure false", "No idea", and "Sure correct". Participants moved a cursor with the same keys as

437    before, and confirmed their confidence judgment with the enter key. No time limit was imposed for both

438    primary choice and confidence rating. Subjects received several practice trials (10 without confidence

439    rating, 14 with confidence rating), before they completed eight experimental blocks of 25 trials.

440    **Experiment 3**

441    Data from this experiment were taken from the confidence database [50], a collection of openly

442    available studies on decision confidence. In this experiment, Prieto, Reyes and Silva (unpublished paper),

443    used the same task as described in Fleming and colleagues [2]. Each participant (N=204 woman, aged 18-

444    35) completed 50 practice trials, followed by 5 blocks of 200 trials.

445    **Data and code availability**

446    All data and analysis code have been deposited online and can be freely accessed (insert link

447    upon publication).

448    **Acknowledgments**

454

## References

455

456  1.  Sanders, J. I., Hangya, B. & Kepecs, A. Signatures of a Statistical Computation in the Human
457      Sense of Confidence. *Neuron* **90**, 499–506 (2016).

458  2.  Fleming, S. M., Weil, R. S., Nagy, Z., Dolan, R. J. & Rees, G. Relating introspective accuracy to
459      individual differences in brain structure. *Science (80-. ).* **329**, 1541–3 (2010).

460  3.  Baird, B., Mrazek, M. D., Phillips, D. T. & Schooler, J. W. Domain-Specific Enhancement of
461      Metacognitive Ability Following Meditation Training. *J. Exp. Psychol. Gen.* **143**, 1972–1979
462      (2014).

463  4.  Rollwage, M. *et al.* Metacognitive failure as a feature of those holding radical political beliefs.
464      *Nat. Hum. Behav.* **2**, 637–644 (2018).

465  5.  Rouault, M., Seow, T., Gillan, C. M. & Fleming, S. M. Psychiatric Symptom Dimensions Are
466      Associated With Dissociable Shifts in Metacognition but Not Task Performance. *Biol. Psychiatry*
467      **84**, 443–451 (2018).

468  6.  Nelson, T. O. A comparison of current measures of the accuracy of feeling-of-knowing
469      predictions. *Psychol. Bull.* **95**, 109–133 (1984).

470  7.  Sandberg, K., Timmermans, B., Overgaard, M. & Cleeremans, A. Measuring consciousness: is
471      one measure better than the other? *Conscious. Cogn.* **19**, 1069–78 (2010).

472  8.  Fleming, S. M. HMeta-d: hierarchical Bayesian estimation of metacognitive efficiency from
473      confidence ratings. *Neurosci. Conscious.* **3**, 1–14 (2017).

474  9.  Maniscalco, B. & Lau, H. A signal detection theoretic approach for estimating metacognitive
475      sensitivity from confidence ratings. *Conscious. Cogn.* **21**, 422–30 (2012).

476  10. Fleming, S. M. & Lau, H. C. How to measure metacognition. *Front. Hum. Neurosci.* **8**, 1–9
477      (2014).

478  11. Mazancieux, A. *et al.* Is There a G Factor for Metacognition? Correlations in Retrospective
479      Metacognitive Sensitivity Across Tasks Is There a G Factor for Metacognition? Correlations in
480      Retrospective Metacognitive Sensitivity Acros. *J. Exp. Psychol. Gen.* (2020).

481  12. Rouault, M., McWilliams, A., Allen, M. G. & Fleming, S. M. Human Metacognition Across
482      Domains: Insights from Individual Differences and Neuroimaging. *Personal. Neurosci.* **1**, (2018).

483  13. McCurdy, L. Y. *et al.* Anatomical coupling between distinct metacognitive systems for memory

484      and visual perception. *J. Neurosci.* **33**, 1897–906 (2013).

485   14.   Fleming, S. M., Ryu, J., Golfinos, J. G. & Blackmon, K. E. Domain-specific impairment in
486          metacognitive accuracy following anterior prefrontal lesions. *Brain* **137**, 2811–2822 (2014).

487   15.   Filevich, E., Dresler, M., Brick, X. T. R. & Ku, S. Metacognitive Mechanisms Underlying Lucid
488          Dreaming. *J. Neurosci.* **35**, 1082–1088 (2015).

489   16.   Rahnev, D. a, Maniscalco, B., Luber, B., Lau, H. C. & Lisanby, S. H. Direct injection of noise to
490          the visual cortex decreases accuracy but increases decision confidence. *J. Neurophysiol.* **107**,
491          1556–1563 (2012).

492   17.   Folke, T., Ouzia, J., Bright, P., De Martino, B. & Filippi, R. A bilingual disadvantage in
493          metacognitive processing. **150**, 119–132 (2016).

494   18.   Scheffers, M. K. & Coles, M. G. Performance monitoring in a confusing world: error-related brain
495          activity, judgments of response accuracy, and types of errors. *J. Exp. Psychol. Hum. Percept.*
496          *Perform.* **26**, 141–151 (2000).

497   19.   Pleskac, T. J. & Busemeyer, J. R. Two-stage dynamic signal detection: A theory of choice,
498          decision time, and confidence. *Psychol. Rev.* **117**, 864–901 (2010).

499   20.   Yu, S., Pleskac, T. J. & Zeigenfuse, M. D. Dynamics of Postdecisional Processing of Confidence.
500          *J. Exp. Psychol. Gen.* **144**, 489–510 (2015).

501   21.   Gold, J. I. & Shadlen, M. N. The neural basis of decision making. *Annu. Rev. Neurosci.* **30**, 535–
502          561 (2007).

503   22.   Ratcliff, R. & McKoon, G. The Diffusion Decision Model□: Theory and Data for Two-Choice
504          Decision Tasks. *Neural Comput.* **20**, 873–922 (2008).

505   23.   Forstmann, B. U. & Wagenmakers, E.-J. Sequential Sampling Models in Cognitive Neuroscience:
506          Advantages, Applications, and Extensions. *Annu. Rev. Psychol.* **67**, 641–666 (2016).

507   24.   Bogacz, R., Wagenmakers, E.-J., Forstmann, B. U. & Nieuwenhuis, S. The neural basis of the
508          speed-accuracy tradeoff. *Trends Neurosci.* **33**, 10–6 (2010).

509   25.   Forstmann, B. U. *et al.* Striatum and pre-SMA facilitate decision-making under time pressure.
510          *Proc. Natl. Acad. Sci. U. S. A.* **105**, 17538–42 (2008).

511   26.   Kiani, R., Corthell, L. & Shadlen, M. N. Choice Certainty Is Informed by Both Evidence and
512          Decision Time. *Neuron* **84**, 1329–1342 (2014).

513   27.   Zylberberg, A., Fetsch, C. R. & Shadlen, M. N. The influence of evidence volatility on choice ,
514          reaction time and confidence in a perceptual decision. *Elife* **5**, 1–31 (2016).

515   28.   Drugowitsch, J., Moreno-Bote, R., Churchland, A. K., Shadlen, M. N. & Pouget, A. The Cost of
516          Accumulating Evidence in Perceptual Decision Making. *J. Neurosci.* **32**, 3612–3628 (2012).

517   29.   Moreno-Bote, R. Decision confidence and uncertainty in diffusion models with partially correlated
518          neuronal integrators. *Neural Comput.* **22**, 1786–1811 (2010).

519   30.   Van Den Berg, R. *et al.* A common mechanism underlies changes of mind about decisions and
520          confidence. *Elife* 1–21 (2016) doi:10.7554/eLife.12192.

521   31.   Desender, K., Donner, T. H. & Verguts, T. Dynamic expressions of confidence within an evidence
522          accumulation framework. *bioRxiv* 1–29 (2020).

523   32.   Rahnev, D. *et al.* The Confidence Database. *Nat. Hum. Behav.* **4**, (2020).

524   33.   Song, C. *et al.* Relating inter-individual differences in metacognitive performance on different
525          perceptual tasks. *Conscious. Cogn.* **20**, 1787–92 (2011).

526   34.   Bor, D., Schwartzman, D. J., Barrett, A. B. & Seth, A. K. Theta-burst transcranial magnetic
527          stimulation to the prefrontal or parietal cortex does not impair metacognitive visual awareness.
528          *PLoS One* **1**, 165–175 (2016).

529   35.   Rounis, E., Maniscalco, B., Rothwell, J. J. C., Passingham, R. R. E. & Lau, H. H. Theta-burst
530          transcranial magnetic stimulation to the prefrontal cortex impairs metacognitive visual awareness.
531          *Cogn. Neurosci.* **1**, 165–175 (2010).

532   36.   Baird, B., Cieslak, M., Smallwood, J., Grafton, S. T. & Schooler, J. W. Regional White Matter
533          Variation Associated with Domain-specific Metacognitive Accuracy. *J. Cogn. Neurosci.* 1–10
534          (2014) doi:10.1162/jocn.

535   37.   Van Maanen, L. *et al.* Neural Correlates of Trial-to-Trial Fluctuations in Response Caution. *J.*
536          *Neurosci.* **31**, 17488–17495 (2011).

537   38.   Evans, N. J., Rae, B., Bushmakin, M., Rubin, M. & Brown, S. D. Need for closure is associated
538          with urgency in perceptual decision-making. *Mem. Cogn.* **45**, 1193–1205 (2017).

539   39.   Freeman, D. *et al.* Delusions and decision-making style: Use of the Need for Closure Scale.
540          *Behav. Res. Ther.* **44**, 1147–1158 (2006).

541   40.   Fleming, S. M. & Daw, N. D. Self-evaluation of decision performance: A general Bayesian

542       framework for metacognitive computation. *Psychol. Rev.* **124**, 1–59 (2016).

543  41.  Pasquali, A., Timmermans, B. & Cleeremans, A. Know thyself: metacognitive networks and
544       measures of consciousness. *Cognition* **117**, 182–90 (2010).

545  42.  Balsdon, T., Wyart, V. & Mamassian, P. Confidence controls perceptual evidence accumulation.
546       *Nat. Commun.* **11**, (2020).

547  43.  Resulaj, A., Kiani, R., Wolpert, D. M. & Shadlen, M. N. Changes of mind in decision-making.
548       *Nature* **461**, 263–266 (2009).

549  44.  Vlassova, A. & Pearson, J. Look Before You Leap: Sensory Memory Improves Decision Making.
550       *Psychol. Sci.* **24**, 1635–1643 (2013).

551  45.  Minson, J. A. & Umphres, C. Confidence in Context: Perceived Accuracy of Quantitative
552       Estimates Decreases With Repeated Trials. *Psychol. Sci.* (2020) doi:10.1177/0956797620921517.

553  46.  Nelson, T. O. & Dunlosky, J. When People's Judgments of Learning are extremely accurate at
554       predicting subsequent recall: The 'Delayed-JOL effect'. *Psychol. Sci.* **2**, 267–270 (1991).

555  47.  Eddelbuettel, D. Seamless R and C++ integration with Rcpp. *Seamless R C++ Integr. with Rcpp*
556       **40**, 1–220 (2013).

557  48.  Mullen, K. M., Ardia, D., Gil, D. L., Windover, D. & Cline, J. DEoptim: An R package for global
558       optimization by differential evolution. *J. Stat. Softw.* **40**, 1–26 (2011).

559  49.  Drescher, L. H., Van den Bussche, E. & Desender, K. Absence without leave or leave without
560       absence: Examining the interrelations among mind wandering, metacognition, and cognitive
561       control. *PLoS One* (2018).

562  50.  Rahnev, D. *et al.* The Confidence Database. *Nat. Hum. Behav.* (2020) doi:10.1038/s41562-019-
563       0813-1.

564