**Measures of metacognitive efficiency across cognitive models of decision confidence**

Manuel Rausch[1,2], Sebastian Hellmann[1], and Michael Zehetleitner[1]

[1] Catholic University of Eichstätt-Ingolstadt

[2] Rhine-Waal University of Applied Sciences

**Author Note**

Manuel Rausch https://orcid.org/0000-0002-5805-5544

Sebastian Hellmann https://orcid.org/0000-0001-6006-5103

Michael Zehetleitner https://orcid.org/0000-0003-3363-2680

Correspondence should be addressed to Manuel Rausch, Hochschule Rhein-Waal. Fakultät für Gesellschaft und Ökonomie, Marie-Curie-Straße 1, 47533 Kleve, Germany. E-Mail: manuel.rausch@hochschule-rhein-waal.de

**Abstract**

Meta-d′/d′ has become the quasi-gold standard to quantify metacognitive efficiency because meta-d′/d′ was developed to control for discrimination performance, discrimination criteria, and confidence criteria even without the assumption of a specific generative model underlying confidence judgments. Using simulations, we demonstrate that meta-d′/d′ is not free from assumptions about confidence models: Only when we simulated data using a generative model of confidence according to which the evidence underlying confidence judgements is sampled independently from the evidence utilized in the choice process from a truncated Gaussian distribution, meta-d′/d′ was unaffected by discrimination performance, discrimination task criteria, and confidence criteria. According to five alternative generative models of confidence, there exist at least some combination of parameters where meta-d′/d′ is affected by discrimination performance, discrimination criteria and confidence criteria. A simulation using empirically fitted parameter sets showed that the magnitude of the correlation between meta-d′/d′ and discrimination performance, discrimination task criteria, and confidence criteria depends heavily on the generative model and the specific parameter set and varies between negligibly small and very large. These simulations imply that a difference in meta-d′/d′ between conditions does not necessarily reflect a difference in metacognitive efficiency but might as well be caused by a difference in discrimination performance, discrimination task criterion, or confidence criteria.

*Keywords:* Metacognition, metacognitive efficiency, confidence, cognitive modelling, signal detection theory, meta-d′/d′

47          **Metacognitive efficiency in cognitive models of decision confidence**

48          A key aspect of metacognition is metacognitive efficiency, defined as a subject's level of

49   metacognition given their discrimination task performance or signal processing capacity

50   (Fleming & Lau, 2014). The gold standard to measure of metacognitive efficiency is meta-d′/d′

51   (Maniscalco & Lau, 2012, 2014). Measuring metacognitive efficiency by meta-d′/d′ has inspired

52   research on many different psychological concepts, including learning (Boldt et al., 2019;

53   Hainguerlot et al., 2018; Taouki et al., 2022), cognitive control (Drescher et al., 2018), vigilance

54   (Maniscalco et al., 2017), memory (Mazancieux et al., 2020; Vandenbroucke et al., 2014),

55   perception (Maniscalco et al., 2016; Odegaard, Chang, et al., 2018), psychopathology (Bhome et

56   al., 2022; Culot et al., 2021; Muthesius et al., 2022; Rouault et al., 2018), beliefs about

57   politicised science (Fischer & Said, 2021; Said et al., 2022), and visual awareness (Charles et al.,

58   2013; Rausch & Zehetleitner, 2016; Vlassova et al., 2014). One reason why the meta-d′/d′

59   method has become so popular is that meta-d′ is believed to provide control over discrimination

60   performance, discrimination task criteria, and confidence criteria (Maniscalco & Lau, 2012,

61   2014), which is a key requirement for measures of metacognitive accuracy (Barrett et al., 2013).

62   Meta-d′ is also popular because it does not explicitly assume a specific generative model for

63   confidence judgments (Maniscalco & Lau, 2014). However, there each exists at least one

64   generative model of confidence which implies that meta-d'/d′ is affected by discrimination

65   performance (Guggenmos, 2021) and confidence criteria (Shekhar & Rahnev, 2021), raising the

66   question how robust meta-d′/d′ is with respect to the control over discrimination performance,

67   discrimination task criteria, and confidence criteria across different generative models of

68   confidence.

69      **The meta-d′/d′ method**

70      The meta-d′/d′ method is based on signal detection theory (Green & Swets, 1966;

71      Peterson et al., 1954; Tanner & Swets, 1954) and type 2 signal detection theory (Clarke et al.,

72      1959; Galvin et al., 2003; Pollack, 1959). The conceptual idea of meta-d′ is to quantify the

73      accuracy of metacognition in terms of discrimination sensitivity in a hypothetical signal

74      detection model describing the primary task, assuming participants had perfect access to the

75      sensory evidence underlying the discrimination choice and were perfectly consistent in placing

76      their confidence criteria (Maniscalco & Lau, 2012, 2014). Using a signal detection model

77      describing the primary task to quantify metacognitive accuracy has the advantage of allowing a

78      direct comparison between metacognitive accuracy and discrimination performance. Meta-d′ can

79      be compared against the estimate of the distance between the two stimulus distributions

80      estimated from discrimination responses, which is referred to as d′: If meta-d′ equals d′, it means

81      that metacognitive accuracy is exactly as good as expected from discrimination performance. If

82      meta-d′ is lower than d′, it means that metacognitive accuracy is worse than expected from

83      discrimination performance (Fleming & Lau, 2014; Maniscalco & Lau, 2012, 2014).

84      The hypothetical signal detection model underlying meta-d′ assumes that the observer

85      selects a binary response $R \in \{-1, 1\}$ about a stimulus characterised by two classes $S \in$

86      $\{-1, 1\}$ as well as a confidence rating out of an ordered set of confidence categories $C \in$

87      $\{1, 2, \ldots, n\}$ (see Table 1 for a list of our mathematical notation). For each presentation of the

88      stimulus, the observer's perceptual system creates sensory evidence delineating the two response

89      options. As there is noise in the system, the sensory evidence is not constant, but modelled as a

90      random sample x out of a separate Gaussian distribution for each of the two stimulus classes (see

91      Fig. 1). The distance d between the two distributions created by the two classes of S is

92    interpreted as the observer's ability to differentiate between the two kinds of S. Participants

93    select a response by comparing the sensory evidence x with a response criterion c, choosing R =

94    -1 if the sensory evidence x is smaller than the response criterion, and R = 1 otherwise.

95    Confidence ratings are chosen by comparing the same sample of sensory evidence $x$ against a set

96    of $2 \times n - 1$ confidence criteria, $\theta_1,\ \theta_2,\ \theta_3, \dots,\ \theta_{2\times n-1}$. For example, if there are four

97    confidence categories, participants are assumed to select a response R of 1 and a confidence level

98    of 3 if the sensory evidence x is smaller than the outermost response criterion $\theta_7$, but at the same

99    time greater than the second outermost response criterion $\theta_6$.

**Table 1**

*Table of mathematical notation and terminology*

| Symbol | Description or terminology |
|---|---|
| $S$ | Stimulus class |
| $R$ | Discrimination response about the stimulus class |
| $C$ | Confidence judgment |
| $n$ | Number of options given by the confidence scale |
| $x$ | Sensory evidence about S |
| $d$ | distance between the two distributions of evidence created by the two different stimulus classes, interpreted as the observer's ability to differentiate between the two stimulus classes |
| $d'$ | Estimate of d based on R |
| $d_{meta}$ | Meta-d′: Estimate of d based on C |
| $c$ | Response criterion for the discrimination judgment |
| $\theta$ | Criterion for confidence judgments |
| $m$ | Metacognitive efficiency parameter within the independent truncated Gaussian model |
| $y$ | Confidence decision variable |

100   **Figure 1**

101   *The hypothetical signal detection theoretic model underlying meta-d′*

102

103    *Note.* The hypothetical signal detection theoretic model describing the primary task underlying

104    meta-d′ (Maniscalco & Lau, 2012, 2014). To estimate meta-d′, it is assumed that the same

105    evidence is available for selecting a response for the discrimination task and for selecting a

106    confidence judgement. Primary task responses and confidence categories are assumed to form an

107    ordered set of responses delineated by a set of criteria θ.

**Meta-d′ vs. generative models of confidence**

109         According to Maniscalco and Lau (2014), the meta-d′/d′ method only makes assumptions

110    about the cognitive architecture underlying the discrimination choice, but meta-d′/d′ does not

111    require an *explicit* assumption about the generative model underlying confidence judgments.

112    However, it should be noted that the hypothetical signal detection model underlying meta-d′ is

113    not dissimilar to the approach taken in studies that aim to identify the generative model

114    underlying confidence judgments. The reason is that the estimation methods available to fit

115    meta–d′ require the computation of the probability of the different levels of confidence given

116    stimulus and discrimination response $p(C|R,S)$. Notably, static generative models of confidence

117    are usually defined by a probability density of confidence ratings and discrimination task

118    responses $p(C, R|S)$ (e.g. Adler & Ma, 2018; Aitchison et al., 2015; Rausch et al., 2018, 2020;

119    Shekhar & Rahnev, 2021). This means what distinguishes the meta-d′ approach from generative

120    models of confidence is whether the probability density is conditioned on the discrimination

121    response or whether the discrimination response is modelled as well. According to both the

122    conditioned maximum likelihood procedure proposed by Maniscalco and Lau (2014) and the

123    Bayesian Markov Chain Monte Carlo (MCMC) method by Fleming (2017), the probability for a

124    specific degree of confidence given stimulus and response $p(C|R, S)$ is given by

$$p(C = i|S, R = -1) = \frac{\int_{\theta_{n-i}}^{\theta_{n-i+1}} \phi_{\mu=d_{meta}\times S\times 0.5}(y)\, dy}{\int_{-\infty}^{\theta_n} \phi_{\mu=d_{meta}\times S\times 0.5}(y)\, dy} \tag{1}$$

$$p(C = i|S, R = 1) = \frac{\int_{\theta_{n+i-1}}^{\theta_{n+i}} \phi_{\mu=d_{meta}\times S\times 0.5}(y)\, dy}{\int_{\theta_n}^{\infty} \phi_{\mu=d_{meta}\times S\times 0.5}(y)\, dy} \tag{2}$$

125            where $\phi$ indicates the Gaussian density function with mean $\mu$ and variance of 1, $\theta_0$ is -∞,

126    $\theta_{2n}$ is ∞, and $d_{meta}$ is meta-d′. According to Maniscalco and Lau (2014), the location of the

127    central confidence criterion $\theta_n$ depends on the perceptual sensitivity of the observer d′ as well as

128    on the primary task criterion c and is given by $\theta_n = c \times d_{meta} \div d'$. According to Fleming's

129    method, $\theta_n$ is identical to c. The formulae (1) and (2) show two important features of the meta-

130    d′/d′ method. First, the formulae for $p(C|S, R)$ are identical to the cumulative truncated gaussian

131    distribution function (Kristensen et al., 2020). Second, the formulae do not include x, the sensory

132    evidence used to make the discrimination choice: This means that the random process underlying

133    confidence judgments only depends on the outcome of the random process underlying the

134    discrimination task decision, i.e., the response $R$, but when conditioned on R, it does not depend

135    on the state of the random process generating the discrimination task decision.
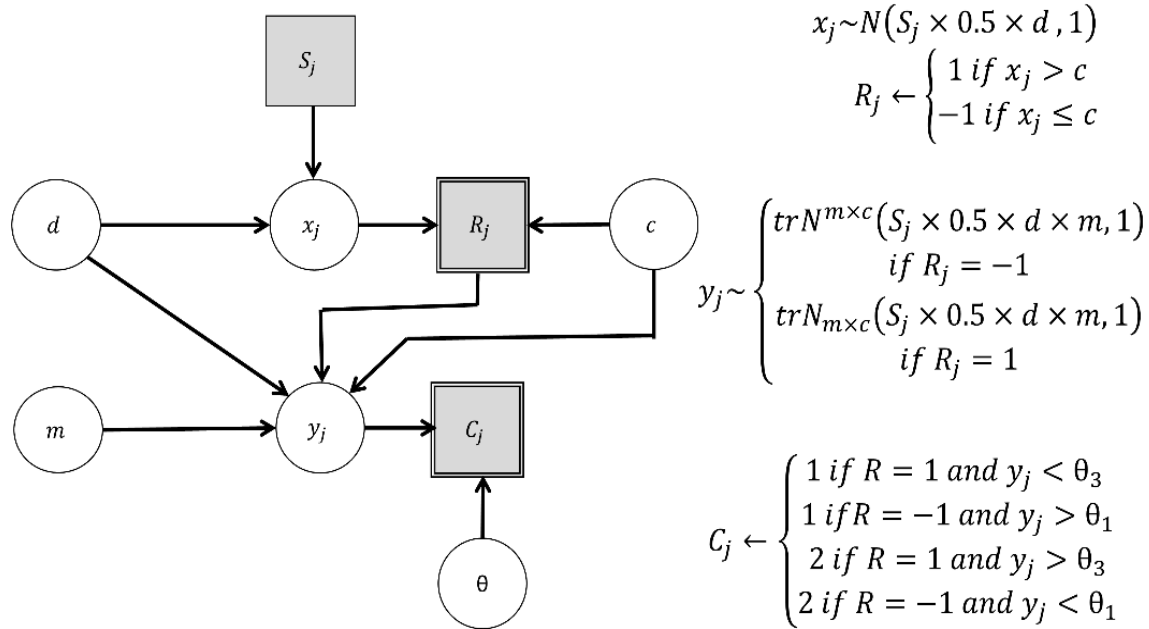
136     **The independent truncated Gaussian model (ITG)**

137          Here, we present a generative model of confidence that is set up to be consistent with the

138     probability functions used to estimate meta-d′: the *independent truncated Gaussian model* (ITG,

139     see Fig. 2). Conceptually, ITG reflects a cognitive mechanism where confidence judgments are

140     based on information generated independently from the sensory evidence used to make the

141     perceptual decision. However, according to ITG, confidence judgments can only be informed by

142     information corroborating the perceptual decision; contradicting information is not available.

143     ITG is identical to standard signal detection theory as far as the discrimination task response is

144     concerned. For the choice about the confidence, according to ITG, there is a separate decision

145     variable for confidence y. The confidence decision variable y is sampled from a truncated

146     Gaussian distribution, with the location parameter equal to $S \times d \times 0.5 \times m$ and a scale

147     parameter of 1. The parameter d quantifies the perceptual ability of the observer and is

148     equivalent to d′ in standard signal detection theory. The parameter m quantifies metacognitive

149     efficiency, which is measured by meta-d′/d′. Notably, y is sampled independently from x, the

150     sensory evidence used in the discrimination decision (see Fig. 3 for a visualisation of the

151     distribution of x and y). The Gaussian distribution of y is truncated in a way that it is impossible

152     to sample evidence that contradicts the original decision: If R = -1, the distribution is truncated to

153     the right of $\theta_n$. If R = 1, the distribution is truncated to the left of $\theta_n$. Because Maniscalco and

154     Lau (2014) and Fleming (2017) defined $\theta_n$ differently, there are also two slightly different

155     versions of ITG. ITG reproduces the probability density of confidence given stimulus and

156     response specified by Maniscalco and Lau (2014) if the distribution of y is truncated at $c \times m$,

157     while to reproduce the probability density of confidence given stimulus and response in Fleming

158     (2017), the distribution must be truncated at c. Just as in the signal detection model, confidence

159    ratings are chosen by comparing the confidence decision variable y against a set of $2 \times n - 1$

160    confidence criteria, $\theta_1, \theta_2, \theta_3, \ldots, \theta_{2 \times n - 1}$.
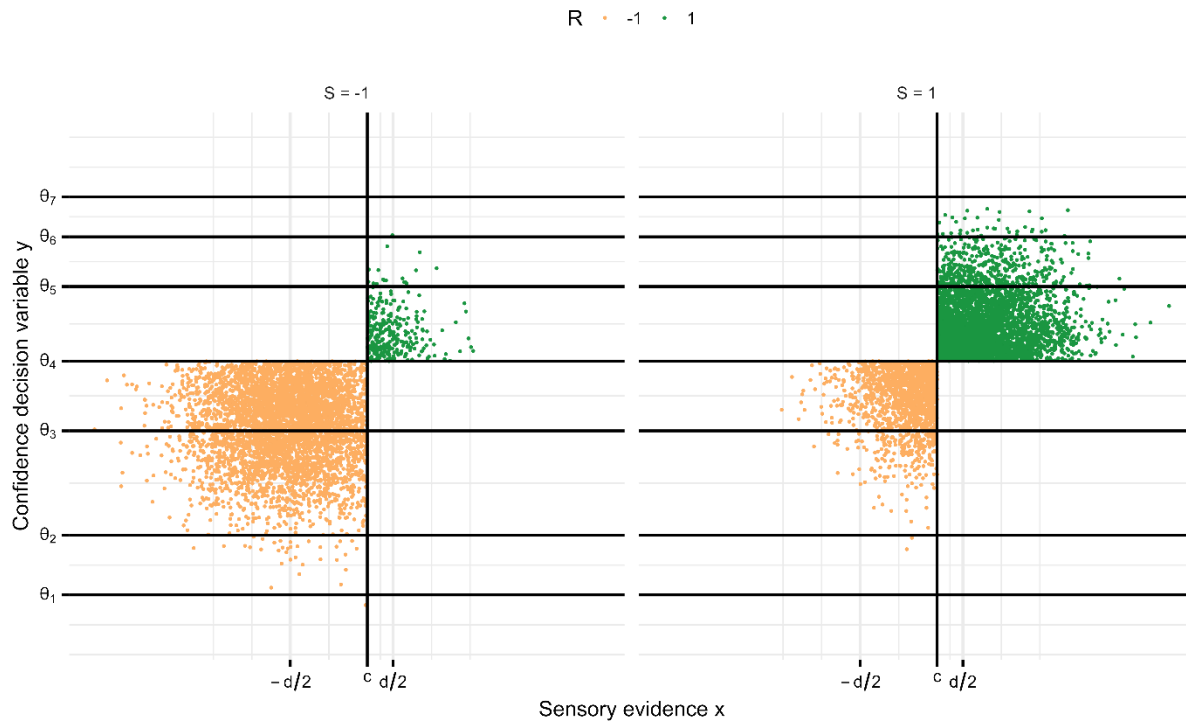
161    **Figure 2**

162    *Bayesian graphical model of the independent truncated Gaussian model (ITG)*

$$x_j \sim N\big(S_j \times 0.5 \times d\,, 1\big)$$

$$R_j \leftarrow \begin{cases} 1 \ if \ x_j > c \\ -1 \ if \ x_j \leq c \end{cases}$$

$$y_j \sim \begin{cases} trN^{m \times c}\big(S_j \times 0.5 \times d \times m, 1\big) \\ \quad if \ R_j = -1 \\ trN_{m \times c}\big(S_j \times 0.5 \times d \times m, 1\big) \\ \quad if \ R_j = 1 \end{cases}$$

$$C_j \leftarrow \begin{cases} 1 \ if \ R = 1 \ and \ y_j < \theta_3 \\ 1 \ if \ R = -1 \ and \ y_j > \theta_1 \\ 2 \ if \ R = 1 \ and \ y_j > \theta_3 \\ 2 \ if \ R = -1 \ and \ y_j < \theta_1 \end{cases}$$

163

164    *Note.* Version of ITG to reproduce the probabilities of confidence categories given stimulus and

165    response underlying the maximum likelihood method devised by Maniscalco and Lau (2014). $S_j$,

166    $R_j$, and $C_j$ are stimulus class, response, and confidence in trial j, respectively , d is the

167    discrimination sensitivity parameter, c is the discrimination criterion, $\theta$ is the confidence

168    criterion, m is the metacognitive efficiency parameter, $x_j$ is the sensory evidence in trial j, and

169    $y_j$ is the confidence decision variable in trial j. $trN_a^b$ indicates a Gaussian distribution which is

170    truncated at the left side and at b at the right side. Following the convention by Lee and

171    Wagenmakers (2013), continuous variables are depicted as circles and discrete variables as

172    squares, observed variables are shaded, unobserved variables not shaded, stochastic dependence

173    is indexed by single borders, and deterministic dependence by double borders.

174    **Figure 3**

175    *Two-dimensional distributions of sensory evidence x and confidence decision variable y*

176    *according to the independent truncated Gaussian model (ITG)*



177

178    *Note.* Fig. 3 is based on a simulation of the ITG model, using Fleming's model specification, and

179    assuming the following parameters: $d = 2$, $c = 0.5$, $m = 0.5$.

180             The implications of the similarity of the meta-d′ method and the ITG model with respect

181    to the interpretation of meta-d′/d′ has to our knowledge not yet been explored: In standard signal

182    detection theory, measures of sensitivity are only guaranteed to be independent from response

183    criteria if the underlying SDT model is a reasonable approximation of the underlying processes

184    (Green & Swets, 1966; Macmillan & Creelman, 2005; Wickens, 2002). Unfortunately, examples

185    of generative models have been presented where meta-d′ is not robust against a variation of

186    discrimination task performance and confidence criteria: According to a model where the

187    confidence criteria are affected by lognormal noise, meta-d′/d′ is influenced by confidence

188    criteria (Shekhar & Rahnev, 2021). According to a Bayesian model where confidence is affected

189    by beta-distributed metacognitive noise, meta-d′/d′ depends on discrimination task performance

190    (Guggenmos, 2021). Thus, the question arises how robust the control that meta-d′/d′ provides

191    over discrimination task performance, discrimination task criterion, and confidence criteria is if

192    the space of different generative models underlying confidence is varied more widely.

193    **Rationale of the present study**

194        In the present study, we investigated whether meta-d′/d′ is influenced by discrimination

195    task performance, discrimination task criterion, and confidence criteria. For this purpose, we

196    simulated artificial data while systematically varying the underlying generative model of

197    confidence. Because the number of generative models of confidence proposed in the literature is

198    far greater than what can be investigated in a single study (e.g. Desender et al., 2021; Fleming &

199    Daw, 2017; Guggenmos, 2022; Mamassian & de Gardelle, 2021; Rausch et al., 2018;

200    Maniscalco & Lau, 2016; Shekhar & Rahnev, 2021; Reynolds et al., 2020; Hellmann et al.,

201    2023; Boundy-Singer et al., 2022; Zhu et al., 2023; Moran et al., 2015; Pereira et al., 2021), for

202    the purpose of the present study, we restricted our analysis to models where the discrimination

203    task decision is made consistent with signal detection theory and thus applying a meta-d′/d′

204    model is considered appropriate (Fleming & Lau, 2014). Besides two versions of the

205    independent truncated gaussian model, one equivalent to the hypothetical SDT models used by

206    Maniscalco and Lau (2014) and one equivalent to the hypothetical SDT models used by Fleming

207    (2017), we used five different models reflecting different cognitive mechanisms how confidence

208    judgments may be generated (see Table 2). For each simulation, we computed meta-d′/d′ using

209    three different methods: 1) the conditioned maximum likelihood method proposed by

210    Maniscalco and Lau (2012, 2014), 2) the Bayesian MCMC method described by Fleming (2017),

211    and 3) conditioned maximum likelihood estimation using Fleming's specification of the

212    hypothetical SDT model.

**Table 2**

*List of cognitive models in which we analyzed the behavior of meta-d'/d'*

| Model | Reference | Conceptual interpretation of the model |
|---|---|---|
| Independent truncated Gaussian model | Maniscalco and Lau (2014) Fleming (2017) | Information used for confidence is generated independently from the evidence used for the choice. Evidence contradicting the original choice cannot be collected. |
| Postdecisional accumulation model | Pleskac and Busemeyer (2010) | After the choice, accumulation of sensory evidence continues for a fixed time interval |
| Gaussian noise model | Maniscalco and Lau (2016) | Confidence is informed by the same sensory evidence as the task decision, but confidence is affected by additive Gaussian noise. |
| Response-congruent evidence model | Maniscalco et al. (2016) Peters et al. (2017) | Confidence is informed only by evidence supporting the selected decision option; evidence in favor of the other option is ignored |
| Confidence boost model | Mamassian and de Gardelle (2021) | Confidence is informed by the evidence used for the choice and by evidence collected in parallel to the choice. In addition, confidence is affected by additive Gaussian noise. |
| Weighted evidence and visibility model | Rausch et al. (2018, 2020, 2021) | Confidence is informed by the evidence used for the choice as well as by a parallel estimate of the difficulty of the task. In addition, confidence is affected by additive Gaussian noise. |

213            We expected that meta-d/d′ is independent from discrimination task performance,

214    discrimination task criteria, and confidence criteria when the generative model is the independent

215    truncated Gaussian model. At least for some of the alternative models, we expected that meta-

216    d′/d′ depends on discrimination task performance, discrimination task criterion, and confidence

217    criteria.

218                                    **Simulation 1**

219    **Method**

220    *Model specification*

221         We simulated data using seven different generative models:

222         i.    the independent truncated Gaussian model with the Gaussian distribution

223               truncated at the discrimination task criterion multiplied with metacognitive

224               efficiency (consistent with the hypothetical SDT model proposed by Maniscalco

225               and Lau, 2014),

226         ii.   the independent truncated Gaussian model with the Gaussian distributions

227               truncated at the discrimination task criterion (consistent with the hypothetical

228               SDT model used by Fleming (2017),

229         iii.   the Gaussian noise model,

230         iv.   the postdecisional accumulation model,

231         v.    the weighted evidence and visibility model,

232         vi.   the confidence boost model, and

233         vii.   the response-congruent evidence model.

234         For all seven models, we assumed that participants select a discrimination response R $\in$

235    {-1, 1} about the stimulus class S $\in$ {-1, 1} as well as a confidence judgment on a five-point

236    scale that the response about the stimulus is correct C $\in$ {1, 2, 3, 3, 5}. According to all seven

237    models, a decision about the stimulus is made by comparing the sensory evidence x against the

238     decision criterion c. Participants respond R = -1 if x < c and R = 1 if x > c. The sensory evidence

239     x is modelled as a random sample from a Gaussian distribution:

240     $$x \sim N(\mu = S \times 0.5 \times d, \sigma = 1)$$

241         The more sensitive the observer is to the stimulus, the greater is the distance d between

242     the centres of the distributions created by the two stimuli. Thus, d is interpreted as the ability of

243     the observer's perceptual system to differentiate between the two kinds of S. The different

244     models are characterised by ways how the confidence decision variable y is generated. A specific

245     degree of confidence is determined by comparing y against a set of confidence criteria. To be

246     consistent with standard SDT, we assumed separate of confidence criteria for each of the two

247     response options. For all models, we assumed for simplicity that confidence criteria are placed

248     symmetrically around the central confidence criterion $\theta_5$ with the placement of criteria

249     determined by the parameter $\tau$. For the version of ITG modelled after Maniscalco and Lau's

250     method, $\theta_5$ was set to $c \times m$. For the version of ITG modelled after Fleming's method, as well as

251     for the five alternative models of confidence, $\theta_5$ was set to c. For R = -1, the other confidence

252     criteria are located at $\theta_1 = \theta_5 - 2 \times \tau$, $\theta_2 = \theta_5 - 1.5 \times \tau$, $\theta_3 = \theta_5 - \tau$, and $\theta_4 = \theta_5 -$

253     $0.5 \times \tau$. For R = 1, the confidence criteria are located at $\theta_6 = \theta_5 + 0.5 \times \tau$, $\theta_7 = \theta_5 + \tau$, $\theta_8 =$

254     $\theta_5 + 1.5 \times \tau$, and $\theta_9 = \theta_5 + 2 \times \tau$. Each criterion delineates between two adjacent confidence

255     criteria, e.g., the observer reports confidence C = 2 if the response R is -1 and y fell between $\theta_1$

256     and $\theta_2$, or if R = 1 and y fell between $\theta_6$ and $\theta_7$. Thus, $\tau$ represents how liberally or

257     conservatively participants place their confidence criteria.

258         **Gaussian noise model.** Conceptually, the Gaussian noise model reflects the idea that

259     confidence is informed by the same sensory evidence as the task decision, but confidence is

260     affected by additive Gaussian noise. Therefore, the confidence decision variable y is also

261    sampled from a Gaussian distribution, with a mean equal to the sensory evidence x and a

262    standard deviation $\sigma_c$, an additional free parameter.

263
$$y \sim N(\mu = x, \sigma = \sigma_c)$$

264    **Postdecisional accumulation model.** The postdecisional accumulation model was

265    inspired by two-stage signal detection theory, according to which accumulation of sensory

266    evidence is continued after the decision for a fixed time interval (Pleskac & Busemeyer, 2010).

267    To ensure comparability with the other models, we used a model that represents the conceptual

268    idea of ongoing accumulation of evidence but does not model reaction time data as well.

269    According to PDA, the confidence decision variable y is sampled from a Gaussian distribution:

270
$$y \sim N(\mu = x + S \times 0.5 \times d \times b, \sigma = \sqrt{b})$$

271    The free parameter b indicates the amount of postdecisional accumulation relative to the

272    amount of evidence available at the time of the discrimination decision. The standard deviation

273    equals the square root of b because both the mean and the variance of the decision variable

274    increase linearly with time in drift diffusion processes (Pleskac & Busemeyer, 2010).

275    **Weighted evidence and visibility model.** The conceptual idea underlying the weighted

276    evidence and visibility model is that the observer combines evidence about the choice-relevant

277    feature of the stimulus with the strength of evidence about choice-irrelevant features to select one

278    out of several confidence categories (Rausch et al., 2018, 2020, 2021). Evidence about choice-

279    irrelevant features of the stimulus can improve confidence judgement because they allow the

280    observer to estimate the reliability of the percept more precisely (Rausch & Zehetleitner, 2019).

281    To express this idea in formal terms, the WEV model assumes that y is sampled from a Gaussian

282    distribution with the standard deviation $\sigma_c$:

283
$$y \sim N(\mu = (1 - w) \times x + w \times d \times R, \sigma = \sigma_c)$$

284       The standard deviation $\sigma_c$ quantifies the amount of unsystematic variability contributing

285   to confidence judgments but not to identification judgments. The unsystematic variability may

286   stem from different sources, including the uncertainty in the estimate of stimulus strength or the

287   noise inherent to metacognitive processes. The factor R ensures that strong stimuli tend to shift

288   the location of the distribution in a way that high confidence is more likely, and likewise, weak

289   stimuli tend to shift the location of the distribution in a way that the probability of low

290   confidence increases.

291       **Confidence boost model.** The confidence boost model represents the idea that the

292   confidence decision variable y is only partially based on the information used during the

293   perceptual decision (Mamassian & de Gardelle, 2021). The confidence boost reflects information

294   used for confidence judgments which was not used for perceptual decision. For this purpose, the

295   model includes the parameter α, which quantifies the degree to which observer base their

296   confidence judgments on information available for the perceptual decision. If α = 0, confidence

297   judgments are exclusively based on information already used for the perceptual decisions; if α =

298   1, the observer has direct access to the original stimulus, and not just the noisy sensory evidence

299   used to make the perceptual decision. In addition, there is again confidence noise superimposed

300   on the confidence decision variable $\sigma_c$. Because Mamassian and de Gardelle (2021) conceived

301   their model for confidence forced choice paradigms, the model was slightly adapted to be

302   applicable for tasks where meta-d′/d′ is typically used. In the version of the model used in the

303   present study, y is sampled from a Gaussian distribution with the standard deviation $\sigma_c$:

304 $$y \sim N(\mu = 0.5 \times S \times d + x \times (1 - \alpha), \sigma = \sigma_c)$$

305       **Response-congruent evidence model.** The model was inspired by the confidence model

306   proposed by Peters et al. (2017). Conceptually, the model represents the idea that observers use

307    all available sensory information to make the primary task decision, but for confidence

308    judgments, they only consider evidence consistent with the selected decision and ignore evidence

309    against the decision (Maniscalco et al., 2016; Odegaard, Grimaldi, et al., 2018; Samaha et al.,

310    2016; Zylberberg et al., 2012). In our version of the model, the response-congruent evidence

311    model assumes two separate samples of sensory evidence collected in each trial, each belonging

312    to one possible identity of the stimulus:

313
$$x_1 \sim N(\mu = (1 - S) \times 0.25 \times d, \sigma = \sqrt{1/2})$$

314
$$x_2 \sim N(\mu = (1 + S) \times 0.25 \times d, \sigma = \sqrt{1/2})$$

315         The sensory evidence used for the discrimination choice is $x = x_2 - x_1$, which implies

316    that the discrimination decision is equivalent to standard signal detection theory. The confidence

317    decision variable depends on the response selected by the observer:

318
$$y = \begin{cases} -x_1, if \ R = -1 \\ \ \ x_2, if \ R = 1 \end{cases}$$

319    ***Simulations***

320         Table 3 lists all parameters we used for our simulations. The parameters were chosen to

321    investigate the behaviour of meta-d′/d′ across a decent range while at the same time avoiding

322    extreme frequencies of events, which are known to lead to unstable behaviour (Barrett et al.,

323    2013). For each generative model, we performed one simulation for each possible combination

324    of parameters. In each simulation, we randomly simulated 10^6 discrimination responses and

325    confidence ratings for both stimuli. Then, we computed meta-d′/d′ using three different methods:

326         i.    the conditioned maximum-likelihood method as described by Maniscalco and Lau

327               (2014),

328         ii.   the Bayesian MCMC method used by Fleming (2017),

329          iii.     a conditioned maximum-likelihood method that uses the specification of the

330                   hypothetical SDT model used by Fleming (2017).

331          A simulation was only included into the results if the estimated standard error of meta-d′

332   was below .005. All analyses were conducted using R (R Core Team, 2020).

**Table 3**

*Parameters for each generative model of confidence*

| Model | Para-meter | values used during simulations | Interpretation of the parameter |
|---|---|---|---|
| All models | d | 0.5, 1.0, 1.5, 2.0, 2.5 | sensitivity of the observer to discriminate between the two stimulus classes |
| | c | 0, 0.25, 0.5, 1, 1.5, 2 | criterion for the primary task response |
| | τ | 0.5, 1.0, 1.5, 2.0, 2.5 | placement of confidence criteria |
| Independent truncated Gaussian model | m | 0.5, 1, 1.5 | Amount of signal available for metacognition relative to the signal available for the discrimination choice |
| Gaussian noise model | $\sigma_c$ | 0.5, 1, 2 | amount of noise superimposed on rating response |
| Postdecisional accumulation model | b | 0.1, 0.5, 1 | amount of postdecisional accumulation relative to the evidence available at the time of the discrimination decision |
| Weighted evidence and visibility model | $\sigma_c$ | 0.5, 2 | amount of Gaussian noise superimposed on rating response |
| | w | 0.25, 0.75 | degree to which confidence relies on sensory evidence about the identity or on strength of evidence about identification-irrelevant features of the stimulus |
| Confidence boost model | $\sigma_c$ | 0.5, 2 | amount of normal noise superimposed on rating response |
| | α | 0.25, 0.75 | degree to which observer has direct access to the original stimulus when making the confidence judgment |

333

334        **Conditioned maximum likelihood estimation of Maniscalco and Lau's model.** To

335    estimate meta-d′ based on conditioned maximum likelihood estimation, we used a translation of

336    the MATLAB code provided by Brian Maniscalco

337    (http://www.columbia.edu/~bsm2105/type2sdt, last accessed 2021-09-20) to R. The algorithm

338    involved the following computational steps: First, the frequency of each confidence category was

339    determined depending on the stimulus class and the accuracy of the response. To correct for

340    extreme proportions, $1/(2n)$ was added to each cell of the frequency table. Second,

341    discrimination sensitivity $d′$ and discrimination criterion c were calculated using standard

342    formulae

$$d' = \Phi^{-1}(\frac{n_{S1R1}}{n_{S1}}) - \Phi^{-1}(\frac{n_{S0R1}}{n_{S0}}) \tag{3}$$

$$c = -\frac{1}{2} \times \left( \Phi^{-1}(\frac{n_{S1R1}}{n_{S1}}) + \Phi^{-1}(\frac{n_{S0R1}}{n_{S0}}) \right) \tag{4}$$

343    with $n_{S1}$ the number of trials when $S = 1$, $n_{S0}$ the number of trials when $S = -1$, $n_{S1R1}$ the

344    number of trials when $S = 1$ and $R = 1$, $n_{S0}$ the number of trials when $S = -1$, $n_{S0R1}$ the number

345    of trials when $S = -1$ and $R = 1$, and $\Phi^{-1}$ the quantile function of the standard Gaussian

346    distribution. The third step involved fitting the meta-d′ model. For this purpose, a maximum

347    likelihood optimization procedure was used with respect to the probability of confidence given

348    stimulus and response as well as the parameters determined at previous steps, i.e., $d′$ and $c$. Model

349    fitting involved a free parameter for meta-d′ $d_{meta}$ as well as the rating criteria $\theta_1, \theta_2, …, \theta_{n-1}$,

350    $\theta_{n+1}, \theta_{n+2}, …, \theta_{2n-1}$. To reproduce the original method by Maniscalco and Lau, $\theta_n$ was fixed at

351    $c \times d_{meta} \div d′$ . To enforce that the criteria were ordered, all free criteria were parametrized as

352    the log of the distance to the adjacent criterion. Model fitting was performed in two steps: First, a

353    coarse grid search was used to identify promising starting values. Second, the five best parameter

354    sets were used as initial values for an Nelder-Mead optimization algorithm as implemented in the

355    R function optim (Nelder & Mead, 1965). We restarted the optimization four times, using the

356    previously found result as initial value for the next iteration to prevent the algorithm from getting

357    stuck in a local minimum. Standard errors associated with the estimate of meta-d′ were obtained

358    by inverting the Hessian matrix returned from optim.

359        **Conditioned maximum likelihood estimation of Fleming's model.** To fit meta-d′/d′

360    using conditioned maximum likelihood estimation and a model specification equivalent to the

361    method used by Fleming (2017), we used the same algorithm as for Maniscalco and Lau's model

362    specification with the exception that $\theta_n$ was fixed at c.

363        **Bayesian Markov Chain Monte Carlo.** To estimate meta-d′/d′ using Bayesian MCMC,

364    we used R code provided by Steve Fleming (https://github.com/metacoglab/HMeta-d, last

365    accessed 2022-10-22), which relies on the free software jags to sample from the posterior

366    distribution (Plummer, 2003). For more details on the underlying Bayesian estimation procedure,

367    see Fleming (2017). Just as for standard meta-d′, discrimination performance d′ and

368    discrimination criterion $c$ were computed first using formulae (3) and (4) and then submitted to

369    jags as constants. The Bayesian estimation procedure was used only for the meta-d′/d′ and

370    confidence criteria. For this purpose, the absolute frequency of each confidence rating given

371    stimulus and response $f(C|S,R)$ was modelled as a multinomial distribution $\mathcal{M}$,

$$f(C|S,R) \sim \mathcal{M}(n = n_{SR}, p = p(C|S,R)) \tag{5}$$

372    where $n_{SR}$ is the number of trials with stimulus S and response R, and $p(C|S,R)$ calculated using

373    formulae (1) and (2). $\theta_n$ was fixed at $c$. $p(C|S,R)$ depends on the free parameters $d_{meta}$ and a

374    set of criteria $\theta$. The priors for the parameters were specified as follows:

$$\theta_{1,2,\ldots,n-1} \sim trN\left(\mu = 0, \quad \sigma = \sqrt{0.5}, \quad a = -\infty, \quad b = c\right) \tag{6}$$

$$\theta_{n+1,n+2,\dots,2\times n-1} \sim tr\mathcal{N}\left(\mu = 0, \quad \sigma = \sqrt{0.5}, \quad a = c, \quad b = \infty\right)$$

$$d_{meta} \sim \mathcal{N}\left(\mu = d', \sigma = \sqrt{2}\right)$$

375   where $\theta_{1,2,\dots,n-1}$ indicates the set of confidence criteria when the response was -1,

376   $\theta_{n+1,n+2,\dots,2\times n-1}$ indicates the set of confidence criteria when the response was 1, $tr\mathcal{N}$ indicates

377   a truncated gaussian distribution with a location parameter $\mu$, scale parameter $\sigma$, lower bound a,

378   and upper bound b, and $d_{meta}$ is meta-d′. These priors reflect the standard settings. Sampling

379   was performed in three separate Markov Chains to allow computation of Gelman and Rubin's

380   convergence diagnostic $\hat{R}$ (Gelman & Rubin, 1992). For each chain, we drew 100,000 samples

381   from the posterior distribution, saving every $10^{th}$ sample to remove autocorrelations in the

382   Markov chain. If $\hat{R}$ was larger than 1.1, the simulation was excluded from the analysis.

383           **Transparency and openness.** All data and analysis code are available at

384   https://osf.io/72uds. This study's design and its analysis were not pre-registered.
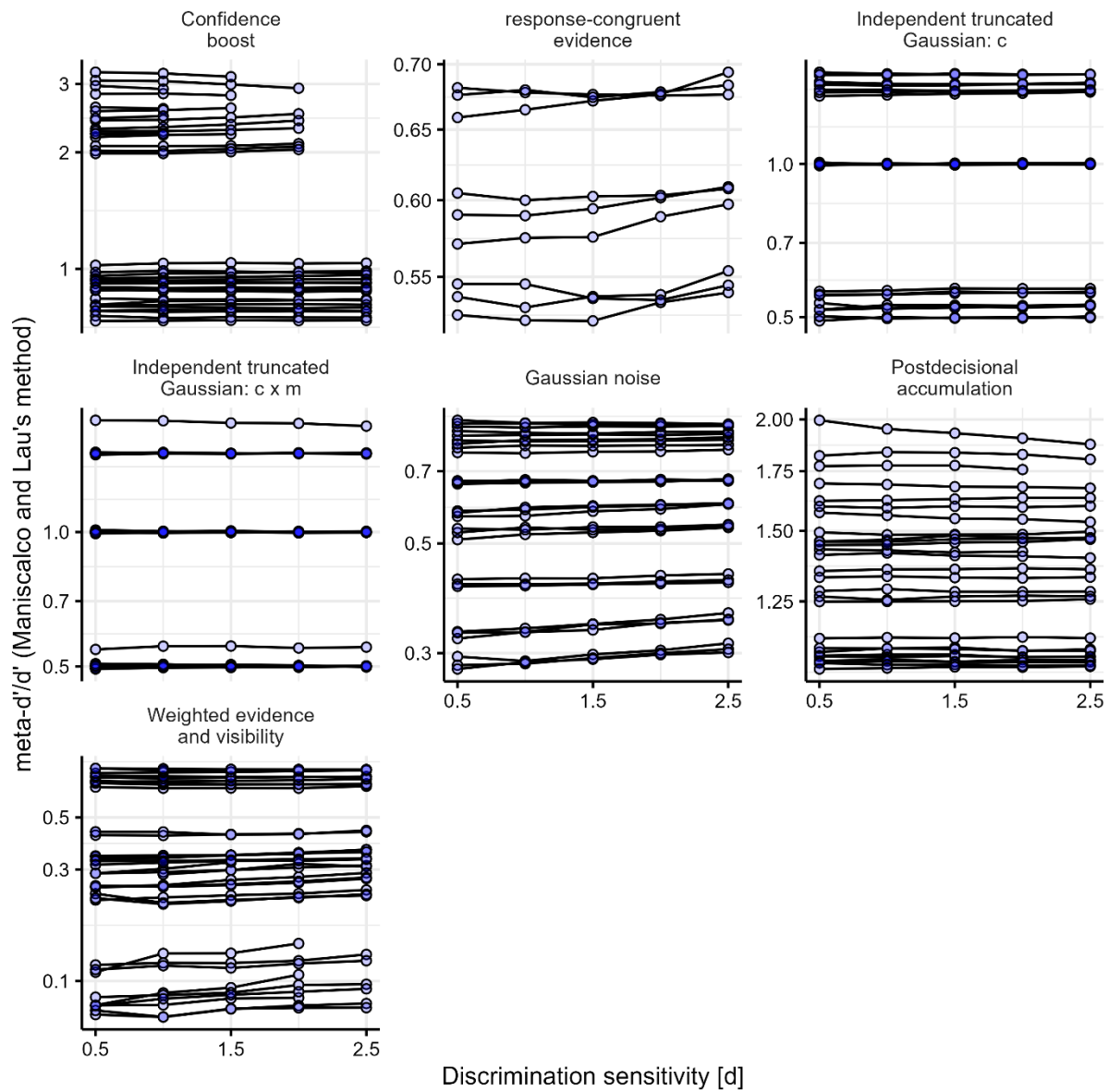
385   **Results**

386   *Discrimination sensitivity*

387           Fig. 4 shows the pattern of meta-d′/d′ as estimated using the conditioned maximum

388   likelihood method proposed by Maniscalco and Lau (2012) as a function of the generative model

389   underlying the simulated data and discrimination sensitivity. Meta-d′/d′ was not perfectly

390   constant across different levels of discrimination sensitivity in any of the seven generative

391   models. For the two independent truncated Gaussian models, meta-d′/d′ was associated with

392   discrimination sensitivity only for a relatively small subset of simulations. In contrast, for the

393   postdecisional accumulation model, the Gaussian noise model, the response-congruent evidence

394   model, and the weighted evidence and visibility model, Fig. 4 shows multiple lines that have a

395    non-zero slope, meaning that meta-d′/d′ depended on discrimination sensitivity for the majority

396    of parameter sets.

397    **Figure 4**

398            *Meta-d'-d' based on conditioned maximum likelihood estimation and model specification*

399    *by Maniscalco and Lau, as function of discrimination sensitivity and generative model of*
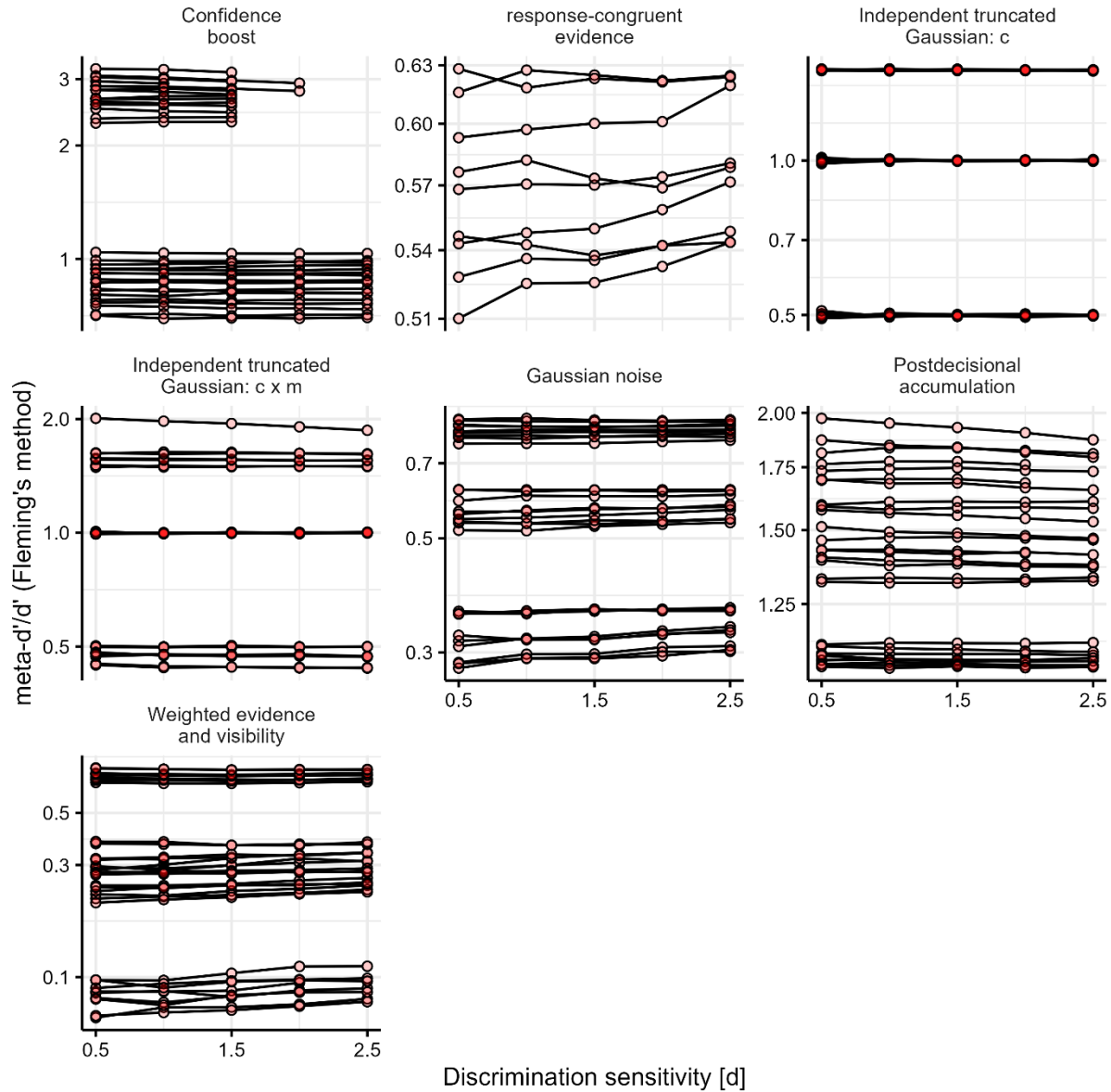
400    *confidence*



401

402     *Note.* Each dot represents one simulation with one combination of parameters. Lines connect

403     simulations that differ only with respect to the parameter quantifying discrimination sensitivity

404     and identical parameter sets otherwise. Lines parallel to the horizontal indicate that meta-d′/d′ is

405     independent from discrimination sensitivity. Note that the y-Axes are different for each

406     generative model of confidence.

407             Fig. 5 shows the pattern of meta-d′/d′ estimated using Fleming's Bayesian MCMC

408     method, again as a function of the generative model underlying the simulated data and

409     discrimination sensitivity. Meta-d′/d′ was constant across levels of discrimination performance

410     when the data was generated according to the independent truncated Gaussian model with

411     distributions truncated at the discrimination criterion c. When the same model was used but with

412     distributions truncated at c × m, there were some parameter sets where discrimination sensitivity

413     affected meta-d′/d′. Again, for the postdecisional accumulation model, the Gaussian noise model,

414     the response-congruent evidence model, and the weighted evidence and visibility model,

415     discrimination sensitivity affected meta-d′/d′ ratios for a large number of parameter sets. When

416     we repeated these analyses using conditioned maximum likelihood estimation but calculating the

417     probability of confidence given stimulus and response following Fleming (2017), the results

418     were visually indistinguishable from Fig. 5.

419     **Figure 5**

420             *Meta-d'/d' based on Bayesian MCMC estimation and Fleming's model specification, as*

421     *function of discrimination sensitivity and generative model of confidence*

422

*Note.* Each dot represents one simulation. Lines connect simulations that differ only with respect

to the parameter quantifying discrimination sensitivity and identical parameter sets otherwise.

Lines parallel to the horizontal indicate that meta-d′/d′ is independent from discrimination

sensitivity. Note that the y-Axes are different for each generative model of confidence.
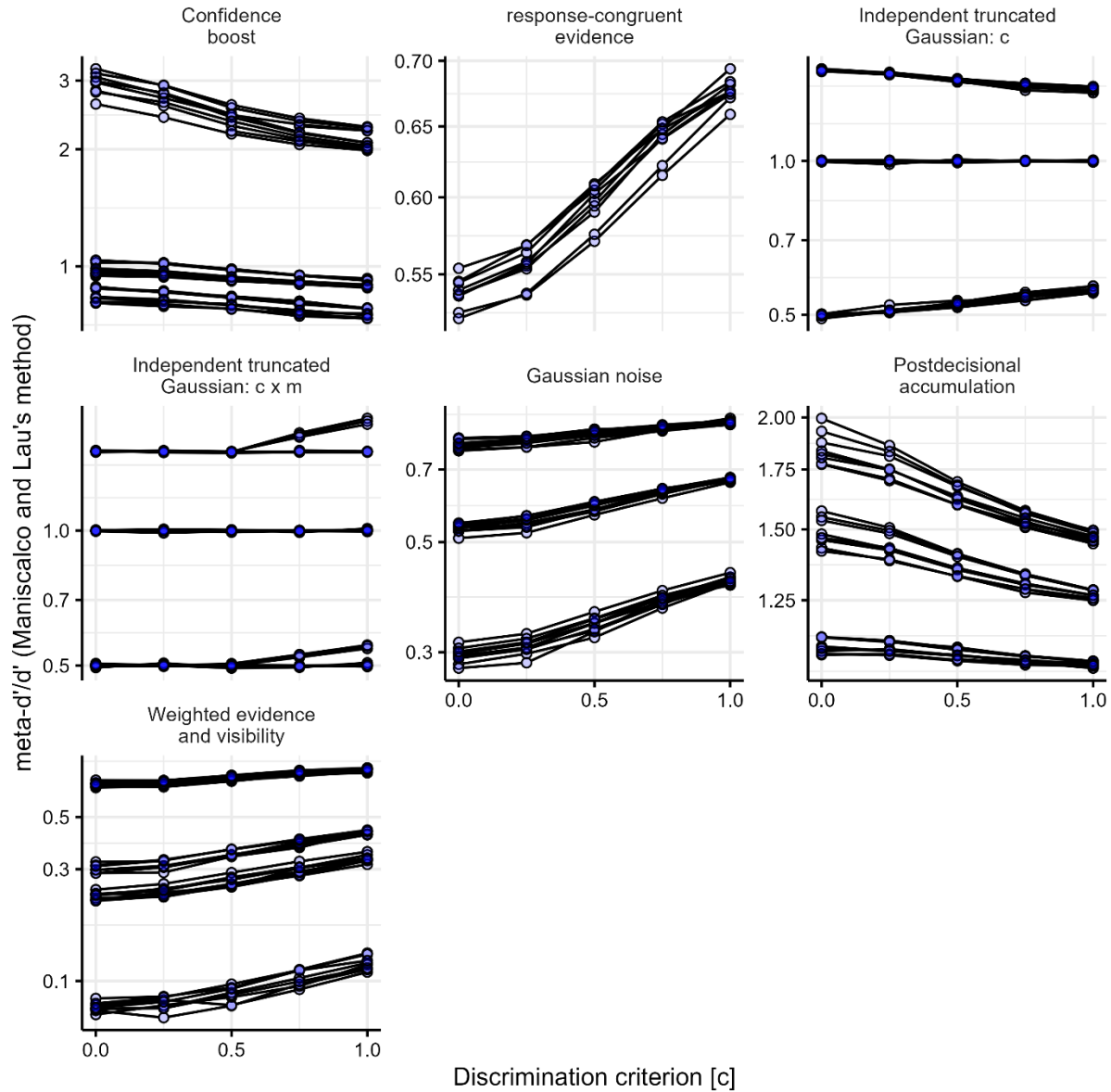
### *Discrimination bias*

The relationship between meta-d′/d′ and discrimination bias across different generative

models is depicted in Fig. 6 for Maniscalco and Lau's original conditioned maximum likelihood

430    method and in Fig. 7 for Fleming's Bayesian MCMC method. Fig. 6 shows that meta-d′/d′

431    estimated using the original method depends on discrimination bias for each single generative

432    model of confidence. Fig. 7 shows that meta-d′/d′ estimated using the Bayesian MCMC method

433    is independent from discrimination bias only if the data is generated according to the

434    independent truncated Gaussian model with the distributions truncated at the discrimination

435    criterion. Again, meta-d′/d′ depends on the discrimination criterion according to all other

436    generative models of confidence. Finally, when meta-d′/d′ was estimated using conditioned

437    maximum likelihood estimation but using the model specification Fleming (2017), the results

438    were the same as in Fig. 6.

439    **Figure 6**

440           *Meta-d′/d′ based on conditioned maximum likelihood estimation and Maniscalco and*

441    *Lau's model specification as function of discrimination bias and generative model of confidence*
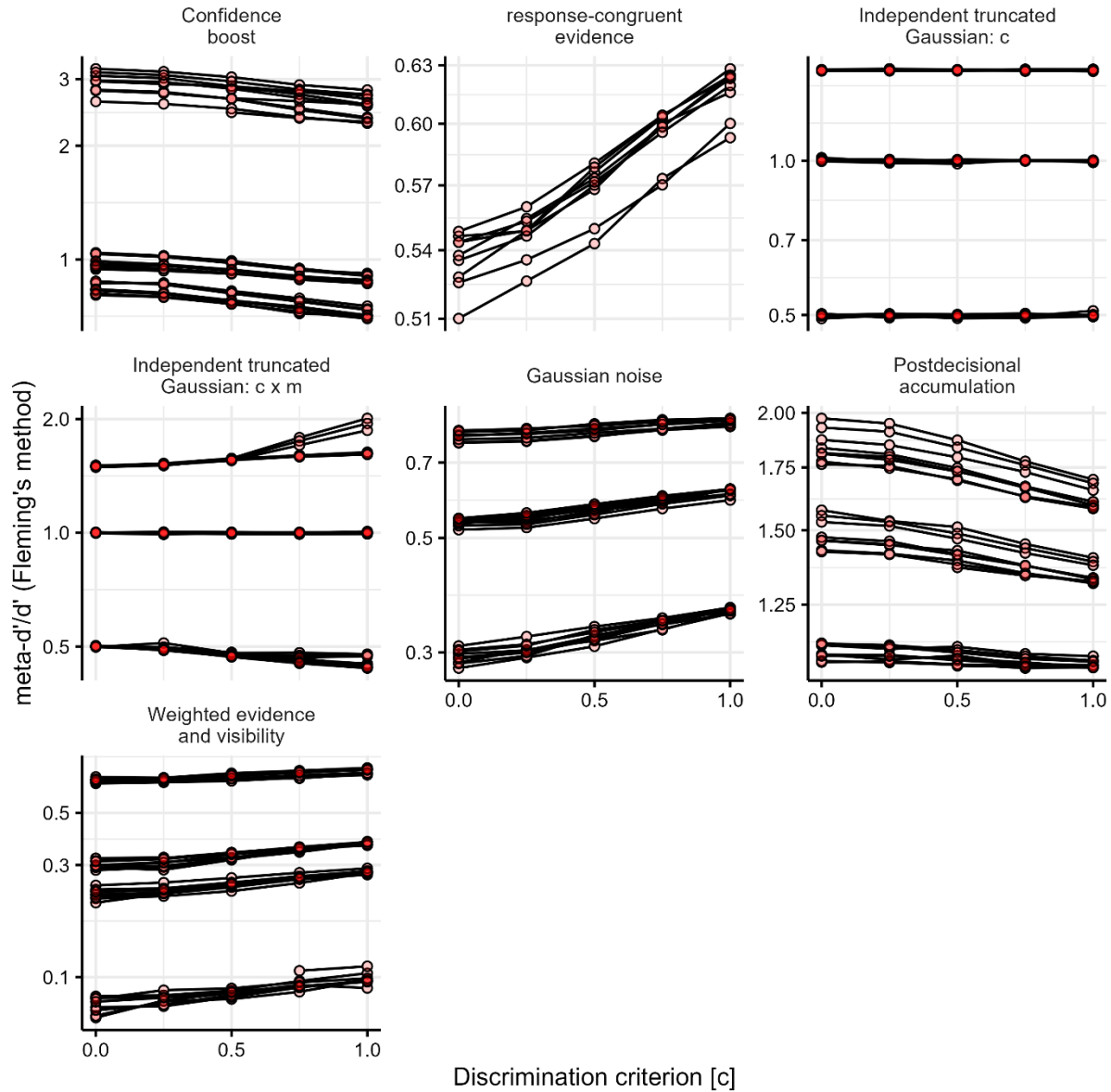
442

443     *Note.* Each dot represents one simulation. Lines connect simulations that differ only with respect

444     to the parameter quantifying discrimination bias and identical parameter sets otherwise. Lines

445     parallel to the horizontal indicate that meta-d′/d′ is independent from discrimination bias. Note

446     that the y-Axes are different for each generative model of confidence.

447     **Figure 7**

448           *Meta-d'/d' based on MCMC estimation and Fleming's model specification as function of*

449     *discrimination bias and generative model of confidence*

451    *Note.* Each dot represents one simulation. Lines connect simulations that differ only with respect

452    to the parameter quantifying discrimination bias and identical parameter sets otherwise. Lines

453    parallel to the horizontal indicate that meta-d′/d′ is independent from discrimination bias. Note

454    that the y-axes are different for each generative model of confidence.
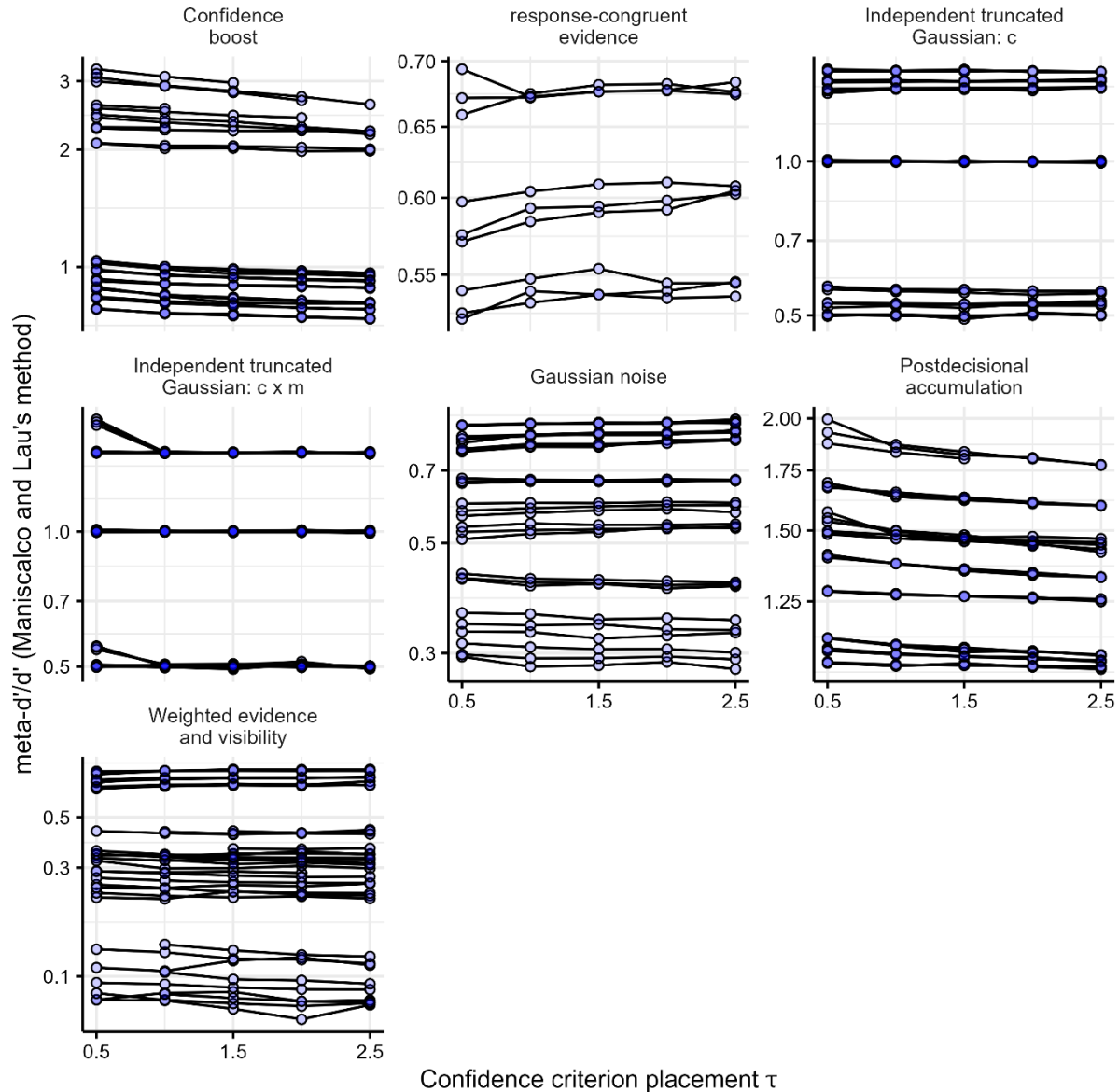
455    ***Confidence criteria***

456           The relationship between meta-d′/d′ and confidence criterion placement across different

457    generative models of confidence is depicted in Fig. 8 for Maniscalco and Lau's original

458    conditioned maximum likelihood method and in Fig. 9 for Fleming's Bayesian MCMC method.

459    Fig. 8 shows that meta-d′/d′ estimated using the original method is never completely independent

460    from confidence criterion placement. Nevertheless, for the two independent truncated Gaussian

461    models, meta-d′/d′ was associated with confidence criterion placement for a relatively small

462    subset of simulated parameter sets. Fig. 9 shows that meta-d′/d′ estimated using Fleming's

463    method is independent from confidence criterion placement only if the data is generated

464    according to the independent truncated Gaussian model with the distributions truncated at the

465    discrimination criterion. For all other generative models of confidence, meta-d′/d′ depends on

466    confidence criterion placement. Finally, when meta-d′/d′ was estimated using conditioned

467    maximum likelihood estimation but with Fleming's model specification, the results were the

468    same as in Fig. 9.

469    **Figure 8**

470         *Meta-d′/d′ based conditioned maximum likelihood estimation and Maniscalco and Lau's*

471    *model specification as function of confidence criterion placement and generative model of*

472    *confidence*

473

*Note.* Each dot represents one simulation. Lines connect simulations that differ only with respect

to the parameter determining confidence criterion placement and identical parameter sets

otherwise. Lines parallel to the horizontal indicate that meta-d′/d′ is independent from confidence

criterion placement. Note that the y-Axes are different for each generative model of confidence.

**Figure 9**

*Meta-d′/d′ based on MCMC estimation and Fleming's model specification as function of*

*confidence criterion placement and generative model of confidence*

481

*Note.* Each dot represents one simulation. Lines connect simulations that differ only with respect

to the parameter determining confidence criterion placement and identical parameter sets

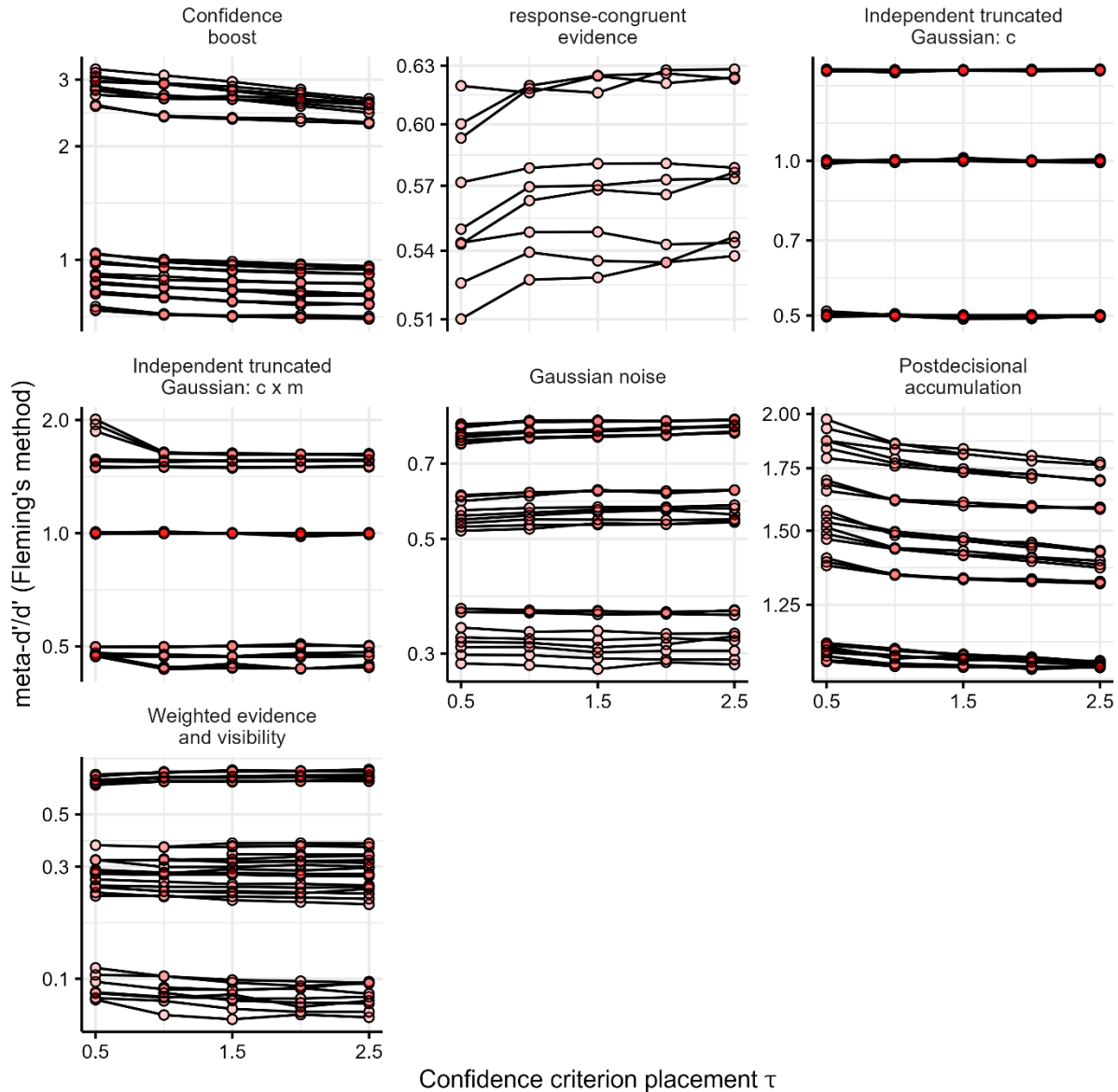otherwise. Lines parallel to the horizontal indicate that meta-d′/d′ is independent from confidence

criterion placement. Note that the y-Axes are different for each generative model of confidence.
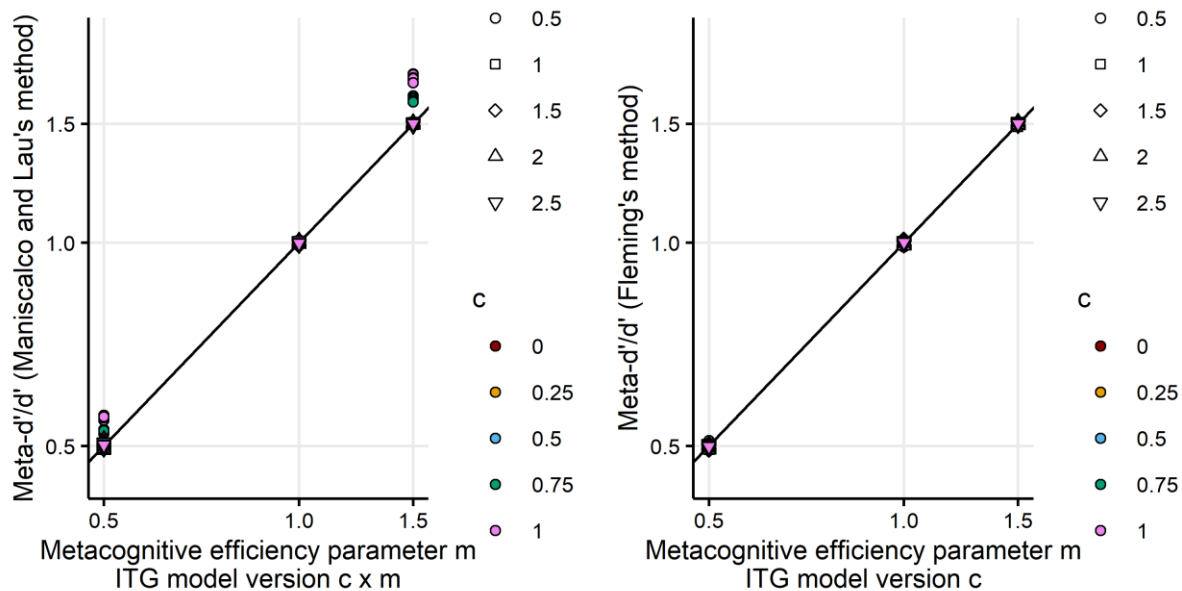
### Recovering metacognitive efficiency parameters

  Finally, we investigated if estimates of meta-d′/d′ recover the metacognitive efficiency

parameter m of the independent truncated Gaussian model. Specifically, meta-d′/d′ estimated

489     using the original SDT model specification by Maniscalco and Lau (2014) was expected to

490     recover the m parameter in the ITG model with the distribution truncated at the objective

491     discrimination criterion c multiplied with m. Meta-d′/d′ estimated using the model specification

492     by Fleming (2017) should recover the m parameter in the ITG model with the distribution

493     truncated at c. Fig. 10 shows that meta-d′/d′ based on Bayesian MCMC estimation and Fleming's

494     model specification indeed recovered the m parameter of the corresponding version of the ITG

495     model. However, meta-d′/d′ using the model specification by Maniscalco and Lau (2014) did not

496     always recover m in the corresponding ITG model. Specifically, meta-d′/d′ overestimated m

497     when the discrimination criterion was at least .75 (i.e., a considerable bias for one of the two

498     stimuli), when τ was 0.5 (i.e. liberal confidence criterion placement), and when m was either 0.5

499     or 1.5 (and thus metacognitive ability and perceptual ability were not the same).

500     **Figure 10**

501          *Meta-d′/d′ as a function of the metacognitive efficiency parameter m, discrimination bias*

502     *parameter θ, and confidence criterion placement parameter τ.*



503

504     *Note.* Left panel: ITG model with distributions truncated at the discrimination criterion c

505     multiplied with m. Accordingly, meta-d′/d′ values on the y-axis were computed using the

506     original method by Maniscalco and Lau (2014). Right panel: ITG model with distributions

507     truncated at the discrimination criterion c. Accordingly, meta-d′/d′ values on the y-axis were

508     computed using the Bayesian MCMC method by Fleming (2017). Colours indicate different

509     objective discrimination criteria. Symbols indicate different placement of confidence criteria.

510     **Discussion**

511          Simulation 1 showed that meta-d′/d′ provides imperfect control over discrimination

512     performance, discrimination bias, and confidence criteria: Only when the data were simulated

513     according to the independent truncated Gaussian model with the distributions truncated at the

514     discrimination bias, and when meta-d′/d′ was estimated using the model specification used by

515     Fleming (2017), meta-d′/d′ was constant across discrimination performance, discrimination bias,

516     and confidence criteria in all simulations. Notably, the control of discrimination sensitivity, bias,

517     and confidence criteria is sensitive to the finer details of model specification: When we simulated

518     data with distributions truncated at the discrimination bias multiplied by the metacognitive

519     efficiency parameter, the generative model consistent with Maniscalco and Lau's method, meta-

520     d′/d′ based on Fleming's model specification was no longer constant as a function of

521     discrimination performance, discrimination bias, and confidence criteria across all simulated

522     parameter sets. When the data were simulated according to one of the other generative models of

523     confidence, meta-d′/d′ was associated with discrimination bias, discrimination sensitivity and

524     confidence criterion placement for numerous simulations.

525          While Simulation 1 shows that meta-d′/d′ depends in principle on discrimination

526     performance, discrimination bias and confidence criteria according to various different models of

527    confidence, it is still unclear whether the effect is large enough to be relevant in practice. In

528    particular, the contamination of meta-d′/d′ by discrimination sensitivity seemed to be relatively

529    small compared to the contamination by discrimination bias and confidence criteria. However, in

530    order to simulate the expected correlations between model parameters and meta-d′/d′ according

531    to different confidence models, it is necessary to specify the distributions of the model

532    parameters across subjects. Unfortunately, the sample sizes of previous modelling studies have

533    been generally too small sample to reasonably estimate the distribution of model parameters

534    across subjects.

<p align="center">**Simulation 2**</p>

536         To investigate how the relationships observed in Simulation 1 may translate into

537    plausible effect sizes, we fitted all seven models of confidence used in Simulation 1 to the data

538    from Experiment 2 by Rouault et al. (2018), an open data set available from the confidence

539    database (Rahnev et al., 2020). Rouault et al. (2018)'s data were chosen because a large sample

540    is necessary for stable estimates of correlation coefficients (Schönbrodt & Perugini, 2013). We

541    then used the parameter sets obtained by model fitting to simulate new data to estimate the

542    correlation between meta-d′/d′ and discrimination sensitivity, discrimination bias and confidence

543    criteria implied by each generative model of confidence.

**Method**

*Experimental task*

546         Rouault et al.'s data consists 497 subjects who participated in an online dot numerosity

547    discrimination task with 210 trials per subject. In each trial, participants were presented with a

548    fixation cross for 1 s. Two black boxes filled with differing numbers of randomly positioned

549    white dots were then presented for 0.3 s. One box was always half-filled (313 dots out of 625

550    positions), while the other box contained an increased number of dots compared to the first box.

551    The position of the box with the higher number of dots was pseudo-randomised across all trials.

552    To maintain a constant level of performance during the experiment and across participants, a

553    staircase was used to adapt the number of extra dots in the target box. The staircase started with a

554    number of 70 extra dots and was a two-down one-up staircase procedure with equal step-sizes

555    for steps up and down. The step-size was calculated in log-space, changing by ± 0.4 for the first

556    5 trials, ± 0.2 for the next 5 trials and ± 0.1 for the rest of the task. After 0.3 s, the dots

557    disappeared, leaving the black boxes on screen until participants indicated which box had the

558    higher number of dots by keyboard button press. Then, subjects were asked to report their

559    confidence in their response on a 6-point rating scale with verbal descriptions (*certainly wrong*,

560    *probably wrong*, *maybe wrong, maybe correct, probably correct, certainly correct*). A detailed

561    description of the study is provided by Rouault et al. (2018).

562    ***Model fitting***

563         All seven generative models of confidence used in Simulation 1 were fitted to the

564    combined distributions of responses and confidence judgments separately for each single

565    participant. The fitting procedure involved the following computational steps: First, the

566    frequency of each confidence level was calculated for each of the two stimulus options and each

567    of the response option. For each model, the set of parameters was determined that minimized the

568    negative log-likelihood of the data given the model. For this purpose, we used a coarse grid

569    search to identify five promising sets of starting values for the optimization procedure. Then,

570    minimization of the negative log-likelihood was performed using a general SIMPLEX

571    minimization routine (Nelder & Mead, 1965) for each set of starting values. To avoid local

572    minima, the optimization procedure was restarted four times.

573        To assess the relative quality of the candidate models, we calculated the Bayes

574    information criterion (Schwarz, 1978) and the AICc (Burnham & Anderson, 2002), a variant of

575    the Akaike information criterion (Akaike, 1974) using the negative likelihood of each model fit

576    with respect to each single participant and the trial number. For statistical testing, we compared

577    the mean AICc and BIC using standard t-tests with *p*-values adjusted for multiple comparisons

578    using Holm's correction.

579    *Simulation*

580        We simulated one new data set for each of the seven generative models of confidence,

581    using the parameter sets we obtained during model fitting, using the same number of subjects as

582    as in the empirical data and 10.000 trials per subject. Then, we estimated meta-d′/d′ two times for

583    each simulated subject using conditioned maximum likelihood estimation, one time with

584    Maniscalco and Lau's model specification, and one time with Fleming's model specification.

585    Because meta-d′/d′ is not normally distributed (Rausch & Zehetleitner, 2023), we assessed the

586    correlation between each parameter of each generative model and the logarithm of meta-d′/d′.

587    We repeated the analysis using unstandardized linear regression slopes with centred parameters

588    as predictors and log(meta-d′/d′) as criterion. All *p*-values were corrected for multiple

589    comparisons using Holm's correction.

590    **Results**

591    *Formal model comparisons*

592        Formal model comparisons revealed that the best fits to the data were obtained by the two

593    versions of the independent truncated Gaussian model, both in terms of $AIC_c$, and BIC. The

594    difference between the two versions of the independent truncated Gaussian model was

595    negligible, $M_{\Delta AIC} = M_{\Delta BIC} = 0.02$, $t(496) = 1.46$, $p = .290$. The fit of both independent truncated

596    Gaussian models was each significantly better than those of the five alternative models in terms

597    of AIC and BIC, all *p's* < .001, although the mean difference was quite small, $M_{\Delta AIC}$'s and

598    $M_{\Delta BIC}$'s ≥ -0.59, all *p's* < .001.

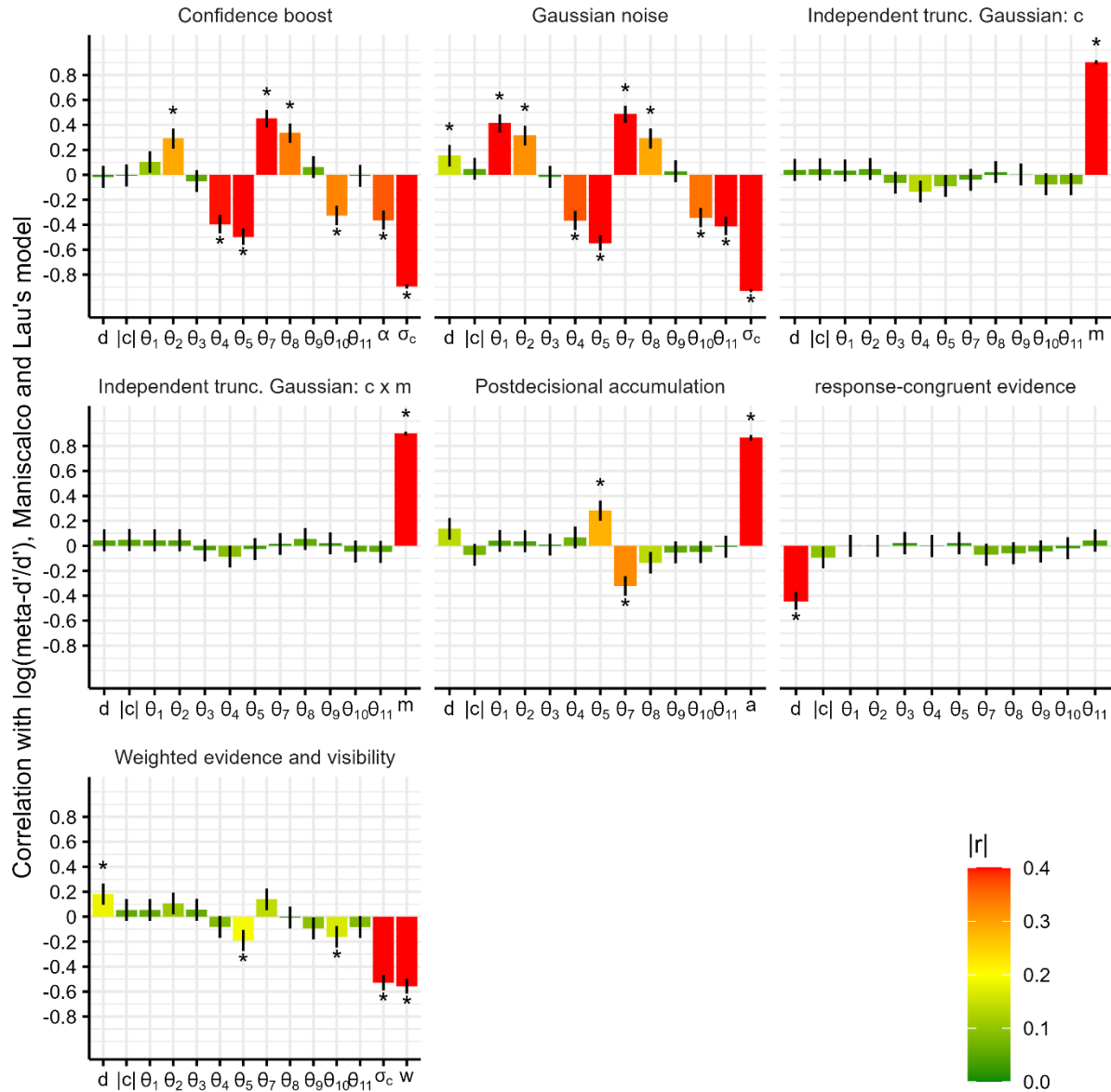### *Correlations between model parameters and simulated meta-d′/d′*

600         Supplementary Table 1 provides the correlation coefficients between each estimated

601    parameter of the different confidence model and log(meta-d′/d′). Figs. 11 and 12 show that as

602    expected, log(meta-d′/d′) is strongly correlated with all model parameters intended to reflect

603    metacognitive efficiency, i.e. $\sigma_c$, m, a, and α. For the two versions of the independent gaussian

604    truncated model, no significant correlation between log(meta-d′/d′) and discrimination sensitivity

605    d, discrimination criterion c, or any of the ten confidence criteria was observed, independently

606    from the specification of the hypothetical SDT model underlying meta-d′/d′.  However, we found

607    a significant large correlation between discrimination sensitivity d and log(meta-d′/d′) for the

608    response-congruent evidence model and a medium-sized correlation between discrimination

609    sensitivity d and log(meta-d′/d′) for the weighted evidence and visibility model. For the Gaussian

610    noise model, a moderate correlation between d and log(meta-d′/d′) was significant only when

611    meta-d′/d′ was estimated using Maniscalco and Lau's model specification, but not for Fleming's

612    model specification. On the contrary, for the postdecisional accumulation model, the correlation

613    between d and log(meta-d′/d′) was significant only when meta-d′/d′ was estimated based on

614    Fleming's model, but not with Maniscalco and Lau's model. A significant medium-sized

615    correlation between log(meta-d′/d′) and discrimination bias c was detected only for the response-

616    congruent evidence model when meta-d′/d′ was estimated using Fleming's model specification.

617    Concerning confidence criteria, we found a very strong correlation between log(meta-d′/d′) and

618    six out of ten confidence criteria for the confidence boost model, seven out of ten for the Gaussian

619     noise model, and two out of ten for the postdecisional accumulation model. In addition, we

620     detected medium-sized correlations between two confidence criteria in the weighted evidence and

621     visibility model.

622     The analysis of regression slopes revealed that for the confidence boost model and the

623     Gaussian noise model, there were only small changes in meta-d′/d′ as a function of confidence

624     criteria, but these changes were very consistent across subjects, resulting in many significant small

625     effects. For the other models and parameters, the interpretation was essentially the same as in the

626     correlation analysis (see Supplementary Table 2).

627     **Figure 11**

628     *Correlation between meta-d'/d' estimated using Maniscalco and Lau's model*

629     *specification and model parameters estimated from Rouault et al. (2018)'s Exp. 2 as a function*

630     *of different generative models of confidence.*

*Note.* Error bars indicate 95% CI.

**Figure 12**

*Correlation between log-transformed meta-d'/d' estimated using Fleming's model*

*specification and model parameters model parameters estimated from Rouault et al. (2018)'s*

*Exp. 2 as a function of different generative models of confidence.*

*Note.* Error bars indicate 95% CI.

**Discussion**

Fitting different models of confidence to Rouault et al. (2018)'s data showed that the two

versions of the independent truncated Gaussian model provide a reasonable fit to confidence in a

dot numerosity discrimination task. Importantly, the model comparisons reported in the present

study should be only interpreted as preliminary, because the data set only included 200 trials per

subject, which is much smaller than the norm in modelling studies. It should also be noted that

645    the statistical properties of different experimental tasks may be quite different, suggesting that

646    the observation that ITG performs well in one data set does not imply that ITG will also perform

647    well in other experimental tasks. Nevertheless, we think that ITG should be considered as a

648    series candidate model in future studies and should be routinely included in future comparisons

649    of confidence models.

650        The simulation using the parameters of the independent truncated Gaussian model

651    obtained during model fitting showed that both versions of meta-d′/d′ were independent of

652    discrimination sensitivity, discrimination bias, and confidence criteria, suggesting that the

653    differences between the two versions of the independent truncated Gaussian model are small

654    enough not to be practically relevant, at least with distributions of parameters as observed in this

655    particular experiment. However, for each of the five alternative models of confidence, we found

656    at least one strong correlation with either discrimination sensitivity or one of the confidence

657    criteria. The correlations with discrimination sensitivity parameters are noteworthy because

658    Rouault et al. used a staircase to keep accuracy constant. This means that staircases still leave

659    enough variance in discrimination sensitivity parameters to produce a large correlation with

660    discrimination sensitivity for the response-congruent evidence model, medium-sized correlations

661    for the weighted evidence and visibility model, or small-to-medium correlations for the gaussian

662    noise model and the postdecisional accumulation model.

663                              **General discussion**

664        The results of the present study suggests that whether or not meta-d′/d′ provides control

665    over discrimination performance, discrimination bias, and confidence criteria strongly depends

666    on the generative model of confidence: Only when the data was simulated according to the

667    independent truncated Gaussian model (ITG) with the distributions truncated at the

668    discrimination bias, and when meta-d′/d′ was estimated using the model specification used by

669    Fleming (2017), meta-d′/d′ was perfectly constant across discrimination performance,

670    discrimination bias, and confidence criteria across all simulations. When we simulated data using

671    the parameters estimated from Rouault et al. (2018)'s Exp. 2, no difference between the two

672    versions of ITG were observed, suggesting that the difference between the two versions of ITG

673    may not always be relevant in practice. However, when the data was simulated not using ITG,

674    but the Gaussian noise model, the postdecisional accumulation model, the weighted evidence and

675    visibility model, the confidence boost model, or the response-congruent evidence model, meta-

676    d′/d′ depended on discrimination sensitivity, discrimination bias, and confidence criterion

677    placement for many simulations. Simulations using parameters obtained by fitting empirical data

678    showed that the expected correlations between meta-d′/d′ and model parameters vary widely

679    across different generative model of confidence and specific parameters. Nevertheless, for each

680    generative model other than ITG, there was at least one medium-sized correlation with either

681    discrimination sensitivity or one of the confidence criteria, suggesting that meta-d′/d′ is

682    associated with discrimination sensitivity and confidence criteria under realistic assumptions

683    about model parameters.

684    **Relation between meta-d′/d′ and generative models of confidence**

685            Meta-d′/d′ has been considered to rely only on the assumption of a specific cognitive

686    architecture underlying the discrimination decision, but to be free from assumptions about the

687    decision variable underlying the confidence decision (Maniscalco & Lau, 2014). In contrast, the

688    main finding of the present study is that meta-d′/d′ is in fact not free from assumptions about the

689    generative model underlying confidence judgments. The reason is that meta-d′/d′ depends on

690    discrimination sensitivity, discrimination bias, and confidence criteria when the data is simulated

691    according to the Gaussian noise model, the weighted evidence and visibility model, the

692    confidence boost model, the postdecisional accumulation model or the response-congruent

693    evidence model. Previous studies revealed two additional models where meta-d′/d′ is

694    confounded, the Bayesian beta-distributed noise model (Guggenmos, 2021) and the lognormal

695    noise model (Shekhar & Rahnev, 2021). Importantly, the present study exceeds those studies in

696    showing that generative models where meta-d′/d′ is contaminated by discrimination sensitivity,

697    discrimination bias and confidence criteria not only exist, but the same result is obtained

698    according to most generative models of confidence. Meta-d′/d′ succeeds in controlling for

699    discrimination sensitivity, discrimination bias and confidence criteria when the data is generated

700    according to the independent truncated gaussian model. Thus, it seems that the control meta-d′/d′

701    provides is highly specific to the independent truncated Gaussian model. Our findings are

702    consistent with the assertion that discrimination sensitivity, discrimination bias and confidence

703    criteria can only be controlled based on estimating the underlying generative model of

704    confidence (Guggenmos, 2022). We cannot prove that no generative model other than ITG exists

705    where meta-d′/d′ performs satisfactorily. However, the control over discrimination sensitivity,

706    discrimination bias, and confidence criteria fails for a large variety of different generative

707    models, which is why it is reasonable to assume that meta-d′/d′ is unlikely to provide effective

708    control in other models which were not examined so far. Overall, this means that meta-d′/d′ from

709    now on should be regarded as a model-based measure of metacognitive efficiency, and

710    researchers who consider using meta-d′/d′ need to ascertain if their data can be adequately

711    described by ITG.

**Evidence for the independent truncated Gaussian model?**

712

713     Because the adequacy of meta-d′/d′ depends on the assumption of ITG as generative

714     model, the question is raised if ITG is a decent models of human confidence judgments. Our

715     analysis of the data of Rouault et al. (2018) is to our knowledge the first (albeit preliminary)

716     evidence that data sets exists which are adequately described by ITG. Unfortunately, previous

717     studies comparing generative models of confidence did not make the link between meta-d′/d′ and

718     generative model of confidence, which is why ITG has not been included into formal model

719     comparisons previously (e.g. Maniscalco & Lau, 2016; Rausch et al., 2018, 2020, 2021; Shekhar

720     & Rahnev, 2021, 2022). Future modelling studies are necessary to investigate how frequently

721     ITG is an adequate description of human confidence. However, there is more evidence for some

722     qualitative predictions of ITG. According to ITG, confidence judgments are subject to a

723     response-congruent confirmation bias because it is impossible to sample a confidence decision

724     variable that contradicts the discrimination decision. In accordance with ITG, previous studies

725     reported that observers' tend to neglect contradictory evidence when they report confidence

726     (Peters et al., 2017; Samaha et al., 2016; Zylberberg et al., 2012), although no evidence for a

727     response-congruent confirmation bias was observed in other experimental paradigms (Rausch et

728     al., 2020; Shekhar & Rahnev, 2022), suggesting a response-congruent confirmation bias many

729     not be a universal feature of human confidence across paradigms. However, there are multiple

730     mathematical ways to represent bias in favour of response-congruent evidence. When we

731     implemented a response-congruent evidence bias in a different way, resulting in the model we

732     refer to as response-congruent evidence model, meta-d′/d′ very strongly correlated with

733     discrimination sensitivity. This finding implies that it is not sufficient that the generative process

734     underlying the confidence data is characterised by a similar conceptual idea as ITG - if meta-d′/d′

735    is to control for by discrimination sensitivity, discrimination bias, and confidence criteria, ITG

736    must be (at least a close approximation of) the generative model of the data.

737            An important limitation shared between ITG and all alternative models investigated in the

738    present study is that the dynamics of the decision process is not accounted for. This is

739    problematic because there is a large body of evidence that confidence judgments depend on the

740    dynamics of decision making (Pleskac & Busemeyer, 2010). Specifically, Pleskac and

741    Busemeyer (2010) showed that when participants are under time pressure when making the

742    decision, metacognitive efficiency is increased. Moran et al. (2015) showed that confidence

743    judgments are related to the reaction time of confidence judgments. Last but not least, there is on

744    average  medium-sized correlation between confidence judgments and reaction time across a

745    wide range of studies (Rahnev et al., 2020). Given the close relationship between decision

746    dynamics and confidence, it may be more apt to model confidence with sequential sampling

747    models rather than signal detection theory (Desender et al., 2022; Hellmann et al., 2023; Pereira

748    et al., 2021; Pleskac & Busemeyer, 2010; Ratcliff & Starns, 2009, 2013; Reynolds et al., 2020).

749    **Measuring metacognitive efficiency using meta-d′/d′**

750            The findings of the present study imply that whenever the independent truncated

751    Gaussian model is a good description of the data, meta-d′/d′ will be the appropriate measure of

752    metacognitive efficiency. However, without any information about the generative model

753    underlying confidence judgments, researchers should not assume that by using meta-d′/d′ to

754    measure metacognitive efficiency, a potential contamination by discrimination sensitivity,

755    discrimination bias, or confidence criteria has been ruled out. We recommend to use meta-d′/d′

756    only for tasks where the independent truncated Gaussian model is a suitable description of the

757    data. There is a limited set of experimental tools available to reduce the potential impact of

758    discrimination sensitivity, discrimination bias, and confidence criteria when measuring

759    metacognitive efficiency using meta-d′/d′. To control for discrimination sensitivity, researchers

760    have used staircases to keep task performance within a specific range (Rahnev & Fleming,

761    2019). However, Simulation 2 suggests that staircases are not sufficient to control for

762    discrimination sensitivity if the data is generated according to the weighted evidence and

763    visibility model or the response-congruent evidence model. It might be possible to reduce the

764    impact of discrimination criteria and confidence criteria by careful instructions and training with

765    the task, although it is unlikely that instruction and training is sufficient to eliminate the effect of

766    criteria.

767            Measuring metacognitive efficiency by meta-d′/d′ is also problematic because meta-d′/d′

768    does not take the dynamics of the decision process into account. Consequently, properties of the

769    dynamical decision process such as response caution might be misinterpreted as effects on

770    metacognitive efficiency (Desender et al., 2022). Overall, the findings of the present study

771    combined with other recent studies (Desender et al., 2022; Guggenmos, 2021; Shekhar &

772    Rahnev, 2021) imply that without any additional information, meta-d′/d′ cannot be

773    unambiguously interpreted in terms of metacognitive efficiency, suggesting that a reanalysis of

774    previously published studies using meta-d′/d′ and possibly a critical reinterpretation is necessary.

775    **Alternatives to meta-d′/d′ for measuring metacognitive efficiency**

776            Whenever ITG is not a decent description of confidence in a particular study, researchers

777    need an alternative to meta-d′-d′ to measure metacognitive efficiency. Traditionally,

778    metacognition has been assessed using measures that also do not explicitly rely on specific

779    generative models of confidence, such as gamma correlation coefficients (Nelson, 1984),

780    confidence slopes (Yates, 1990), phi correlations (Rounis et al., 2010), or area under type 2 ROC

781    curves (Fleming et al., 2010). However, none of these measures is designed to control for

782    discrimination performance and thus, by definition, none of these measures are measures of

783    metacognitive efficiency.

784         There are several model-based alternative measures of metacognitive efficiency: First,

785    one available method is to fit a lognormal noise model, in which metacognitive ability is

786    quantified by the lognormal noise parameter $\sigma_{meta}$ (Shekhar & Rahnev, 2021, 2022). The

787    lognormal noise model provides a decent account for confidence in a low contrast orientation

788    discrimination task as well as a letter numerosity discrimination task (Shekhar & Rahnev, 2022).

789    Second, in two-alternative forced choice confidence paradigms, it is possible to quantify

790    metacognitive efficiency using the confidence boost model (Mamassian & de Gardelle, 2021).

791    The measure of metacognitive efficiency $\eta$ is computed by dividing the variance of the

792    confidence noise of a hypothetical ideal observer by the variance of confidence noise estimated

793    for the participant. Besides, two-alternative forced choice confidence paradigms may be an

794    attractive way to eliminate the impact of confidence criteria (Barthelmé & Mamassian, 2009).

795    Finally, relying on two-stage signal detection theory (Pleskac & Busemeyer, 2010; Yu et al.,

796    2015), Desender et al. (2022) proposed the v-ratio to measure metacognitive efficiency. The v-

797    ratio divides the drift rate estimated from confidence judgments by the drift rate estimated from

798    discrimination responses and reaction time.

799         Notably, just as meta-d′/d′ is only a good measure of metacognitive efficiency when the

800    data confirm to the independent truncated Gaussian model, $\sigma_{meta}$, $\eta$, and v-ratio are expected to

801    control for discrimination sensitivity, discrimination bias and confidence criteria only when the

802    data confirm to the corresponding generative model. To our knowledge, it has not yet been

803    investigated how sensitive $\sigma_{meta}$, $\eta$, and v-ratio are to a contamination from discrimination

804    sensitivity, discrimination bias and confidence criteria are when generative model underlying

805    confidence judgment is varied. The findings of the present study are consistent with the view that

806    measures of metacognitive efficiency provide control over discrimination sensitivity,

807    discrimination bias and confidence criteria only if the generative model of confidence is

808    correctly identified and the corresponding measure of metacognitive efficiency is used

809    (Guggenmos, 2022). Unfortunately, for the time being, there is no consensus about the

810    computational principles underlying confidence judgments (Rahnev et al., 2022). This means

811    that a good practice for future studies will be to first use cognitive modelling to identify the

812    generative model underlying confidence judgments in a specific paradigm empirically, and then

813    use the corresponding model-based measure of metacognitive efficiency (Guggenmos, 2021;

814    Mamassian & de Gardelle, 2021; Shekhar & Rahnev, 2021). When data in a specific task is well

815    accounted for by the independent truncated Gaussian model, meta-d′/d′ is the appropriate way to

816    measure metacognitive efficiency. However, when data is better described by an alternative

817    model of confidence, researchers need to use a measure of metacognitive efficiency that

818    corresponds to the model that is the best explanation of the data. Because researchers have

819    implicitly fitted versions of the independent truncated Gaussian model all along when they used

820    meta-d′/d′, it does not seem too far-fetched that researchers will begin to regularly fit alternative

821    generative models of confidence as well. It will be necessary to develop open and easy-to-use

822    software packages to make fitting a variety of confidence models available to a larger part of the

823    field (e.g., Rausch & Hellmann, 2023). Sometimes it will be impossible to identify the true

824    generative model underlying confidence judgments for a specific data set, either because the

825    number of trials is too low or because of model mimicry. In these cases, it will be prudent to

826    perform a robustness analysis to show that the results of the study do not depend on specific

827    analysis decisions (Gelman & Loken, 2014; Steegen et al., 2016). This means that the modelling

828    analysis needs to be repeated with all models of confidence that cannot be ruled out empirically

829    to show that results are robust across models of confidence.

830        It is very difficult, and perhaps impossible, to come up with a novel measure of

831    metacognitive efficiency with all the attractive properties that meta-d′/d′ was supposed to have,

832    i.e., controlling for discrimination sensitivity, discrimination bias, and confidence criteria

833    without requiring a specific generative model of confidence. The present study does not rule out

834    the possibility that a future study will be able to find such a measure. However, given the results

835    of the present study, we are sceptical that such a measure can ever be found; we recommend

836    rigorous testing of whether any newly proposed measure of metacognitive efficiency effectively

837    controls for discrimination performance, discrimination bias, and confidence criteria.

838    **Conclusion**

839        We showed that meta-d′/d′ is not free from assumptions about the generative model

840    underlying confidence judgments. Only if the data is generated according to the independent

841    truncated gaussian model, meta-d′/d′ guarantees control over discrimination performance,

842    discrimination bias, and confidence criteria. The control fails according to a wide range of

843    alternative generative models of confidence; the expected correlation with discrimination

844    sensitivity and confidence criteria varies across alternative generative model but can be very

845    large. Consequently, researchers who want to measure metacognitive efficiency using meta-d′/d′

846    need to examine if their data can be reasonably described by the independent truncated Gaussian

847    model.

848                                    **References**

849    Adler, W. T., & Ma, W. J. (2018). Comparing Bayesian and non-Bayesian accounts of human

850        confidence reports. *PLOS Computational Biology*, *14*(11), e1006572.

851        https://doi.org/10.1371/journal.pcbi.1006572

852    Aitchison, L., Bang, D., Bahrami, B., & Latham, P. E. (2015). Doubly Bayesian Analysis of

853        Confidence in Perceptual Decision-Making. *PLoS Computational Biology*, *11*(10),

854        e1004519. https://doi.org/10.1371/journal.pcbi.1004519

855    Akaike, H. (1974). A New Look at the Statistical Model Identification. *IEEE transactions on*

856        *automatic control*, *AC-19*(6), 716–723. https://doi.org/10.1007/978-1-4612-1694-0_16

857    Barrett, A. B., Dienes, Z., & Seth, A. K. (2013). Measures of metacognition on signal-detection

858        theoretic models. *Psychological Methods*, *18*(4), 535–552.

859        https://doi.org/10.1037/a0033268

860    Barthelmé, S., & Mamassian, P. (2009). Evaluation of Objective Uncertainty in the Visual

861        System. *PLoS Computational Biology*, *5*(9), e1000504. https://doi.org/10.1371/Citation

862    Bhome, R., McWilliams, A., Price, G., Poole, N. A., Howard, R. J., Fleming, S. M., & Huntley,

863        J. D. (2022). Metacognition in functional cognitive disorder. *Brain Communications*,

864        *4*(2), fcac041. https://doi.org/10.1093/braincomms/fcac041

865    Boldt, A., Blundell, C., & De Martino, B. (2019). Confidence modulates exploration and

866        exploitation in value-based learning. *Neuroscience of Consciousness*, *2019*(1), niz004.

867        https://doi.org/10.1093/nc/niz004

868    Boundy-Singer, Z. M., Ziemba, C. M., & Goris, R. L. T. (2022). Confidence reflects a noisy

869        decision reliability estimate. *Nature Human Behaviour*, *7*(1), 142–154.

870        https://doi.org/10.1038/s41562-022-01464-x

871    Burnham, K. P., & Anderson, D. R. (2002). *Model selection and multimodel inference: A*

872        *practical information-theoretic approach* (2. Aufl.). Springer.

873    Charles, L., Van Opstal, F., Marti, S., & Dehaene, S. (2013). Distinct brain mechanisms for

874        conscious versus subliminal error detection. *NeuroImage*, *73*, 80–94.

875        https://doi.org/10.1016/j.neuroimage.2013.01.054

876    Clarke, F. R., Birdsall, T. G., & Tanner, W. P. (1959). Two Types of ROC Curves and

877        Definitions of Parameters. *The Journal of the Acoustical Society of America*, *31*(5), 629–

878        630. https://doi.org/10.1121/1.1907764

879    Culot, C., Fantini-Hauwel, C., & Gevers, W. (2021). The influence of sad mood induction on

880        task performance and metacognition. *Quarterly Journal of Experimental Psychology*,

881        *74*(9), 1605–1614. https://doi.org/10.1177/17470218211004205

882    Desender, K., Donner, T. H., & Verguts, T. (2021). Dynamic expressions of confidence within

883        an evidence accumulation framework. *Cognition*, *207*(104522), 1–11.

884        https://doi.org/10.1016/j.cognition.2020.104522

885    Desender, K., Vermeylen, L., & Verguts, T. (2022). Dynamic influences on static measures of

886        metacognition. *Nature Communications*, *13*(1), 1–30. https://doi.org/10.1038/s41467-

887        022-31727-0

888    Drescher, L. H., Van den Bussche, E., & Desender, K. (2018). Absence without leave or leave

889        without absence: Examining the interrelations among mind wandering, metacognition

890        and cognitive control. *PLOS ONE*, *13*(2), e0191639.

891        https://doi.org/10.1371/journal.pone.0191639

892    Fischer, H., & Said, N. (2021). Importance of domain-specific metacognition for explaining

893          beliefs about politicized science: The case of climate change. *Cognition*, *208*, 104545.

894          https://doi.org/10.1016/j.cognition.2020.104545

895    Fleming, S. M. (2017). HMeta-d: Hierarchical Bayesian estimation of metacognitive efficiency

896          from confidence ratings. *Neuroscience of Consciousness*, *1*, 1–14.

897          https://doi.org/10.1093/nc/nix007

898    Fleming, S. M., & Daw, N. D. (2017). Self-evaluation of decision-making: A general Bayesian

899          framework for metacognitive computation. *Psychological Review*, *124*(1), 91–114.

900          https://doi.org/10.1037/rev0000045

901    Fleming, S. M., & Lau, H. C. (2014). How to measure metacognition. *Frontiers in human*

902          *neuroscience*, *8*(443), 1–9. https://doi.org/10.3389/fnhum.2014.00443

903    Fleming, S. M., Weil, R. S., Nagy, Z., Dolan, R. J., & Rees, G. (2010). Relating introspective

904          accuracy to individual differences in brain structure. *Science*, *329*(5998), 1541–1543.

905          https://doi.org/10.1126/science.1191883

906    Galvin, S. J., Podd, J. V., Drga, V., & Whitmore, J. (2003). Type 2 tasks in the theory of signal

907          detectability: Discrimination between correct and incorrect decisions. *Psychonomic*

908          *Bulletin & Review*, *10*(4), 843–876.

909    Gelman, A., & Loken, E. (2014). The statistical Crisis in science. *American Scientist*, *102*(6),

910          460–465. https://doi.org/10.1511/2014.111.460

911    Gelman, A., & Rubin, D. B. (1992). Inference from iterative simulation using multiple

912          sequences. *Statistical Science*, *7*(4), 457–511.

913    Green, D. M., & Swets, J. A. (1966). *Signal detection theory and psychophysics*. Wiley.

914 Guggenmos, M. (2021). Measuring metacognitive performance: Type 1 performance dependence

915        and test-retest reliability. *Neuroscience of Consciousness*, *7*(1), 1–14.

916        https://doi.org/10.1093/nc/niab040

917 Guggenmos, M. (2022). Reverse engineering of metacognition. *eLife*, *11*, 1–29.

918        https://doi.org/10.7554/eLife.75420

919 Hainguerlot, M., Vergnaud, J.-C., & de Gardelle, V. (2018). Metacognitive ability predicts

920        learning cue-stimulus associations in the absence of external feedback. *Scientific Reports*,

921        *8*(1), 5602. https://doi.org/10.1038/s41598-018-23936-9

922 Hellmann, S., Zehetleitner, M., & Rausch, M. (2023). Simultaneous modeling of choice,

923        confidence, and response time in visual perception. *Psychological Review*.

924        https://doi.org/10.1037/rev0000411

925 Kristensen, S. B., Sandberg, K., & Bibby, B. M. (2020). Regression methods for metacognitive

926        sensitivity. *Journal of Mathematical Psychology*, *94*(102297), 1–17.

927        https://doi.org/10.1016/j.jmp.2019.102297

928 Lee, M. D., & Wagenmakers, E.-J. (2013). *Bayesian Cognitive Modeling: A Practical Course*.

929        Cambridge University Press.

930 Macmillan, N. A., & Creelman, C. D. (2005). *Detection Theory. A user's guide*. Lawrence

931        Erlbaum Associates.

932 Mamassian, P., & de Gardelle, V. (2021). Modeling Perceptual Confidence and the Confidence

933        Forced-Choice Paradigm. *Psychological Review*, 1–23.

934        https://doi.org/10.1037/rev0000312

935   Maniscalco, B., & Lau, H. (2012). A signal detection theoretic method for estimating

936         metacognitive sensitivity from confidence ratings. *Consciousness and Cognition*, *21*(1),

937         422–430.

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

939         reports of sensory awareness. *Neuroscience of Consciousness*, *1*, 1–17.

940         https://doi.org/10.1093/nc/niw002

941   Maniscalco, B., & Lau, H. C. (2014). Signal Detection Theory Analysis of Type 1 and Type 2

942         Data: Meta-d', Response- Specific Meta-d', and the Unequal Variance SDT Model. In S.

943         M. Fleming & C. D. Frith (Hrsg.), *The Cognitive Neuroscience of Metacognition* (S. 25–

944         66). Springer. https://doi.org/10.1007/978-3-642-45190-4_3

945   Maniscalco, B., McCurdy, L. Y., Odegaard, B., & Lau, H. (2017). Limited Cognitive Resources

946         Explain a Trade-Off between Perceptual and Metacognitive Vigilance. *The Journal of*

947         *Neuroscience*, *37*(5), 1213–1224. https://doi.org/10.1523/JNEUROSCI.2271-13.2016

948   Maniscalco, B., Peters, M. A. K., & Lau, H. (2016). Heuristic use of perceptual evidence leads to

949         dissociation between performance and metacognitive sensitivity. *Attention, Perception &*

950         *Psychophysics*, *78*, 923–937. https://doi.org/10.3758/s13414-016-1059-x

951   Mazancieux, A., Dinze, C., Souchay, C., & Moulin, C. J. A. (2020). Metacognitive domain

952         specificity in feeling-of-knowing but not retrospective confidence. *Neuroscience of*

953         *Consciousness*, *2020*(1), niaa001. https://doi.org/10.1093/nc/niaa001

954   Moran, R., Teodorescu, A. R., & Usher, M. (2015). Post choice information integration as a

955         causal determinant of confidence: Novel data and a computational account. *Cognitive*

956         *Psychology*, *78*, 99–147. https://doi.org/10.1016/j.cogpsych.2015.01.002

957    Muthesius, A., Grothey, F., Cunningham, C., Hölzer, S., Vogeley, K., & Schultz, J. (2022).

958            Preserved metacognition despite impaired perception of intentionality cues in

959            schizophrenia. *Schizophrenia Research: Cognition*, *27*, 100215.

960            https://doi.org/10.1016/j.scog.2021.100215

961    Nelder, J. A., & Mead, R. (1965). A simplex method for function minimization. *The Computer*

962            *Journal*, *7*(4), 308–313. https://doi.org/10.1093/COMJNL/7.4.308

963    Nelson, T. O. (1984). A Comparison of Current Measures of the Accuracy of Feeling-of-

964            Knowing Predictions. *Psychological Bulletin*, *95*(1), 109–133.

965    Odegaard, B., Chang, M. Y., Lau, H., & Cheung, S.-H. (2018). Inflation versus filling-in: Why

966            we feel we see more than we actually do in peripheral vision. *Phil. Trans. R. Soc. B*, *373*,

967            20170345. https://doi.org/10.1098/rstb.2017.0345

968    Odegaard, B., Grimaldi, P., Cho, S. H., Peters, M. A. K., Lau, H., & Basso, M. A. (2018).

969            Superior colliculus neuronal ensemble activity signals optimal rather than subjective

970            confidence. *Proceedings of the National Academy of Sciences*, *115*(7), E1588–E1597.

971            https://doi.org/10.1073/pnas.1711628115

972    Pereira, M., Megevand, P., Tan, M. X., Chang, W., Wang, S., Rezai, A., Seeck, M., Corniola,

973            M., Momjian, S., Bernasconi, F., Blanke, O., & Faivre, N. (2021). Evidence

974            accumulation relates to perceptual consciousness and monitoring. *Nature*

975            *Communications*, *12*(1), 3261. https://doi.org/10.1038/s41467-021-23540-y

976    Peters, M. A. K., Thesen, T., Ko, Y. D., Maniscalco, B., Carlson, C., Davidson, M., Doyle, W.,

977            Kuzniecky, R., Devinsky, O., Halgren, E., & Lau, H. (2017). Perceptual confidence

978            neglects decision-incongruent evidence in the brain. *Nature Human Behaviour*, *1*(0139),

979            1–21. https://doi.org/10.1038/s41562-017-0139

980    Peterson, W. W., Birdsall, T. G., & Fox, W. C. (1954). The theory of signal detectability.

981          *Transactions of the IRE Professional Group on Information Theory*, *4*(4), 171–212.

982          https://doi.org/10.1109/TIT.1954.1057460

983    Pleskac, T. J., & Busemeyer, J. R. (2010). Two-Stage Dynamic Signal Detection: A Theory of

984          Choice , Decision Time, and Confidence. *Psychological Review*, *117*(3), 864–901.

985          https://doi.org/10.1037/a0019737

986    Plummer, M. (2003). JAGS: A Program for Analysis of Bayesian Graphical Models Using Gibbs

987          Sampling. *Proceedings of the 3rd International Workshop on Distributed Statistical*

988          *Computing (DSC 2003)*. http://www.ci.tuwien.ac.at/Conferences/DSC-2003/

989    Pollack, I. (1959). On Indices of Signal and Response Discriminability. *Journal of the Acoustical*

990          *Society of America*, *31*, 1031.

991    R Core Team. (2020). *R: A Language and Environment for Statistical Computing*. R Foundation

992          for Statistical Computing. https://www.r-project.org/

993    Rahnev, D., Balsdon, T., Charles, L., de Gardelle, V., Denison, R., Desender, K., Faivre, N.,

994          Filevich, E., Fleming, S. M., Jehee, J., Lau, H., Lee, A. L. F., Locke, S. M., Mamassian,

995          P., Odegaard, B., Peters, M. A. K., Reyes, G., Rouault, M., Sackur, J., … Zylberberg, A.

996          (2022). Consensus Goals in the Field of Visual Metacognition. *Perspectives on*

997          *Psychological Science*, *17*(6), 1746–1765. https://doi.org/10.1177/174569162210756

998    Rahnev, D., Desender, K., Lee, A. L. F., Adler, W. T., Aguilar-Lleyda, D., Akdoğan, B.,

999          Arbuzova, P., Atlas, L. Y., Balcı, F., Bang, J. W., Bègue, I., Birney, D. P., Brady, T. F.,

1000          Calder-Travis, J., Chetverikov, A., Clark, T. K., Davranche, K., Denison, R. N., Dildine,

1001          T. C., … Zylberberg, A. (2020). The Confidence Database. *Nature Human Behaviour*, *4*,

1002          317–325. https://doi.org/10.1038/s41562-019-0813-1

Rahnev, D., & Fleming, S. M. (2019). How experimental procedures influence estimates of metacognitive ability. *Neuroscience of Consciousness*, *5*(1), 1–9. https://doi.org/10.1093/nc/niz009

Ratcliff, R., & Starns, J. J. (2009). Modeling Confidence and Response Time in Recognition Memory. *Psychological Review*, *116*(1), 59–83. https://doi.org/10.1037/a0014086

Ratcliff, R., & Starns, J. J. (2013). Modeling Confidence Judgments, Response Times, and Multiple Choices in Decision Making: Recognition Memory and Motion Discrimination. *Psychological Review*, *120*(3), 697–719. https://doi.org/10.1037/a0033152

Rausch, M., & Hellmann, S. (2023). *statConfR: Models of Decision Confidence and Metacognition* (0.0.1) [R]. https://cran.r-project.org/web/packages/statConfR/index.html

Rausch, M., Hellmann, S., & Zehetleitner, M. (2018). Confidence in masked orientation judgments is informed by both evidence and visibility. *Attention, Perception, and Psychophysics*, *80*(1), 134–154. https://doi.org/10.3758/s13414-017-1431-5

Rausch, M., Hellmann, S., & Zehetleitner, M. (2021). Modelling visibility judgments using models of decision confidence. *Attention, Perception & Psychophysics*, *83*, 3311–3336. https://doi.org/10.3758/s13414-021-02284-3

Rausch, M., & Zehetleitner, M. (2016). Visibility is not equivalent to confidence in a low contrast orientation discrimination task. *Frontiers in Psychology*, *7*(591), 1–15. https://doi.org/10.3389/fpsyg.2016.00591

Rausch, M., & Zehetleitner, M. (2019). The folded X-pattern is not necessarily a statistical signature of decision confidence. *PLoS Computational Biology*, *15*(10), e1007456. https://doi.org/10.1371/journal.pcbi.1007456

1025    Rausch, M., & Zehetleitner, M. (2023). Evaluating false positive rates of standard and

1026          hierarchical measures of metacognitive accuracy. *Metacognition and Learning*.

1027          https://doi.org/10.1007/s11409-023-09353-y

1028    Rausch, M., Zehetleitner, M., Steinhauser, M., & Maier, M. E. (2020). Cognitive modelling

1029          reveals distinct electrophysiological markers of decision confidence and error monitoring.

1030          *NeuroImage*, *218*(116963), 1–14. https://doi.org/10.1016/j.neuroimage.2020.116963

1031    Reynolds, A., Kvam, P. D., Osth, A. F., & Heathcote, A. (2020). Correlated racing evidence

1032          accumulator models. *Journal of Mathematical Psychology*, *96*, 102331.

1033          https://doi.org/10.1016/j.jmp.2020.102331

1034    Rouault, M., Seow, T., Gillan, C. M., & Fleming, S. M. (2018). Psychiatric Symptom

1035          Dimensions Are Associated With Dissociable Shifts in Metacognition but Not Task

1036          Performance. *Biological Psychiatry*, *84*(6), 443–451.

1037          https://doi.org/10.1016/j.biopsych.2017.12.017

1038    Rounis, E., Maniscalco, B., Rothwell, J. C., Passingham, R. E., & Lau, H. (2010). Theta-burst

1039          transcranial magnetic stimulation to the prefrontal cortex impairs metacognitive visual

1040          awareness. *Cognitive Neuroscience*, *1*(3), 165–175.

1041          https://doi.org/10.1080/17588921003632529

1042    Said, N., Fischer, H., & Anders, G. (2022). Contested science: Individuals with higher

1043          metacognitive insight into interpretation of evidence are less likely to polarize.

1044          *Psychonomic Bulletin & Review*, *29*(2), 668–680. https://doi.org/10.3758/s13423-021-

1045          01993-y

1046    Samaha, J., Barrett, J. J., Sheldon, A. D., LaRocque, J. J., & Postle, B. R. (2016). Dissociating

1047          perceptual confidence from discrimination accuracy reveals no influence of

metacognitive awareness on working memory. *Frontiers in Psychology*, *7*(851), 1–8.

https://doi.org/10.3389/fpsyg.2016.00851

Schönbrodt, F. D., & Perugini, M. (2013). At what sample size do correlations stabilize ? *Journal*

*of Research in Personality*, *47*(5), 609–612. https://doi.org/10.1016/j.jrp.2013.05.009

Schwarz, G. (1978). Estimating the dimensions of a model. *The Annals of Statistics*, *6*(2), 461–

464. https://doi.org/10.1214/aos/1176348654

Shekhar, M., & Rahnev, D. (2021). The Nature of Metacognitive Inefficiency in Perceptual

Decision Making. *Psychological Review*, *128*(1), 45–70.

https://doi.org/10.1037/rev0000249

Shekhar, M., & Rahnev, D. (2022). How do humans give confidence? A comprehensive

comparison of process models of metacognition. *PsyArXiv*.

https://doi.org/10.31234/osf.io/cwrnt

Steegen, S., Tuerlinckx, F., Gelman, A., & Vanpaemel, W. (2016). Increasing Transparency

Through a Multiverse Analysis. *Perspectives on Psychological Science*, *11*(5), 702–712.

https://doi.org/10.1177/1745691616658637

Tanner, W. P., Jr., & Swets, J. A. (1954). A decision-making theory of visual detection.

*Psychological Review*, *61*(6), 401–409. https://doi.org/10.1037/h0058700

Taouki, I., Lallier, M., & Soto, D. (2022). The role of metacognition in monitoring performance

and regulating learning in early readers. *Metacognition and Learning*, 58.

https://doi.org/10.1007/s11409-022-09292-0

Vandenbroucke, A. R. E., Sligte, I. G., Barrett, A. B., Seth, A. K., Fahrenfort, J. J., & Lamme, V.

A. F. (2014). Accurate Metacognition for Visual Sensory Memory Representations.

*Psychological Science*, *25*(4), 861–873. https://doi.org/10.1177/0956797613516146

1071    Vlassova, A., Donkin, C., & Pearson, J. (2014). Unconscious information changes decision

1072        accuracy but not confidence. *Proceedings of the National Academy of Sciences*, *111*(45),

1073        16214–16218. https://doi.org/10.1073/pnas.1403619111

1074    Wickens, T. D. (2002). *Elementary signal detection theory*. Oxford University Press.

1075    Yates, J. F. (1990). *Judgment and decision making*. Prentice Hall.

1076    Yu, S., Pleskac, T. J., & Zeigenfuse, M. D. (2015). Dynamics of Postdecisional Processing of

1077        Confidence. *Journal of Experimental Psychology: General*, *144*(2), 489–510.

1078        https://doi.org/10.1037/xge0000062

1079    Zhu, J.-Q., Sundh, J., Spicer, J., Chater, N., & Sanborn, A. N. (2023). The autocorrelated

1080        Bayesian sampler: A rational process for probability judgments, estimates, confidence

1081        intervals, choices, confidence judgments, and response times. *Psychological Review*.

1082        https://doi.org/10.1037/rev0000427

1083    Zylberberg, A., Barttfeld, P., & Sigman, M. (2012). The construction of confidence in a

1084        perceptual decision. *Frontiers in integrative Neuroscience*, *6*(79), 1–10.

1085        https://doi.org/10.3389/fnint.2012.00079

1086