

1 Running head: logistic regression and metacognition

2

3

4

5

6 **Should metacognition be measured by logistic regression?**

7

8

9

10 Manuel Rausch ^{1,2} and Michael Zehetleitner ^{1,2}

11

12 ¹ Katholische Universität Eichstätt-Ingolstadt, Eichstätt, Germany

13

14 ² Ludwig-Maximilians-Universität München, Munich, Germany

15

16

17 Correspondence should be addressed at:

18 Manuel Rausch

19 Katholische Univerität Eichstätt-Ingolstadt

20 Psychologie II

21 Ostenstraße 25, "Waisenhaus"

22 85072 Eichstätt

23 Germany

24 Phone: +49 8421 93 21639

25 Email: manuel.rausch@ku.de

26 <http://www.ku.de/ppf/psychologie/psych2/mitarbeiter/m-rausch/>

27

28 This manuscript was accepted for publication in *Consciousness and Cognition*. Please cite
29 this work as: Rausch, M., Zehetleitner, M. (2017). Should metacognition be measured by
30 logistic regression? *Consciousness and Cognition*. 49, 291-312. The final publication is
31 available at <http://dx.doi.org/10.1016/j.concog.2017.02.007>

32 © 2017. This manuscript version is made available under the CC-BY-NC-ND 4.0 license

33 <http://creativecommons.org/licenses/by-nc-nd/4.0/>

34

35

Abstract

36 Are logistic regression slopes suitable to quantify metacognitive sensitivity, i.e. the efficiency
37 with which subjective reports differentiate between correct and incorrect task responses? We
38 analytically show that logistic regression slopes are independent from rating criteria in one
39 specific model of metacognition, which assumes (i) that rating decisions are based on sensory
40 evidence generated independently of the sensory evidence used for primary task responses
41 and (ii) that the distributions of evidence are logistic. Given a hierarchical model of
42 metacognition, logistic regression slopes depend on rating criteria. According to all
43 considered models, regression slopes depend on the primary task criterion. A reanalysis of
44 previous data revealed that massive numbers of trials are required to distinguish between
45 hierarchical and independent models with tolerable accuracy. It is argued that researchers
46 who wish to use logistic regression as measure of metacognitive sensitivity need to control
47 the primary task criterion and rating criteria.

48 *Keywords:* metacognition; metacognitive sensitivity, logistic regression; signal detection
49 theory; type 2 signal detection theory; generalized linear regression, cognitive modelling

50

51 1 Introduction

52 Metacognitive sensitivity, also called type 2 sensitivity or resolution of confidence, refers to
 53 the efficiency with which participants' subjective reports during an experimental task
 54 discriminate between correct and incorrect responses in a primary task (Baranski & Petrusic,
 55 1994; Fleming & Lau, 2014; Galvin, Podd, Drga, & Whitmore, 2003). It relates to a key
 56 aspect of metacognition: If participants possessed any knowledge about their performance in
 57 the task, their subjective reports about the task should differentiate between correct and
 58 erroneous trials. Consequently, measures of metacognitive sensitivity are relevant for all
 59 research areas where quantifying participants' insight into their task performance is of
 60 interest, including consciousness research (Dienes, 2004). Given the theoretical importance
 61 of metacognitive sensitivity, a universally accepted measure is desirable. However, various
 62 competing measures of metacognitive sensitivity were proposed in the literature:

- 63 • gamma correlation coefficients (Nelson, 1984),
- 64 • $a'/\text{type } 2 \text{ } d'$ (Kunimoto, Miller, & Pashler, 2001)
- 65 • type-2 receiver operating characteristic (Fleming, Weil, Nagy, Dolan, & Rees,
66 2010)
- 67 • meta- d' (Maniscalco & Lau, 2012)
- 68 • logistic regression analysis (Sandberg, Timmermans, Overgaard, &
69 Cleeremans, 2010)

70 Logistic regression has been widely used in empirical studies as measure of the association
 71 between verbal reports and task accuracy (Rausch, Müller, & Zehetleitner, 2015; Rausch &
 72 Zehetleitner, 2014; Sandberg et al., 2010; Siedlecka, Paulewicz, & Wierzchoń, 2016;
 73 Wierzchoń, Asanowicz, Paulewicz, & Cleeremans, 2012; Wierzchoń, Paulewicz, Asanowicz,
 74 Timmermans, & Cleeremans, 2014). However, while gamma correlations, a' , and meta- d'
 75 have been extensively examined using both empirical and analytical methods (Barrett,
 76 Dienes, & Seth, 2013; Evans & Azzopardi, 2007; Galvin et al., 2003; Masson & Rotello,
 77 2009), the conditions for logistic regression to be an appropriate measure of metacognitive
 78 sensitivity have never been systematically investigated.

79 1.1 Logistic regression as measure of metacognitive sensitivity

80 Logistic regression is a specific case of a generalized linear regression model (GLM). In
 81 general, it is a method to quantify the relationship between a binary outcome variable and one
 82 or several dichotomous or continuous predictors. The standard approach to quantify
 83 metacognitive sensitivity by means of logistic regression is to model the probability of being
 84 correct in the primary task $P(T)$ as a linear function of a subjective report C , e.g. a confidence
 85 judgment or a visibility rating. A linear relationship between predictors and outcome is
 86 obtained by transforming the probability of being correct into the logarithm of the odds of the
 87 primary response being correct to being incorrect:

$$\log\left(\frac{P(T)}{1 - P(T)}\right) = a + b * C \quad (1)$$

88 As can be seen from Fig. 1, metacognitive sensitivity is indexed by the slope b of the
 89 regression line: the steeper the regression line, the stronger are subjective reports associated
 90 with the probability of being correct (Sandberg et al., 2010). Logistic regression is also used
 91 to quantify the minimal criteria participants apply when they make a subjective report: The
 92 more conservative participants' reporting strategy is, the better they perform while still giving

93 the lowest possible subjective report. As the intercept a is just the transformed accuracy when
94 the subjective report is zero, it is interpreted as measure of criterion (Wierchoń et al., 2012).

95 Quantifying metacognition by logistic regression is tempting due to three reasons: First, the
96 hierarchical structure of the data often found in behavioral experiments can be explicitly
97 included into the model by using nested random effects (Sandberg, Bibby, & Overgaard,
98 2013; Siedlecka et al., 2016): For example, two experimental groups with several participants
99 each contributing a number of trials can be described by a random effect of trial nested within
100 a random effect of participant nested within groups. As such an analysis can be conducted on
101 a single trial level without the need for summary statistics to be computed for each
102 participant, logistic regression may also be a promising way to increase statistical power
103 (Sandberg et al., 2010). Second, using random effects allows the data to be unbalanced, i.e.
104 the number of observations can vary between conditions or there can be empty cells in the
105 design matrix (Rausch et al., 2015; Siedlecka et al., 2016). Consequently, slopes on the group
106 level can be obtained even when not all participants made errors in all experimental
107 conditions. This is particularly useful in studies of metacognitive sensitivity because the
108 number of errors may vary heavily between participants and conditions. Finally, it has been
109 argued that logistic regression, unlike SDT-measures, does not make any assumptions about
110 the sources of the evidence involved in making subjective reports (Siedlecka et al., 2016).

111 However, logistic regression as measure of metacognition may also suffer from at least two
112 drawbacks: First, it depends on the assumption that the relation between subjective reports
113 and transformed accuracy is linear. A non-linear relationship implies that there is no single
114 slope of the regression line that could be interpreted as measure of metacognitive sensitivity.
115 A previous analysis suggested non-linear trends between subjective reports and logit-
116 transformed accuracy occur relatively frequently, although there were also some data sets
117 where the assumption of a linear relationship appeared to be justified (Rausch et al., 2015). It
118 should be noted that this critique only applies to rating scales with more than two response
119 options because two data from two response options always can be connected by a straight
120 line.

121 The present analysis explores a potential second and more principal problem of logistic
122 regression: An adequate measure of metacognitive sensitivity should depend exclusively on
123 the amount of evidence available for subjective reports, and should not be confounded by
124 other factors such as rating criteria and response biases (Barrett et al., 2013). Although the
125 distinction between slopes and intercepts appears superficially similar to a separation
126 between metacognitive sensitivity and criteria, it has never been systematically investigated
127 what exactly are the assumptions that have to be fulfilled so that slopes are independent from
128 criteria.

129 **1.2 Models of metacognition**

130 A systematic investigation of the impact of rating criteria and primary task criterion on
131 logistic regression models requires a mathematically formulated model that accounts for both
132 task responses as well as subjective reports. One of the most prominent models of perceptual
133 decision making under uncertainty is signal detection theory (SDT) (Green & Swets, 1966;
134 Macmillan & Creelman, 2005; Wickens, 2002). Standard SDT applies to tasks where
135 participants are instructed to correctly classify a binary stimulation S . For the purpose of
136 present analysis, it is not relevant whether the two variants of the stimulation are interpreted
137 as signal and noise or as two different stimuli. Each stimulus provides the participants with
138 sensory evidence which of the two response options he or she should select. Participants

139 select their responses based on a comparison of the sensory evidence with a primary task
140 criterion θ . θ represents the degree to which observers tend towards one response option
141 independent of the sensory evidence. Participants respond 0 if the sensory evidence is smaller
142 than θ and 1 otherwise. As there is noise in the system, the sensory evidence is not always the
143 same at each presentation of the stimulus, but instead is modeled as a random sample out of
144 two distributions, one for each variant of S. If the observer's perceptual system was unable to
145 differentiate between the variants of S, the two distributions would be identical. The more
146 sensitive the observer is to the stimulus, the greater is the distance d between the centers of
147 the two distributions. This distance d can therefore be interpreted as the ability of the
148 observer's perceptual system to differentiate between the two kinds of S. The SDT model can
149 be extended to include subjective reports by assuming that task responses and ratings are
150 considered to form an ordered set of responses such as "I'm sure it's A", "I guess A", "I
151 guess B", "I'm sure it's B". The different response options are delineated by a series of
152 criteria, c_1, c_2, \dots, c_n . Participants select one response out of the set of responses by comparing
153 the sensory evidence against the set of criteria. For example, they respond "I guess A" when the
154 sensory evidence falls between that criterion separating the response "I'm sure it's A" from
155 "I guess A" and that criterion separating "I guess A" from "I guess B". An important
156 implication of the SDT model is that the full evidence available to task responses is also
157 available to subjective reports (Macmillan & Creelman, 2005), which is why the model is
158 sometimes referred to as "ideal observer model" (Barrett et al., 2013). While the SDT model
159 has been successfully applied to a vast number of different experiments over the last decades
160 (Macmillan & Creelman, 2005; Wickens, 2002), in recent years, more recent experiments
161 both found support for the SDT model (Peters & Lau, 2015), but also situations where
162 subjective reports were not as optimal as expected from SDT (Maniscalco & Lau, 2012,
163 2016) or even better than expected (e.g. Rausch & Zehetleitner, 2016).

164 1.2.1 Hierarchical model

165 Measuring metacognition only makes sense in models where metacognition is not necessarily
166 perfect. For the purpose of the present analysis, we consider two models of how the SDT
167 model can be extended to account for imperfect metacognition, the *hierarchical model* and
168 the *independent model*. In both of these two models, the task response is selected by a
169 comparison between sensory evidence and the task criterion, just as in the SDT model. A
170 summary of all free parameters of the two models is found in Table 1.

171 TABLE 1 ABOUT HERE

172 The hierarchical model assumes that the sensory evidence involved in selecting the task
173 response is also involved in selecting a rating category (see Fig. 2a). However, in contrast to
174 standard SDT, it is not assumed that the sensory evidence involved in the task response
175 completely determines the subjective report as well. Instead, the sensory evidence is read out
176 by metacognitive processes, whereas the read-out can be incomplete or distorted (Maniscalco
177 & Lau, 2016). One way to express this mathematically is by assuming that the sensory
178 evidence is overlaid by random additive noise, characterized by its standard deviation, σ (cf.
179 Maniscalco & Lau, 2016). The additive noise σ can be interpreted as the degree of distortion
180 between the task process and the metacognitive processes. When the additive noise is absent,
181 the model is identical to the SDT model. For simplicity, we also assume that two rating criteria
182 are placed symmetrically around the primary task criterion θ . The distances of the rating
183 criteria to the task criterion are controlled by the same parameter τ . Conceptually, the spread
184 of rating criteria τ expresses whether subjective reports are made more liberally or more
185 conservatively: When τ is large, the rating criteria are further away from the center of the

186 distribution of sensory evidence. Thus, it is rather unlikely that the sensory evidence is more
187 extreme than the rating criteria; therefore, participants will not often report high confidence.

188 1.2.2 *Independent model*

189 The present study will propose that logistic regression is closely related to another model of
190 metacognition, the *independent model*. The independent model is a new variant of so-called
191 dual channel models (cf. Cul, Dehaene, Reyes, Bravo, & Slachevsky, 2009; Maniscalco &
192 Lau, 2016; Rahnev, Maniscalco, Luber, Lau, & Lisanby, 2012). Dual channel models assume
193 that the rating decision does not depend on the sensory evidence considered for the task
194 decision. Instead, rating decisions are based on a second sample of evidence (see Fig. 2b). In
195 contrast to previous flavors of the dual channel model, the independent model discussed here
196 assumes that there is no interaction between the two samples of evidence except that
197 observers know which of two the task responses they choose. As a consequence, they respond
198 high confidence or high visibility when the evidence sampled in parallel confirms the
199 response option they decide for. It is assumed that the distributions from which the parallel
200 samples of evidence are drawn are characterized by the same shape as those distributions of
201 evidence involved in the primary task, except that the distance between the two distributions
202 may deviate from the distance of distributions in the primary task and is denoted by the rating
203 sensitivity parameter d_m . Conceptually, the rating sensitivity parameter d_m can be interpreted
204 as the amount of evidence available to metacognitive processes that predict whether a
205 decision is correct. Again, we assume that rating criteria are placed symmetrically around the
206 primary task criterion θ , with the distances of the rating criteria to the task criterion controlled
207 by the same parameter τ . The independent model is able to accommodate patterns of data that
208 the hierarchical model struggles to explain: First, the hierarchical model cannot account for
209 participants successfully detecting their own errors (Yeung & Summerfield, 2012). Second,
210 the independent model is able to account for blind insight, e.g. cases when participants
211 perform at chance, but their confidence responses are able to differentiate between correct
212 and incorrect trials (Scott, Dienes, Barrett, Bor, & Seth, 2014). While the hierarchical model
213 achieved better fits to the data than several variants of dual channel models in a metacontrast
214 masking task (Maniscalco & Lau, 2016), the independent model has never been formally
215 assessed with empirical data.

216 1.2.3 *Rationale of the present study*

217 The present analysis was performed to investigate the eligibility of logistic regression as a
218 measure of metacognition. For this purpose, we computed analytically whether logistic
219 regression slopes depend on parameters conceptually associated with metacognition, i.e. the
220 internal noise σ in the hierarchical model and the distance between distributions d_m in the
221 independent model. To test whether logistic regression slopes is biased by task and rating
222 criteria, we varied task bias θ and the spread of rating criteria τ in both models, and
223 investigated the effects on logistic regression slopes. If logistic regression slopes were
224 suitable measures of metacognitive sensitivity, they should be associated with internal noise
225 in the hierarchical model and rating sensitivity in the independent model, and also be
226 independent from task bias and report criteria. To examine the generality of these effects, we
227 also varied the shape of the distributions of evidence generated by the two stimuli. In
228 addition, we varied the link functions, i.e. the transformations to relate subjective reports and
229 task performance. In addition, we performed an analogous analysis to investigate if rating
230 criteria can be assessed by regression intercepts. Finally, data obtained in a low-contrast
231 orientation discrimination task (Rausch & Zehetleitner, 2016) was reanalyzed to investigate if
232 it is possible to differentiate between the hierarchical and the independent model empirically.

233 **2 Logistic regression in a hierarchical model of metacognition**

234 All analyses were conducted using the free software R (R Core Team, 2014). The analysis
 235 code and all reported results are freely available at the Open Science Framework
 236 (<https://osf.io/72aqe/>) to facilitate reproduction of the present study and replication of its
 237 results (Ince, Hatton, & Graham-Cumming, 2012; Morin et al., 2012).

238 **2.1 Calculation of GLM slopes**

239 Analytical closed-form solutions exist for the coefficients of logistic regression when there is
 240 one categorical predictor (Lipovetsky, 2015), but this approach generalizes to other GLMs as
 241 well. In case of metacognition, the GLM is given by the formula

$$g(P(T = 1|C)) = a + b * C \quad (2)$$

242 with g denoting the link function, which describes the transformation used to relate accuracy
 243 and predictors, $P(T = 1|C)$ the probability of being correct in the primary task conditioned
 244 on participant's subjective report, a as intercept, b as slope, and $C \in \{0, 1\}$ as subjective
 245 report. The intercept is

$$a = g(P(T = 1|C = 0)) \quad (3)$$

246 and the slope, which is indicative of metacognitive sensitivity, is given by

$$b = g(P(T = 1|C = 1)) - g(P(T = 1|C = 0)) \quad (4)$$

247 Consequently, GLM slopes can be computed analytically in all situations when
 248 $P(T = 1|C = 0)$ and $P(T = 1|C = 1)$ are known.

249 **2.2 Model description**

250 A graphical model of the hierarchical model is found in Fig. 3. We assume that in each trial
 251 of the experiment, participants are presented with one out of two manifestations of the
 252 stimulus $S \in \{0, 1\}$ which both occur with the same probability. Participants select a response
 253 $R \in \{0, 1\}$ which of the two stimuli was presented. Accuracy of the response T is 1 when $S =$
 254 R , and $T = 0$ otherwise. In each trial, participants select a response based on a single sample
 255 of sensory evidence x . The sensory evidence x is a random sample out of a distribution that
 256 depends on the stimulus. The two distributions corresponding to the two stimulus variants
 257 have the same standard deviation of 1. However, their location depends on the sensitivity
 258 parameter d : When $S = 0$, the mean of the distribution is $-0.5 d$, and when $S = 1$, then the
 259 mean of the distribution of evidence is $0.5 d$. Participants' response R is 0 when x is smaller
 260 than the primary task criterion θ , and 1 otherwise. The sensory evidence x is overlaid by
 261 internal noise, which is randomly sampled out of a distribution with a mean of 0 and a
 262 standard deviation determined by the rating noise parameter σ . The sum of x and internal
 263 noise is called decision variable z and determines subjective reports by a comparison with
 264 rating criteria c_0 and c_1 : When $R = 1$ and the z is greater than the rating criterion c_1 , then the
 265 subjective report C is 1. When $R = 1$ and the z is smaller than the rating criterion c_1 , then the
 266 subjective report C is 0. Likewise, when $R = 0$, then C is 1 if z is smaller than the rating
 267 criterion c_0 , and C is 0 if z is greater than the rating criterion c_0 . For simplicity, we assume
 268 that the distance between θ and c_0 is the same as the distance between θ and c_1 . This distance
 269 is controlled by the parameter τ , which reflects the conservativeness of rating criteria. The

270 formulae for computing $P(T = 1|C = 0)$ and $P(T = 1|C = 1)$ in the hierarchical model can
271 be found in the Appendix A.

272 We considered two different distributions of evidence, the Gaussian distribution and the
273 logistic distribution. The Gaussian distribution is often motivated by the central limits
274 theorem and the averaging of many events (DeCarlo, 1998). The logistic distribution can be
275 motivated from Choice Theory (Luce & Suppes, 1965; Macmillan & Creelman, 2005). For
276 most SDT applications, results obtained based on logistic and Gaussian distributions are very
277 similar (DeCarlo, 1998; Wickens, 2002).

278 **2.3 Link functions**

279 The link function that transforms the probability of being correct to the logarithm of the odds
280 of being correct to being incorrect is called the logit link and is the defining feature of logistic
281 regression. In the present analysis, three different link functions are examined:

- 282 1) the logit link
- 283 2) the probit link
- 284 3) the half-logit link.

285 The logit link was the default choice in previous studies. The probit link was included into
286 the analysis because SDT models can be seen as a subclass of generalized linear regression
287 models where the inverse link function corresponds to the cumulative distribution function of
288 the evidence (Brockhoff & Christensen, 2010; DeCarlo, 1998). Logistic regression assumes a
289 logistic distribution of errors, while probit regression assumes a standard normal distribution
290 of errors. Consequently, it appears necessary to examine if logistic regression is only valid in
291 logistic SDT models, and probit regression in Gaussian SDT models. The adjusted link was
292 proposed as an adjustment in tasks with a finite guessing probability (Brockhoff & Müller,
293 1997). The use of the logit and probit transform implies that the probability of being correct
294 varies between 0 and 1; however, in binary tasks, the probability of being correct is bounded
295 by the guessing probability 0.5 (Rausch et al., 2015). To account for the guessing probability,
296 a link function can be used to ensure that the transformed accuracy is free to vary in the full
297 range between $-\infty$ and ∞ . In tasks with two choices, this can be achieved by the adjusted link
298 function $g(x) = \log((x - 0.5)/(1 - x))$. This function was referred to as half-logit link (cf.
299 Williams, Ramaswamy, & Oulhaj, 2006).

300 **2.4 Results**

301 To investigate if GLM slopes are sensitive to metacognition and unbiased by primary task
302 criterion and by rating criteria assuming the hierarchical model, we calculated logistic
303 regression slopes as a function of internal noise σ as well as primary task criterion θ (Fig. 4)
304 and spread of rating criteria τ (Fig. 5). These calculations were repeated using Gaussian and
305 logistic distributions of evidence and with logit, probit, and half-logit link functions.

306 *2.4.1 Are GLM slopes sensitive to metacognition?*

307 In Fig. 4 and 5, separate lines indicate different degrees of internal noise σ . Greater amounts
308 of internal noise are associated with lower regression slopes at each level of primary task
309 criterion (Fig. 4), and at each level of rating criteria spread (Fig. 5). This pattern holds
310 independently from distributions of evidence and link functions (separate panels of Fig. 4 and

311 5). When the amounts of noise are extreme, i.e. when metacognition is effectively absent,
312 regression slopes also tend towards zero. Overall, GLM analysis is sensitive to metacognition
313 according to the hierarchical model.

314 2.4.2 *Do GLM slopes depend on the primary task criterion?*

315 As can be seen from each panel of Fig. 4, regression slopes depend heavily on the primary
316 task criterion according to the hierarchical model. The precise form of the relationship
317 between primary task criterion and slopes depends on a complex interaction between the
318 amount internal noise, the shape of the distributions of evidence, and the link functions.
319 When the distributions are logistic, when the half-logit transform is used, or when the internal
320 noise is not small, logistic regression slopes increase monotonously with primary task
321 criterion: Consequently, greater regression slopes are not necessarily due to metacognition,
322 but could also be due to a stronger bias towards one of the task alternatives. However, when
323 Gaussian distributions are assumed and the amount of noise is small, the relationship between
324 regression slopes and primary task criterion is u-shaped: Slopes are maximal when observers
325 are either not biased at all or extremely biased towards one of the task alternatives. Therefore,
326 a primary task criterion may not only increase, but also decrease regression slopes. Overall,
327 regression slopes depend on the primary task criterion, but the direction and magnitude of the
328 effect is strongly dependent on the other model parameters of the hierarchical model.

329 FIG.4 ABOUT HERE

330 2.4.3 *Do GLM slopes depend on rating criteria?*

331 As can be seen from Fig. 5, logistic and probit regression slopes increase with rating criteria
332 spread τ , i.e. when more conservative rating criteria are set (top and central panels). The
333 relationship between half-logit regression slopes and rating criteria can be u-shaped or
334 decreasing (bottom panels). However, the relationships between slopes and rating criteria are
335 moderated by internal noise as well the shape of the distributions: When the distributions are
336 Gaussian, the effect imposed by rating criteria will be smaller the more internal noise is
337 superimposed on the sensory evidence. When the distributions are logistic, the effect imposed
338 by rating criteria is maximal at medium level of internal noise. Overall, according to
339 hierarchical models, differences between regression slopes can not only be caused by
340 metacognition and task criterion, but also by the way participants set rating criteria. Again,
341 the effect is strongly dependent on the other model parameters of the hierarchical model.

342 FIG. 5 ABOUT HERE

343 2.4.4 *Can rating criteria be assessed by regression intercepts?*

344 To investigate if regression intercepts are sensitive to rating criteria and independent from
345 metacognition and primary task criterion in the hierarchical model, we calculated intercepts
346 as a function of internal noise σ , primary task criterion θ , rating criteria spread τ , shape of the
347 distributions, and link functions. Consistent across distributions and link functions, intercepts
348 were negatively related to the primary task criterion θ and positively correlated to the internal
349 noise σ (see Supplementary Fig. 1). However, intercepts were only sensitive to the
350 spread of rating criteria τ when the amount of internal noise was low (see Supplementary Fig.
351 2). This means that intercept effects cannot uniquely be attributed to rating criteria, but also
352 to metacognition or due to primary task criterion. Even more, when rating criteria are
353 different between two conditions, it will not be possible to detect the effect using intercepts
354 when the amount of internal noise is high.

355 3 Logistic regression in the independent model of metacognition

356 3.1 Calculation of GLM slopes and link functions

357 Calculation of GLM slopes and link functions were identical to the hierarchical model.

358 3.2 Model description

359 The independent model can be expressed by the graphical model in Fig. 6. It is analogous to
360 the hierarchical model with the following differences: Participants select only the primary
361 task response based on the sensory evidence x . For subjective reports, a second sample of
362 sensory evidence y is created, which is stochastically independent from x . The standard
363 deviation of the distribution of y is assumed to be 1. The location of the distribution of y
364 depends on the stimulus in the current trial as well as on the rating sensitivity parameter d_m :
365 When $S = 0$, the mean of the distribution is $-0.5 d_m$, and when $S = 1$, then the mean of the
366 distribution of evidence is $0.5 d_m$. When $R = 1$ and y is greater than the rating criterion c_1 ,
367 then the participant's subjective report C is 1. When $R = 1$ and y is smaller than the rating
368 criterion c_1 , then the subjective report C is 0. Likewise, when $R = 0$, the C is 1 if y is smaller
369 than the rating criterion c_0 , and C is 0 if y is greater than the rating criterion c_0 . Again, it is
370 assumed for simplicity that the distance between θ and c_0 is the same as the distance between
371 θ and c_1 , controlled by the parameter τ reflecting the conservativeness of rating criteria. The
372 formulae for computing $P(T = 1|C = 0)$ and $P(T = 1|C = 1)$ in the independent model can
373 be found in the Appendix B.

374 3.3 Results

375 To investigate if GLM slopes are sensitive to metacognition and unbiased by the primary task
376 criterion and by rating criteria according to the independent model, slopes were calculated as
377 a function of rating sensitivity d_m as well as primary task criterion θ (Fig. 7) and spread of
378 rating criteria τ (Fig. 8). Again, these calculations were performed using Gaussian and
379 logistic distributions of evidence and using logit, probit, and half-logit link functions.

380 3.3.1 *Are GLM slopes sensitive to metacognition?*

381 In Fig. 7 and 8, separate lines indicate different degrees of rating sensitivity d_m . Greater rating
382 sensitivity was associated with increasing regression slopes at each level of primary task
383 criterion (Fig. 7), and at each level of rating criteria spread (Fig. 8). This pattern holds across
384 the different distributions of evidence and link functions (separate panels of Fig. 7 and 8).
385 When rating sensitivity is 0, i.e. when the sensory evidence available to metacognition does
386 not differentiate between the two stimuli, regression slopes become 0 when observers are
387 unbiased towards the two response options. However, when rating sensitivity is 0 and when
388 there is a bias, the slope will become negative. Overall, these results indicate GLM analysis is
389 sensitive to metacognition in the independent model, but the sign of the slope should not be
390 interpreted without consideration of the primary task criterion.

391 3.3.2 *Do GLM slopes depend on the primary task criterion?*

392 As can be seen from each panel of Fig. 7, slopes are influenced by the primary task criterion
393 θ according to the independent model as well. While there is always an effect of primary task

394 criterion, direction and magnitude of the effect depends on a complex interaction between the
395 amount of metacognition as indexed by the rating sensitivity d_m , the shape of the distributions
396 of evidence, as well as the link functions. When the distributions of evidence are Gaussian
397 and when a probit or logit link function is used, the relationship between primary task
398 criterion θ and slopes appears to be u-shaped and similar across different levels of
399 metacognition: Slopes reach a minimum at medium primary task criterion θ , and increase
400 when θ is either 0 or maximal (see Fig. 7 upper and central panel to the left). When the
401 distributions of evidence are logistic and when a probit or logit link functions is applied,
402 GLM slopes decrease with increasing θ (see Fig. 7 upper and central panel to the right).
403 When the half-logit link function is used, the slopes increase exponentially with θ for
404 medium-to-large rating sensitivities. However, when rating sensitivity is low, slopes decrease
405 with increasing θ . Overall, these observations indicate that slopes as measures of
406 metacognition in the independent model can be biased by the primary task criterion; the kind
407 of bias however depends on the other model parameters as well as on the choice of the link
408 function.

409 FIG.7 ABOUT HERE

410 3.3.3 Do GLM slopes depend on rating criteria?

411 Fig. 8 shows that GLM slopes are independent from the spread of rating criteria τ in one case:
412 When the evidence is assumed to be logistically distributed, logistic regression slopes are
413 independent from rating criteria (see Fig. 8 upper right panel). Probit regression slopes
414 decrease when more conservative rating criteria are used. In contrast, when the evidence is
415 assumed to be Gaussian, both logistic and probit regression slopes increase when more
416 conservative rating criteria are applied (see Fig. 8 upper and central panel to the left). These
417 effects are moderated by the amount of evidence available to ratings, i.e. rating sensitivity d_m .
418 The larger d_m is, the more pronounced is the effect of rating criteria on GLM slopes. For the
419 half-logit link function, slopes increase massively when the spread of rating criteria is very
420 small (see Fig. 8 lower row). In summary, logistic regression slopes are unbiased by the
421 spread of the rating criteria in the independent model when the distributions of evidence are
422 logistic. When other link functions and Gaussian distributions are assumed, slopes can be
423 heavily influenced by rating criteria.

424 FIG. 8 ABOUT HERE

425 3.3.4 Can rating criteria be assessed by regression intercepts?

426 To investigate if regression intercepts are sensitive to rating criteria and independent from
427 metacognition and primary task criterion in the independent model, we calculated intercepts
428 as a function of rating sensitivity d_m , primary task criterion θ , rating criteria spread τ , shape of
429 the distributions, and link functions. Consistent across distributions and link functions,
430 intercepts were negatively associated with the primary task criterion θ and positively
431 associated with the rating sensitivity d_m (see Supplementary Fig. 3). The intercepts were only
432 sensitive to the spread of rating criteria τ when d_m was above zero (see Supplementary Fig.
433 4). Overall, this means that according to the independent model – just as the hierarchical
434 model - intercept effects could be due to rating criteria, degree of metacognition, and primary
435 task criterion. In addition, true effects on rating criteria will remain undetected when
436 metacognition is low.

437 **4 Model fits of the independent and the hierarchical model in a low-contrast** 438 **orientation discrimination task**

439 The present analysis implies that logistic regression slopes are unbiased by rating criteria
 440 according to only one specific model of metacognition, namely the logistic independent
 441 model. In all models, the slopes are dependent on primary task criteria. Consequently, if
 442 researchers intend to use logistic regression as measure of metacognition, it would be useful
 443 to identify the cognitive model underlying the rating data. While measures of primary task
 444 criteria are readily available from SDT (Green & Swets, 1966; Macmillan & Creelman, 2005;
 445 Wickens, 2002), it might be a challenge to differentiate the independent model and logistic
 446 distributions from the hierarchical model and Gaussians. We reanalyze confidence ratings
 447 obtained in a recent low-contrast orientation discrimination experiment (Rausch &
 448 Zehetleitner, 2016) to investigate if this is possible using cognitive modeling and the
 449 maximum likelihood procedure.

450 **4.1 Reanalysis**

451 *4.1.1 Experimental task*

452 20 participants, all of which provided written informed consent, performed one training block
 453 and nine experimental blocks of 42 trials each of a low contrast orientation discrimination
 454 task. First, participants were presented with a fixation cross for 1 s. Then, the target stimulus,
 455 a binary grating oriented either horizontally or vertically, was presented for 200 ms with
 456 varying contrast levels of 0, 2.2, 3.9, 5.0, 5.5, and 6.9%. The screen remained blank
 457 afterwards until participants made a non-speeded discrimination response by key press
 458 whether the target had been horizontal or vertical. After each discrimination response,
 459 participants made two subjective reports, one regarding their visual experience of the
 460 stimulus, and one regarding their confidence in being correct in the discrimination task. For
 461 that, each question was displayed on the screen, which was: “How clearly did you see the
 462 grating?” or “How confident are you that your response was correct?” The sequence of
 463 questions was balanced across participants. Participants delivered subjective reports on a
 464 visual analog scale using a joystick, which means that participants selected a position along a
 465 continuous line between two end points by moving a cursor. The end points were labeled as
 466 “unclear” and “clear” for the experience scale and “unconfident” and “confident” for the
 467 confidence scale, i.e. observers indicated their experience or confidence by the selected
 468 cursor position on the continuous scale. If the discrimination response was erroneous, the trial
 469 ended by displaying the word “error” for 1 s on the monitor. There was no feedback with
 470 respect to the subjective report. Please refer to Rausch and Zehetleitner (2016) for a more
 471 detailed description of the experiment.

472 *4.1.2 Models*

473 We fitted eight different models to the data, which were characterized by all possible
 474 combinations of the following three features:

- 475 (i) The model could be either a hierarchical or an independent model as outlined
 476 above.
- 477 (ii) The distributions of evidence and noise could be either Gaussian or logistic.
- 478 (iii) The primary task criterion θ was either treated as a free parameter or fixed at
 479 0.

480 In all eight models, we assumed that the discrimination sensitivity d varied across contrast
 481 levels, while the other parameters were assumed to be constant across contrast levels. Thus,
 482 each model involved six different sensitivity parameters $d_1 - d_6$, one for each contrast level.
 483 Moreover, each model involved a series of 11 rating criteria spread parameters $\tau_1 - \tau_{11}$. These
 484 parameters described how close the rating criteria were located to the primary task criterion
 485 θ : τ_1 denotes the location of the closest pair of criteria at both sides of θ ; τ_2 referred to the
 486 second pair, and so on. In hierarchical models, the degree of metacognition was denoted by
 487 the internal noise parameter σ . For independent models, the rating sensitivity d_m was assumed
 488 to be a constant fraction of the discrimination sensitivity d , denoted by the parameter a .
 489 Overall, the models had 18 or 19 free parameters depending on if the primary task criterion θ
 490 was fixed at zero.

491 4.1.3 Model fitting

492 Model fitting was performed separately for each single participant. The fitting procedure
 493 involved the following computational steps. First, the continuous confidence ratings were
 494 discretized by dividing the continuous scale into equal 12 partitions. Analyses using four or
 495 eight bins gave similar results when used on the empirical data; however, 12 bins improved
 496 the recovery of model-generated data (see 4.1.6), which is why results based on 12 bins are
 497 reported. Second, we computed the frequency of each rating bin given orientation of the
 498 stimulus and the orientation response. Third, for each of the 8 models, the set of parameters
 499 was determined that maximized the likelihood of the data. To compute the likelihood, we
 500 made two widespread assumptions in SDT modeling (Dorfman & Alf, 1969; Maniscalco &
 501 Lau, 2016): (i) responses in each trial were assumed to be independent from each other and
 502 (ii) the joint probability of a task response and a subjective report given the stimulus was
 503 constant across trials. Formally, the likelihood of a set of parameters given primary task
 504 responses and subjective reports $\mathcal{L}(p|R, C)$ is given by

$$\mathcal{L}(p|R, C) \propto \prod_{i,j,k} P(R_i, C_j|S_k, p)^{n(R_i, C_j|S_k)} \quad (5)$$

505 where $P(R_i, C_j|S_k, p)$ denotes the probability of a primary task response in conjunction with a
 506 specific subjective report given the stimulus and the set of parameters, and $n(R_i, C_j|S_k)$
 507 indicates the frequency how often the participants gave a specific response in conjunction
 508 with a specific subjective report given the stimulus. The set of parameters with the maximum
 509 likelihood was determined by minimizing the negative log likelihood as the latter is
 510 computationally more stable:

$$\log(\mathcal{L}(p|R, C)) \propto \sum_{i,j,k} \log(P(R_i, C_j|S_k, p)) \times n(R_i, C_j|S_k) \quad (6)$$

511 To minimize the negative log likelihood, we used a general SIMPLEX minimization routine
 512 implemented in the R function `optim` (Nelder & Mead, 1965). The formulae for
 513 $P(R_i, C_j|S_k, p)$ are found in the Appendix as formulae (A6) – (A13) and (B5) – (B12).
 514 Internal noise was parametrized as the log of the standard deviation of the noise to allow
 515 negative values of the parameter during the fitting process. To maintain a fixed sequence of
 516 rating criteria, we did not directly fit $\tau_1 - \tau_{11}$, instead, optimization was performed on the log
 517 distance to the nearest rating criterion closer to the primary task criterion.

518 4.1.4 Model selection

519 Following the fitting procedure, we assessed the relative quality of the eight candidate models
 520 using the Bayes information criterion (BIC, Schwarz, 1978) and the Akaike information
 521 criterion (AIC, Akaike, 1974). Conceptually, the BIC measures the degree of belief that a
 522 certain model is the true data-generating model relative to the other models under
 523 comparison, assuming that the true generative model is among the set of candidate models. In
 524 contrast, the AIC measures the loss of information when the true generative model is
 525 approximated by the candidate model. We used AIC_c , a variant of AIC that corrects for finite
 526 sample sizes (Burnham & Anderson, 2002). BIC and AIC_c take into account descriptive
 527 accuracy (i.e. goodness of fit) and parsimony (i.e. smallest number of parameters), but the
 528 BIC favors parsimony more heavily than the AIC_c does. BIC and AIC_c are given by the
 529 following formulae:

$$BIC = -2 \log(\mathcal{L}(p|R, C)) + k \log(n) \quad (7)$$

$$AIC_c = -2 \log(\mathcal{L}(p|R, C)) + 2k + \left(\frac{2k(k+1)}{(n-k-1)} \right) \quad (8)$$

530 where k indicates the number of parameters and n the number of observations.

531 4.1.5 Statistical testing

532 In experiments with a standard number of trials, even if the independent logistic model with
 533 fixed task criterion was true, it cannot be expected that models can be correctly identified for
 534 each single participant. However, in this case, it would still be expected that independent
 535 models obtained the best fit more frequently than hierarchical models, that models based on
 536 the logistic distribution would achieve the best fit more often than Gaussian models, and
 537 likewise that models with the free primary task criterion fixed at 0 would obtain better fits
 538 than models with a free primary task criterion.

539 Therefore, we determined for each participant which of the eight candidate models achieved
 540 the minimal BIC and minimal AIC_c . Then, we performed three tests if those model features
 541 that imply that logistic regression slopes are independent from criteria are more likely to
 542 result in the best BIC or AIC_c : First, we tested if independent models were more likely to
 543 achieve the best BIC or AIC_c than hierarchical models. Second, we assessed if models with
 544 the primary task criterion fixed at 0 achieved the best BIC/ AIC_c with a greater probability
 545 than models with the primary task criterion as free parameter. Finally, we examined if models
 546 based on logistic distributions attained minimal BIC and AIC_c more frequently than models
 547 based on Gaussians.

548 Statistical testing was based on Bayes factors for proportions implemented in the R library
 549 BayesFactor (Morey & Rouder, 2015). Bayes factors provide continuous measures of how
 550 the evidence supports the alternative hypothesis over the null hypothesis and vice versa
 551 (Dienes, 2011; Rouder, Speckman, Sun, Morey, & Iverson, 2009). Specifically, the Bayes
 552 factor indicates how the prior odds about alternative hypothesis and null hypothesis need to
 553 be multiplied to obtain the posterior odds of the two hypotheses. As null hypothesis, we
 554 assumed that both variants of the model were equally likely to achieve the best fit. As one-
 555 sided alternative hypothesis, we assumed a logistic distribution of the logits of the probability
 556 around 0 over the interval 0.5 and 1 with a scale parameter of 0.5.

557 4.1.6 Recovery of model generated data

558 To investigate the reliability of the model selection process, we generated random data sets
 559 based on each participant's parameter set obtained during the fitting process (analogous

560 procedure to Maniscalco & Lau, 2016). We used only the parameter sets of two models: The
561 independent logistic model with primary task criterion fixed at 0 was used because it implies
562 that slopes are independent of confounds by criteria and thus logistic regression slopes can be
563 used as measure of metacognition without concern. The alternative model was the
564 hierarchical Gaussian model with free primary task criterion because it is maximally different
565 from the independent logistic model with primary task criterion fixed at 0. This procedure
566 yielded a mock replica of each participant's behavioral data. We then fitted all eight models
567 to these simulated data and performed model selection using AICc and BIC. If the model
568 selection methodology is reliable, then independent, logistic and fixed θ models should be
569 selected when the data is generated according to the independent logistic model with primary
570 task criterion fixed at 0. Likewise, when the hierarchical Gaussian model with free primary
571 task criterion is used to create the data, hierarchical, Gaussian models with a free primary
572 task parameter should be preferred during model selection. We replicated the analysis using
573 100, 200, 378 (the same number of trials as in the real experiment), 600, 1200, 2500, 5000,
574 10000, 20000 and 50000 trials to estimate the number of trials required to perform a reliable
575 model selection.

576 4.2 Results

577 The fitting procedure converged for all participants except for one participant, where the
578 fitting of the hierarchical Gaussian model did not converge. On average, the best model fit
579 was obtained by the hierarchical logistic model with fixed primary task criterion both in
580 terms of BIC ($M = 1432.8$) as well as AICc ($M = 1363.7$). According to the BIC, the strength
581 of evidence in favor the hierarchical logistic model with fixed primary task criterion as
582 indexed by the difference in BIC was always substantial (all other models: M 's ≥ 1437.7). In
583 contrast, according to the AICc, there was only a small benefit of the hierarchical logistic
584 model with fixed primary task criterion compared to the independent logistic model with free
585 task criterion parameter ($M = 1364.5$) and the hierarchical logistic model including a free task
586 criterion parameter ($M = 1364.9$). As can be seen from Fig. 9, both the hierarchical logistic
587 model with task criterion fixed at zero and the hierarchical logistic model including a free
588 task criterion parameter provided qualitatively good accounts of the distributions of
589 discrimination accuracy and confidence ratings.

590 4.2.1 *Is there support for the independent model?*

591 Independent models of metacognition only provided the best model fit in 30% of the
592 participants according to the BIC, and in 40% of the participants according to the AICc. The
593 Bayes factor analysis indicated evidence against the hypothesis that independent models are
594 more likely to attain the best fit both in terms of BIC, $BF_{10} = 0.19$, as well as in terms of
595 AICc, $BF_{10} = 0.29$. This result implies that prior beliefs that a flavor of the independent
596 models is the generative model of the present data should be attenuated. Accordingly,
597 estimates of metacognition in the present data set based on GLM slopes are likely to be
598 affected by rating criteria.

599 4.2.2 *Can the primary task criterion be fixed at 0?*

600 Models with the primary task criterion fixed at zero attained the best model fit in 60% of the
601 participants according to the BIC, and in 40% according to the AICc. The Bayes factor
602 analysis does not provide any support in favor of the hypothesis that models with the primary
603 task criterion fixed at zero are associated with better BICs, $BF_{10} = 1.02$. However, there was

604 evidence against the hypothesis that models with the primary task criterion fixed at zero are
605 associated with better AIC_c s, $BF_{10} = 0.29$. Consequently, there is some indication that GLM
606 slopes in the present data would also be influenced by the primary task criterion, but the
607 evidence is not consistent.

608 4.2.3 *Is the evidence distributed logistically?*

609 Models based on the assumption of logistic distributions achieved the minimal BIC in 85% of
610 the participants according to BIC and even in 90% according to AIC_c . The Bayes factors
611 indicated strong evidence that the logistic models were more likely to produce better fits than
612 Gaussian models, both in terms of BIC, $BF_{10} = 56.1$, as well with regard to AIC_c , $BF_{10} =$
613 231.6 . This means that only one out of the three conditions for logistic regression slopes to be
614 independent of criteria was met in the present data set.

615 4.2.4 *Can the model underlying simulated data be recovered?*

616 Fig. 10 shows the results of the model recovery analysis of data simulated according to the
617 independent logistic model with fixed task criterion (squares) and data generated according to
618 the hierarchical Gaussian model with a free task criterion (circles).

619 As can be seen from the panels on the left, when the data was simulated according to the
620 independent logistic model with fixed task criterion, the model was nearly always correctly
621 identified as being one of the independent models. However, when the true model was the
622 hierarchical Gaussian model with free task criterion, the model was relatively often
623 misclassified as independent: When the trial number was small ($N \leq 200$), model recovery
624 was even below chance. 5000 trials were required to obtain a tolerable recovery rate of
625 approximately 70%. Increasing the trial number even more did not substantially improve
626 model recovery.

627 The central panels of Fig. 10 show model recovery with respect to the free primary task
628 criterion vs. the primary task criterion fixed at 0. For the BIC as model selection criterion,
629 5000 trials were necessary to detect the free task criterion parameter with a tolerable accuracy
630 of 70%. In contrast, for the AIC_c , 1200 trials were sufficient to reach 70%. Notably, the AIC_c
631 does not favor models with smaller number of parameter as heavily as the BIC does. It can
632 also be seen that model recovery accuracy decreased with trial number for the independent
633 logistic model with fixed task criterion (dotted lines). However, at the same time, model
634 recovery accuracy increased with sample size for the hierarchical Gaussian model with free
635 task criterion (straight lines). An explanation for this pattern may be that AIC_c and BIC both
636 favor more parsimonious models. For small sample sizes, both AIC_c and BIC prefer models
637 with a smaller number of parameters, which is why models are selected where the primary
638 task criterion is fixed at 0. The bias towards smaller results will result in a correct
639 classification with respect to the primary task criterion when the true model is the
640 independent logistic model with fixed task criterion, but an error will occur when the true
641 model is the hierarchical Gaussian model with free task criterion.

642 Finally, model recovery with respect to logistic and Gaussian distributions is depicted in the
643 right panels of Fig. 10. For large sample sizes, logistic distributions were more often correctly
644 identified than Gaussians. However, 1200 trials were sufficient to obtain a tolerable recovery
645 rate of approximately 70% for both distributions.

646 5 Discussion

647 The analysis presented here investigated whether GLM slopes as measures of metacognition
648 are biased by the spread of rating criteria and the primary task criterion. We showed
649 analytically that logistic regression slopes are independent from rating criteria only according
650 to one specific model of metacognition: the independent model based on logistic
651 distributions. When other distributions were assumed, when other link functions were used,
652 or when a hierarchical model was adopted, regression slopes always depended on the spread
653 of rating criteria. The direction and magnitude of these effects depended on the other model
654 parameters. The primary task criterion was related to regression slopes in all considered
655 models. Depending on the model parameters, the relationship between slopes and task
656 criterion were increasing, decreasing, or even u-shaped. An analysis of regression intercepts
657 revealed that intercepts were insensitive to rating criteria when the amount of metacognition
658 was too low. In addition, we examined if these models can be identified empirically on an
659 existing data set, observing that a massive number of trials is required to distinguish between
660 hierarchical and independent models with tolerable accuracy.

661 **5.1 Is logistic regression a biased measure of metacognitive sensitivity?**

662 When the aim of a study is to estimate the degree of metacognitive sensitivity, it is generally
663 accepted that a suitable measure should be independent from the primary task criterion and
664 rating criteria (Barrett et al., 2013). However, the present study revealed that logistic
665 regression slopes depend on the primary task criterion independent of the underlying model
666 of metacognition. Logistic regression slopes also depend on rating criteria except for one
667 specific model of metacognition: When the sensory evidence considered for primary task
668 responses is stochastically independent from the sensory evidence used in rating decisions,
669 and when these two types of evidence both form logistic distributions, then logistic regression
670 slopes are independent from rating criteria. This means that when researchers encounter an
671 empirical effect on logistic regression slopes, without knowledge about the underlying model
672 of metacognition, there will be at least three possible explanations for the effect: (i) the effect
673 can be mediated by participants' degree of metacognition of the processes engaged in
674 performing the task, (ii) the effect might also be due to participants' bias towards one of the
675 task alternatives, and (iii) the effect may also depend entirely on differences how liberal or
676 conservative participants' rating criteria are. Likewise, researchers might be unable to
677 observe real effects on participants' degree of metacognition when there is also a difference
678 between participants' bias or between rating criteria because the effects of metacognition and
679 criteria could balance out.

680 **5.2 Does the present analysis generalize to other models of decision-making and** 681 **metacognition?**

682 The present analysis was based on specific assumptions about the decision process as well as
683 the cognitive model underlying metacognition. Concerning the decision process, we assumed
684 the standard SDT model. Concerning metacognition, only two models, the hierarchical and
685 the independent model were considered. Is it reasonable to assume that the characteristics of
686 GLM slopes outlined here generalize to other models of decision-making and metacognition?

687 As the literature provides a multitude of competing models of decision making, confidence
688 and / or metacognition, it is not feasible to investigate the eligibility of logistic regression for
689 each model proposed in the literature. Examples for different models include the bounded
690 accumulation model (Kiani, Corthell, & Shadlen, 2014), the collapsing confidence boundary

691 model (Moran, Teodorescu, & Usher, 2015), the consensus model (Paz, Insabato, Zylberberg,
692 Deco, & Sigman, 2016), the reaction time account (Ratcliff, 1978), the self-evaluation model
693 (Fleming & Daw, 2016), and two-stage signal detection theory (Pleskac & Busemeyer, 2010).
694 The attempt seems even futile because the number of models in the literature is continuously
695 increasing. However, the two models tested here represent two complementary prototypes of
696 models of metacognition: The hierarchical model is the simplest possible variant of models
697 where the evidence used for the decision between the primary task alternatives also informs
698 the decision between different rating criteria. The view that the evidence used in the decision
699 process is also involved in metacognition is a standard tenet of theories about metacognition.
700 In contrast, the independent model assumes that evidence used for the rating decision is
701 sampled entirely independently from the evidence used for the task response. It is an open
702 question whether datasets exist that can be conveniently described by the independent model.

703 The majority of existing models in the literature are closely related to one of the two models
704 or the models constitute a combination of the two models. For example, post-decisional
705 accumulation models can be seen as a combination of the hierarchical and the independent
706 model, where a second, independent sample of evidence is acquired later in time than the
707 sensory evidence used for performing the task (Moran et al., 2015; Pleskac & Busemeyer,
708 2010). When a model assumes that rating decisions are informed by both evidence considered
709 in the primary task as well as evidence sampled in parallel, it is reasonable to expect that
710 biases apparent in the hierarchical and the independent model persist when the two sources of
711 evidence are combined.

712 Of course, it is still possible that there will be new theories which imply that logistic
713 regression is not affected by primary task criteria and rating criteria. However, the
714 implications of the present study do not require that such a model does not exist. What the
715 study implies is that according to plausible models of metacognition, logistic regression is
716 affected from primary task criterion and rating criteria. As a consequence, researchers who
717 wish to use logistic regression to measure metacognitive sensitivity need to show that their
718 effects cannot be alternatively explained by rating criteria and primary task criteria.

719 **5.3 Can the independent model be empirically identified?**

720 To exclude the possibility that effects on regression slopes are not caused or masked by rating
721 criteria, it would be useful to identify the underlying model of metacognition. Unfortunately,
722 the present model recovery analysis revealed that the amount of trials required to correctly
723 classify a true underlying hierarchical Gaussian model is massive. Moreover, the hierarchical
724 Gaussian model was still occasionally misclassified as independent or logistic even with
725 extreme trial numbers. In a similar comparison between different models of metacognition,
726 two thirds of the participants were excluded to reduce noise in the data and improve model
727 selection, implying that these models are also not trivial to distinguish based on other data
728 sets (Maniscalco & Lau, 2016). Likewise, Gaussian and logistic distributions are known to
729 produce similar results in many applications (DeCarlo, 1998; Wickens, 2002). Consequently,
730 when researchers intend to model the cognitive architecture of metacognition, they will need
731 to ensure that both the sample size and the trial number are sufficiently large to ensure that
732 classification errors are outnumbered by correct classifications. However, for those
733 researchers who are only interested in measuring metacognition, the standard application of
734 logistic regression as measure of metacognitive sensitivity, cognitive modeling will usually
735 not be a feasible option because the number of trials in standard experiments is typically too
736 small. Future studies might be able to provide more efficient methods to distinguish between
737 the hierarchical Gaussian and the independent logistic model.

738 **5.4 What other methods can be used to estimate metacognitive sensitivity?**

739 What are the options to avoid the confound of metacognitive sensitivity with task criteria and
740 rating criteria? There are two options: (i) researchers can resort to measures of metacognitive
741 sensitivity other than GLM slopes, or (ii) they can control for primary task and rating criteria
742 statistically.

743 What are the alternatives to logistic regression slopes (cf. Fleming & Lau, 2014)? A method
744 that recently received a considerable amount of attention is meta-d'. Meta-d' quantifies the
745 degree of metacognition in terms of a standard SDT model where task responses and
746 subjective reports are made based on identical evidence (Maniscalco & Lau, 2012). There are
747 several reasons to use meta-d': Meta-d' is reasonably robust to changes of primary task
748 criteria and rating criteria in an optimal observer SDT model, a decreasing signal SDT model
749 and an increasing signal SDT model (Barrett et al., 2013). The decreasing signal SDT model
750 is closely related to the hierarchical model in the present study, while the increasing signal
751 SDT model can be seen as combination of the hierarchical model and the independent model
752 in the present study. In addition, meta-d' provides control over performance in the primary
753 task. Meta-d' has even the unique advantage of allowing comparisons with primary task
754 performance as metacognitive sensitivity and primary task performance are measured on the
755 same scale (Fleming & Lau, 2014; Maniscalco & Lau, 2012). Thus, meta-d' is able to assess
756 imperfect metacognition as well as metacognition better than expected from task
757 performance.

758 An argument against the use of meta-d' is that meta-d' requires assumptions about the
759 distributions of evidence during the decision process. The most common choice are equal
760 Gaussian distributions (Barrett et al., 2013; Maniscalco & Lau, 2012), but other distributions
761 have been implemented as well (Rausch et al., 2015). However, to our knowledge, no study
762 so far has investigated how often the distributional assumptions of meta-d' are in fact
763 violated, or how sensitive meta-d' is to violations of these distributional assumptions.
764 Moreover, meta-d' has also never been assessed in an independent model of metacognition.

765 SDT approaches that rely on distributional assumptions are often criticized because it is often
766 hard to test of these assumptions are justified. Indeed, in a classical paper, Nelson (1984)
767 recommended gamma correlations to avoid the distributional assumptions made by
768 parametric SDT methods. Although many researchers have since used gamma correlations,
769 simulations suggested that rating criteria strongly impact on gamma correlations, making
770 results obtained by gamma ambiguous as they could be due to metacognition or due to
771 criterion setting (Masson & Rotello, 2009).

772 An non-parametric SDT approach to estimate metacognitive sensitivity is by means of type 2
773 ROC-curves (Fleming et al., 2010). This method provides a measure of metacognitive
774 sensitivity free of bias by rating criteria and distributional assumptions. However, the number
775 of trials required to estimate type 2 ROC-curves can be massive (Nelson, 1984), and type 2
776 ROC curves do not provide any control over performance in the primary task and the
777 associated criteria (Fleming & Lau, 2014).

778 Overall, it appears to us that meta-d' is the most useful measure of metacognitive sensitivity,
779 although more research on the distributional assumptions required by meta-d' would be desirable.
780 Meanwhile, when the underlying distributions of evidence cannot be ascertained, it may be
781 useful to check if meta-d' and type 2 ROC-curve converge to the same results.

782 **5.5 How can criteria be controlled?**

783 When researchers do not wish to resort to other methods than logistic regression slopes, they
784 should at least control primary task criteria and rating criteria statistically. This approach
785 requires of course the assessment of criteria independent of metacognitive sensitivity. As
786 logistic regression will often be used when researchers do not wish to make explicit
787 assumptions about the underlying model of metacognition, the measure of criteria should also
788 be model-free.

789 Logistic regression intercepts, which are occasionally used as measure of rating criteria
790 (Wierzchon et al., 2012), do not fulfill these requirements. According to the present analysis,
791 intercepts also depend on model parameters assumed to reflect the degree of metacognition.
792 Moreover, when metacognition is low, intercepts are no longer sensitive to rating criteria.
793 Overall, rating criteria need to be controlled by other measures than regression intercepts.

794 A model based approach to assess rating criteria is based on the standard SDT model used to
795 estimate meta-d'. When the standard SDT model is assumed, rating criteria and the primary
796 task criterion can be directly estimated from the model (Rausch & Zehetleitner, 2016). The
797 advantage of this approach is that it allows to control for primary task criterion and rating
798 criteria at the same time. Unfortunately, this approach again requires assumptions about the
799 distributions of evidence. In addition, it has never been investigated if these estimates are
800 unbiased when the data-generating model is not the standard SDT model or a hierarchical
801 model of metacognition.

802 SDT theory provides numerous other indices for response bias (Green & Swets, 1966;
803 Macmillan & Creelman, 2005; Wickens, 2002). A non-parametric measure of primary task
804 criterion is β_K , which is calculated from the empirical receiver operating characteristic
805 (Kornbrot, 2006). When researchers construct type 2 ROC-curves, they can compute B_{roc} , the
806 analog of β_K for type 2 ROC-curves, as a measure of conservative and liberal rating criteria
807 (Fleming et al., 2010). A disadvantage of these measures is again that the number of trials
808 required for ROC-curves is large. Nevertheless, when distributional assumptions need to be
809 avoided, β_K and B_{roc} appear to be the most promising approach.

810 It should be noted that in order to control the impact of criteria on regression slopes, it is not
811 sufficient to demonstrate that the mean primary task criterion and the mean rating criteria are
812 the same between two conditions. The reason is that the effects of criteria on GLM slopes
813 often follow non-linear trends. Thus, when the variance of criteria is greater in one condition,
814 the effect of the maximal and minimal criteria will not necessarily balance out. Consequently,
815 slopes can be different between two conditions solely due to different variances of primary
816 task criteria or rating criteria. Researchers who would like to rule out than an effect on slopes
817 is not due to criteria should assess if the distributions of primary task criteria and rating
818 criteria are the same between conditions.

819 **6 Conclusion**

820 Logistic regression slopes as measures of metacognitive sensitivity always depend on the
821 primary task criterion and are independent from rating criteria according to only one specific
822 model of metacognition, namely the independent logistic model. It is argued that researchers
823 who want to quantify metacognitive sensitivity using logistic regression should provide
824 evidence that the underlying model is the independent logistic model where the primary task
825 criterion is fixed at zero. However, this will often not be feasible as number of trials required
826 to allow accurate model selection is massive. Alternatively, they can also control the primary
827 task criterion and rating criteria statistically or use alternative methods to measure
828 metacognitive sensitivity.

829 **7 Acknowledgements**

830 This research was supported by grant ZE 887/3-1 of the Deutsche Forschungsgesellschaft
831 (DFG) (to MZ). The funders had no role in study design, data collection, analysis, decision to
832 publish, or preparation of the manuscript. We are grateful to Sebastian Hellmann for helpful
833 comments.

834 **8 References**

- 835 Akaike, H. (1974). A New Look at the Statistical Model Identification. *IEEE Transactions on*
836 *Automatic Control*, 19(6), 716–723.
- 837 Baranski, J. V., & Petrusic, W. M. (1994). The calibration and resolution of confidence in
838 perceptual judgments. *Perception & Psychophysics*, 55(4), 412–28.
- 839 Barrett, A. B., Dienes, Z., & Seth, A. K. (2013). Measures of metacognition on signal-
840 detection theoretic models. *Psychological Methods*, 18(4), 535–52.
841 <http://doi.org/10.1037/a0033268>
- 842 Brockhoff, P. B., & Christensen, R. H. B. (2010). Thurstonian models for sensory
843 discrimination tests as generalized linear models. *Food Quality and Preference*, 21(3),
844 330–338. <http://doi.org/10.1016/j.foodqual.2009.04.003>
- 845 Brockhoff, P. M., & Müller, H. G. (1997). Random effect threshold models for dose-response
846 relationships with repeated measurements. *Journal of the Royal Statistical Society Series*
847 *B-Methodological*, 59(2), 431–446.
- 848 Burnham, K. P., & Anderson, D. R. (2002). *Model selection and multimodel inference: a*
849 *practical information-theoretic approach* (2nd ed.). New York: Springer.
- 850 Cul, A. Del, Dehaene, S., Reyes, P., Bravo, E., & Slachevsky, A. (2009). Causal role of
851 prefrontal cortex in the threshold for access to consciousness. *Brain*, 132, 2531–2540.
852 <http://doi.org/10.1093/brain/awp111>
- 853 DeCarlo, L. T. (1998). Signal detection theory and generalized linear models. *Psychological*
854 *Methods*, 3(2), 186–205. <http://doi.org/10.1037//1082-989X.3.2.186>
- 855 Dienes, Z. (2004). Assumptions of Subjective Measures of Unconscious Mental States.
856 *Journal of Consciousness Studies*, 11(9), 25–45.
- 857 Dienes, Z. (2011). Bayesian Versus Orthodox Statistics: Which Side Are You On?
858 *Perspectives on Psychological Science*, 6(3), 274–290.
859 <http://doi.org/10.1177/1745691611406920>
- 860 Dorfman, D. D., & Alf, E. (1969). Maximum-likelihood estimation of parameters of signal-
861 detection theory and determination of confidence intervals—Rating-method data.
862 *Journal of Mathematical Psychology*, 6, 487–496. [http://doi.org/10.1016/0022-](http://doi.org/10.1016/0022-2496(69)90019-4)
863 [2496\(69\)90019-4](http://doi.org/10.1016/0022-2496(69)90019-4)
- 864 Evans, S., & Azzopardi, P. (2007). Evaluation of a “bias-free” measure of awareness. *Spatial*
865 *Vision*, 20(1), 61–77. <http://doi.org/10.1163/156856807779369742>
- 866 Fleming, S. M., & Daw, N. D. (2016). Self-evaluation of decision performance: A general
867 Bayesian framework for metacognitive computation. *Psychological Review*, 1–59.
868 <http://doi.org/10.1002/sml.1>)
- 869 Fleming, S. M., & Lau, H. (2014). How to measure metacognition. *Frontiers in Human*

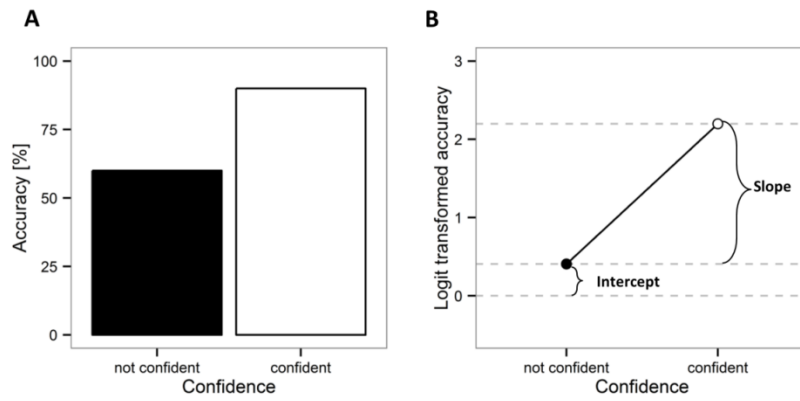
- 870 *Neuroscience*, 8, 1–9. <http://doi.org/10.3389/fnhum.2014.00443>
- 871 Fleming, S. M., Weil, R. S., Nagy, Z., Dolan, R. J., & Rees, G. (2010). Relating introspective
872 accuracy to individual differences in brain structure. *Science*, 329(5998), 1541–1543.
873 <http://doi.org/10.1126/science.1191883>
- 874 Galvin, S. J., Podd, J. V., Drga, V., & Whitmore, J. (2003). Type 2 tasks in the theory of
875 signal detectability: discrimination between correct and incorrect decisions.
876 *Psychonomic Bulletin & Review*, 10(4), 843–76.
- 877 Green, D. M., & Swets, J. A. (1966). *Signal detection theory and psychophysics*. New York:
878 Wiley.
- 879 Ince, D. C., Hatton, L., & Graham-Cumming, J. (2012). The case for open computer
880 programs. *Nature*, 482, 485–488. <http://doi.org/10.1038/nature10836>
- 881 Kiani, R., Corthell, L., & Shadlen, M. N. (2014). Choice certainty is informed by both
882 evidence and decision time. *Neuron*, 84(6), 1329–1342.
883 <http://doi.org/10.1016/j.neuron.2014.12.015>
- 884 Kornbrot, D. E. (2006). Signal detection theory, the approach of choice: model-based and
885 distribution-free measures and evaluation. *Perception & Psychophysics*, 68(3), 393–414.
- 886 Kuniyoto, C., Miller, J., & Pashler, H. (2001). Confidence and accuracy of near-threshold
887 discrimination responses. *Consciousness and Cognition*, 10(3), 294–340.
888 <http://doi.org/10.1006/ccog.2000.0494>
- 889 Lipovetsky, S. (2015). Analytical closed-form solution for binary logit regression by
890 categorical predictors. *Journal of Applied Statistics*, 42(1), 37–49.
891 <http://doi.org/10.1080/02664763.2014.932760>
- 892 Luce, R. D., & Suppes, P. (1965). Preference, utility, and subjective probability. In & E. G.
893 R. D. Luce, R. Bush (Ed.), *Handbook of Mathematical Psychology*. New York: Wiley.
- 894 Macmillan, N. A., & Creelman, C. D. (2005). *Detection Theory. A user's guide*. Mahwah,
895 NY: Lawrence Erlbaum Associates.
- 896 Maniscalco, B., & Lau, H. (2012). A signal detection theoretic approach for estimating
897 metacognitive sensitivity from confidence ratings. *Consciousness and Cognition*, 21(1),
898 422–30. <http://doi.org/10.1016/j.concog.2011.09.021>
- 899 Maniscalco, B., & Lau, H. (2016). The signal processing architecture underlying subjective
900 reports of sensory awareness. *Neuroscience of Consciousness*, (November 2015), 1–17.
901 <http://doi.org/10.1093/nc/niw002>
- 902 Masson, M. E. J., & Rotello, C. M. (2009). Sources of bias in the Goodman-Kruskal gamma
903 coefficient measure of association: implications for studies of metacognitive processes.
904 *Journal of Experimental Psychology. Learning, Memory, and Cognition*, 35(2), 509–27.
905 <http://doi.org/10.1037/a0014876>
- 906 Moran, R., Teodorescu, A. R., & Usher, M. (2015). Post choice information integration as a
907 causal determinant of confidence: Novel data and a computational account. *Cognitive*
908 *Psychology*, 78, 99–147. <http://doi.org/10.1016/j.cogpsych.2015.01.002>
- 909 Morey, R. D., & Rouder, J. N. (2015). BayesFactor: Computation of Bayes factors for
910 common designs. R package version 0.9.10-1.
- 911 Morin, A., Urban, J., Adams, P. D., Foster, I., Sali, A., Baker, D., & Sliz, P. (2012). Shining
912 Light into Black Boxes. *Science*, 336(6078), 159–160.
913 <http://doi.org/10.1126/science.1218263>

- 914 Nelder, J. A., & Mead, R. (1965). A simplex method for function minimization. *The*
915 *Computer Journal*, 7, 308--313.
- 916 Nelson, T. O. (1984). A comparison of current measures of the accuracy of feeling-of-
917 knowing predictions. *Psychological Bulletin*, 95(1), 109–133.
918 <http://doi.org/10.1037//0033-2909.95.1.109>
- 919 Paz, L., Insabato, A., Zylberberg, A., Deco, G., & Sigman, M. (2016). Confidence through
920 consensus: a neural mechanism for uncertainty monitoring. *Scientific Reports*, 6, 21830.
921 <http://doi.org/10.1038/srep21830>
- 922 Peters, M. A. K., & Lau, H. (2015). Human observers have optimal introspective access to
923 perceptual processes even for visually masked stimuli. *eLife*, 4(OCTOBER2015), 1–30.
924 <http://doi.org/10.7554/eLife.09651>
- 925 Pleskac, T. J., & Busemeyer, J. R. (2010). Two-stage dynamic signal detection: a theory of
926 choice, decision time, and confidence. *Psychological Review*, 117(3), 864–901.
927 <http://doi.org/10.1037/a0019737>
- 928 R Core Team. (2014). R: A Language and Environment for Statistical Computing. Vienna,
929 Austria: R Foundation for Statistical Computing. Retrieved from [http://www.r-](http://www.r-project.org/)
930 [project.org/](http://www.r-project.org/)
- 931 Rahnev, D. a, Maniscalco, B., Lubner, B., Lau, H., & Lisanby, S. H. (2012). Direct injection of
932 noise to the visual cortex decreases accuracy but increases decision confidence. *Journal*
933 *of Neurophysiology*, 107(6), 1556–63. <http://doi.org/10.1152/jn.00985.2011>
- 934 Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review*, 85(2), 59–108.
- 935 Rausch, M., Müller, H. J., & Zehetleitner, M. (2015). Metacognitive sensitivity of subjective
936 reports of decisional confidence and visual experience. *Consciousness and Cognition*,
937 35, 192–205. <http://doi.org/10.1016/j.concog.2015.02.011>
- 938 Rausch, M., & Zehetleitner, M. (2014). A comparison between a visual analogue scale and a
939 four point scale as measures of conscious experience of motion. *Consciousness and*
940 *Cognition*, 28, 126–140. <http://doi.org/10.1016/j.concog.2014.06.012>
- 941 Rausch, M., & Zehetleitner, M. (2016). Visibility is not equivalent to confidence in a low
942 contrast orientation discrimination task. *Frontiers in Psychology*, 7, 591.
943 <http://doi.org/10.3389/fpsyg.2016.00591>
- 944 Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests
945 for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review*, 16(2),
946 225–37. <http://doi.org/10.3758/PBR.16.2.225>
- 947 Sandberg, K., Bibby, B. M., & Overgaard, M. (2013). Measuring and testing awareness of
948 emotional face expressions. *Consciousness and Cognition*, 22(3), 806–9.
949 <http://doi.org/10.1016/j.concog.2013.04.015>
- 950 Sandberg, K., Timmermans, B., Overgaard, M., & Cleeremans, A. (2010). Measuring
951 consciousness: Is one measure better than the other? *Consciousness and Cognition*,
952 19(4), 1069–1078. <http://doi.org/10.1016/j.concog.2009.12.013>
- 953 Schwarz, G. (1978). Estimating the dimensions of a model. *The Annals of Statistics*, 6(2),
954 461–464. <http://doi.org/10.1214/aos/1176348654>
- 955 Scott, R. B., Dienes, Z., Barrett, A. B., Bor, D., & Seth, A. K. (2014). Blind insight:
956 Metacognitive discrimination despite chance task performance. *Psychological Science*,
957 25, 1–20. <http://doi.org/10.1177/0956797614553944>

- 958 Siedlecka, M., Paulewicz, B., & Wierzchoń, M. (2016). But I Was So Sure! Metacognitive
959 Judgments Are Less Accurate Given Prospectively than Retrospectively. *Frontiers in*
960 *Psychology*, 7(February). <http://doi.org/10.3389/fpsyg.2016.00218>
- 961 Wickens, T. D. (2002). *Elementary signal detection theory*. New York: Oxford University
962 Press.
- 963 Wierzchoń, M., Asanowicz, D., Paulewicz, B., & Cleeremans, A. (2012). Subjective
964 measures of consciousness in artificial grammar learning task. *Consciousness and*
965 *Cognition*, 21(3), 1141–53. <http://doi.org/10.1016/j.concog.2012.05.012>
- 966 Wierzchoń, M., Paulewicz, B., Asanowicz, D., Timmermans, B., & Cleeremans, A. (2014).
967 Different subjective awareness measures demonstrate the influence of visual
968 identification on perceptual awareness ratings. *Consciousness and Cognition*, 27, 109–
969 120. <http://doi.org/10.1016/j.concog.2014.04.009>
- 970 Williams, J., Ramaswamy, D., & Oulhaj, A. (2006). 10 Hz flicker improves recognition
971 memory in older people. *BMC Neuroscience*, 7, 21. [http://doi.org/10.1186/1471-2202-7-](http://doi.org/10.1186/1471-2202-7-21)
972 21
- 973 Yeung, N., & Summerfield, C. (2012). Metacognition in human decision-making: confidence
974 and error monitoring. *Philosophical Transactions of the Royal Society of London. Series*
975 *B, Biological Sciences*, 367(1594), 1310–21. <http://doi.org/10.1098/rstb.2011.0416>
- 976
- 977

Table 1. Parameters of the hierarchical and the independent model of metacognition

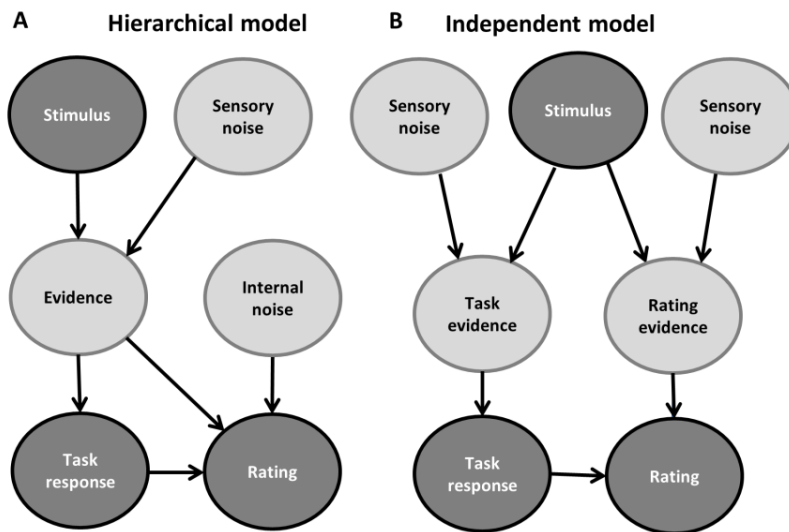
Name	Symbol	Conceptual interpretation	Part of which models?
Sensitivity	d	Objective discrimination ability of the observer between the two stimulus alternatives	hierarchical and independent
Primary task criterion	θ	Bias of the observer towards one of the two stimulus alternatives	hierarchical and independent
Rating criteria spread	τ	Degree of conservativeness of subjective reports	hierarchical and independent
Internal noise	σ	Amount of distortion during metacognitive read-out of sensory evidence	hierarchical
Rating sensitivity	d_m	Ability of the second, metacognitive channel to discern between the two stimulus alternatives	independent



979

980 *Fig. 1.* Quantifying the relationship between trial accuracy and subjective reports by logistic
 981 regression. (A): Data of a hypothetical experiment. Task accuracy in % correct is plotted
 982 separately for two categories of subjective reports, “not confident” and “confident”. (B):
 983 Same data, but accuracy transformed into logits. Logistic regression is based on fitting a
 984 linear function on such transformed data. The slope of the regression line is interpreted as
 985 metacognitive sensitivity, the intercept as criterion.

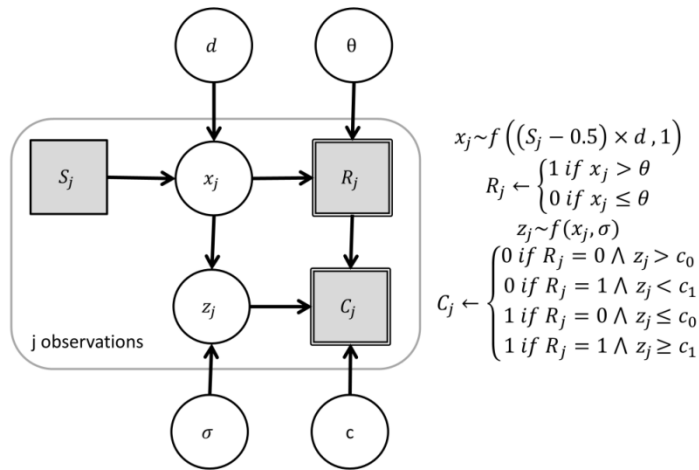
986



987

988 *Fig. 2.* The hierarchical and the independent model of metacognition. According to the
 989 hierarchical model, the rating is generated by the same evidence as the response to the
 990 primary task, but the evidence is distorted by internal noise. The task response determines
 991 which criteria are applied to select a subjective report. According to the independent model,
 992 evidence is created in parallel by two channels. The task response is selected based on the
 993 evidence selected in one of the channels. The rating response depends on whether the
 994 evidence sampled independently in the second channel confirms the response.

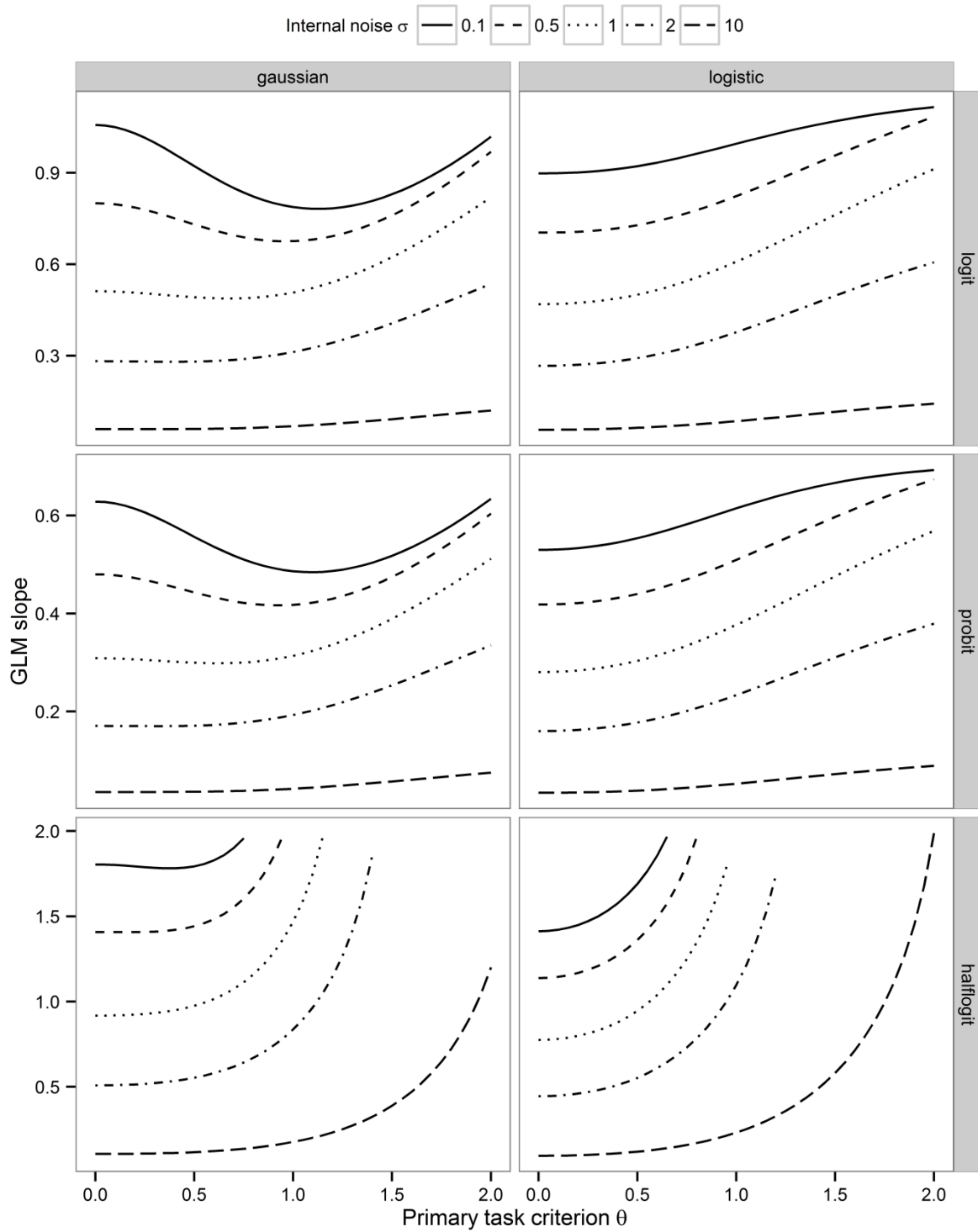
995



996

997 *Fig. 3.* Graphical model of the hierarchical model of metacognition. In each trial, the observer
 998 is faced with sensory evidence x , which depends on the stimulus S as well as the observers'
 999 sensitivity to discriminate between the two stimuli d . The response R is selected based on a
 1000 comparison between x and the task bias θ . In addition, the decision variable z for selecting
 1001 one out of several rating options depends on x as well as on unsystematic noise σ . The rating
 1002 C depends on the decision variable z , the response R , and the rating criterion c .

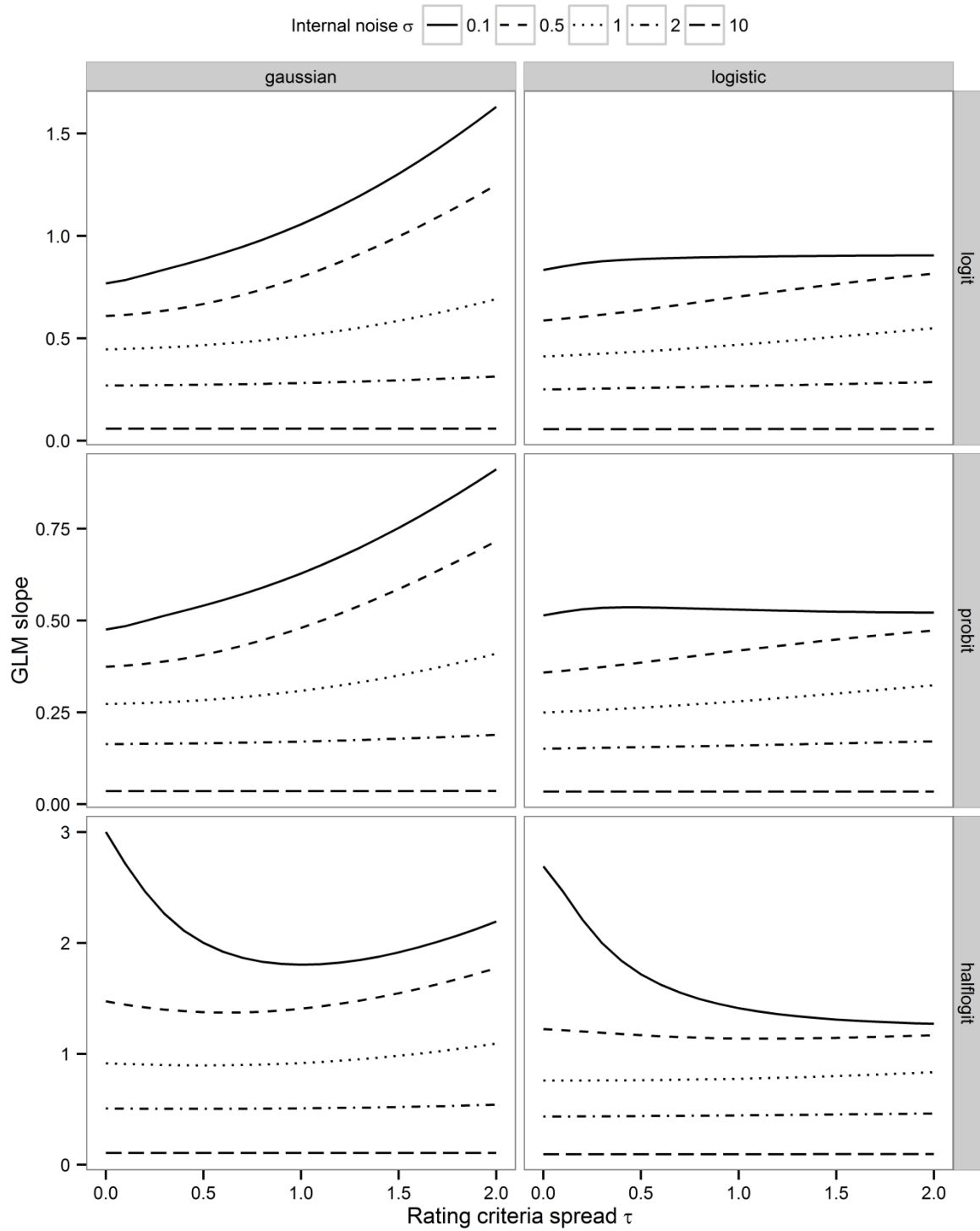
1003



1004

1005 *Fig. 4.* Generalized linear regression slopes according to the hierarchical model of
 1006 metacognition as a function of primary task criterion θ (X-Axis), internal noise σ (different
 1007 lines), distributions of evidence (different columns) and link function (different rows). The
 1008 other parameters were fixed to $d = 1$ and $\tau = 1$.

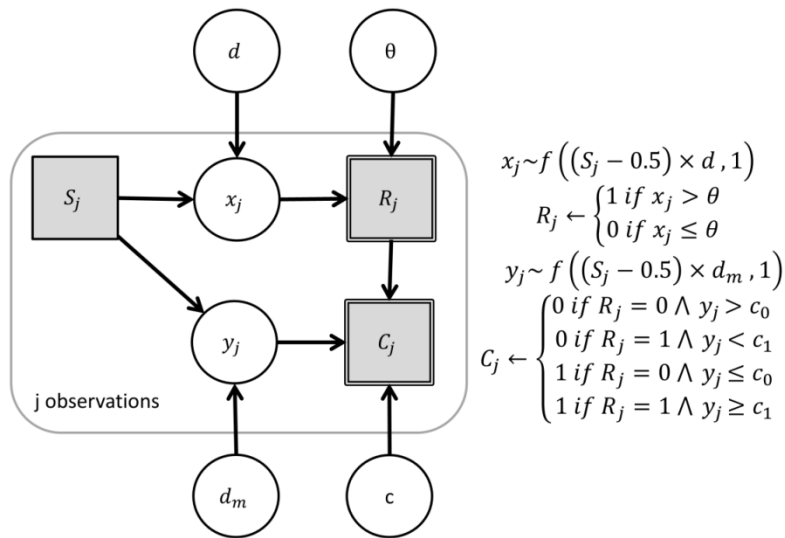
1009



1010

1011 *Fig. 5.* Generalized linear regression slopes according to the hierarchical model of
 1012 metacognition as a function of rating criteria spread τ (X-Axis), internal noise σ (different
 1013 lines), distributions of evidence (different columns) and link function (different rows). The
 1014 other parameters were fixed: $d = 1$, $\theta = 0$.

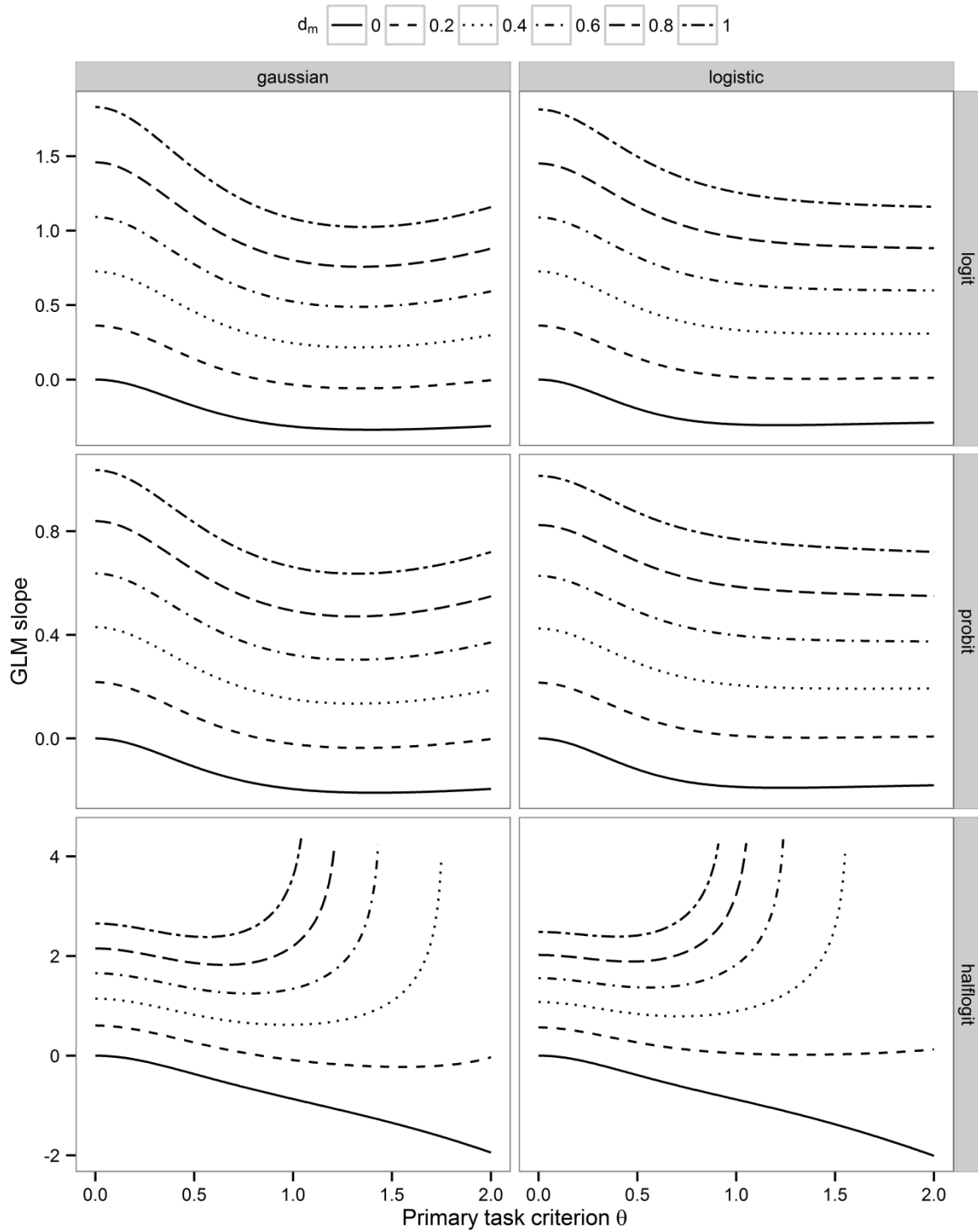
1015



1016

1017 *Fig. 6.* Graphical model representing the dual evidence model of metacognition. Each trial,
 1018 the observer is faced with sensory evidence x , which depends on the stimulus S as well as the
 1019 observers' sensitivity to discriminate between the two stimuli d . The response R is selected
 1020 based on a comparison between x and the task bias θ . In contrast to the single evidence
 1021 model, the decision variable y for selecting one out of the several rating options depends on
 1022 the stimulus S as well as on the metacognitive access parameter d_m . The rating C depends on
 1023 the decision variable y , the response R , and the rating criterion c .

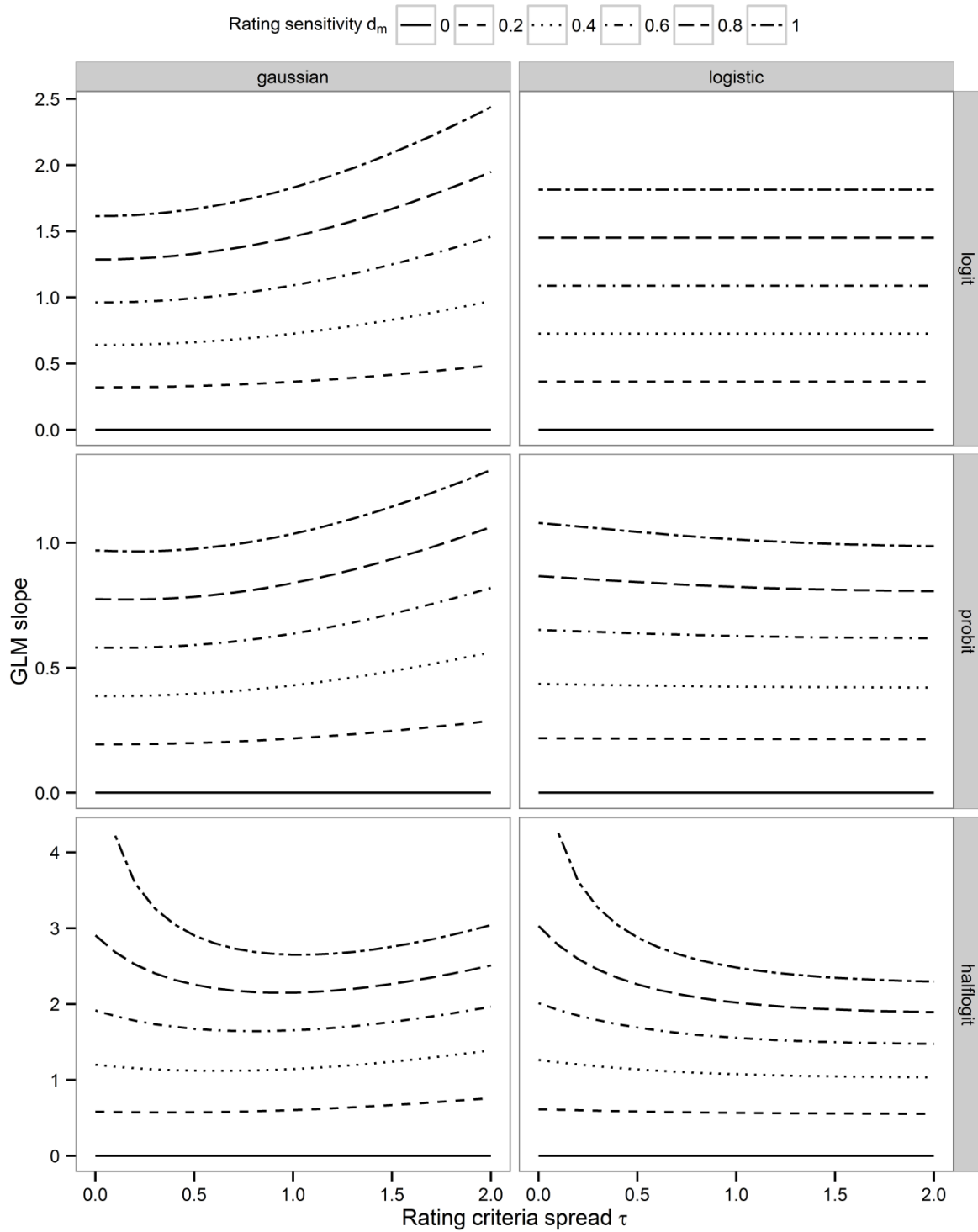
1024



1025

1026 *Fig. 7.* Generalized linear regression slopes according to the independent model of
 1027 metacognition as a function of primary task criterion θ (X-Axis), rating sensitivity d_m
 1028 (different lines), distributions of evidence (different columns) and link function (different
 1029 rows). The other parameters were fixed: $d = 1$, $\tau = 1$.

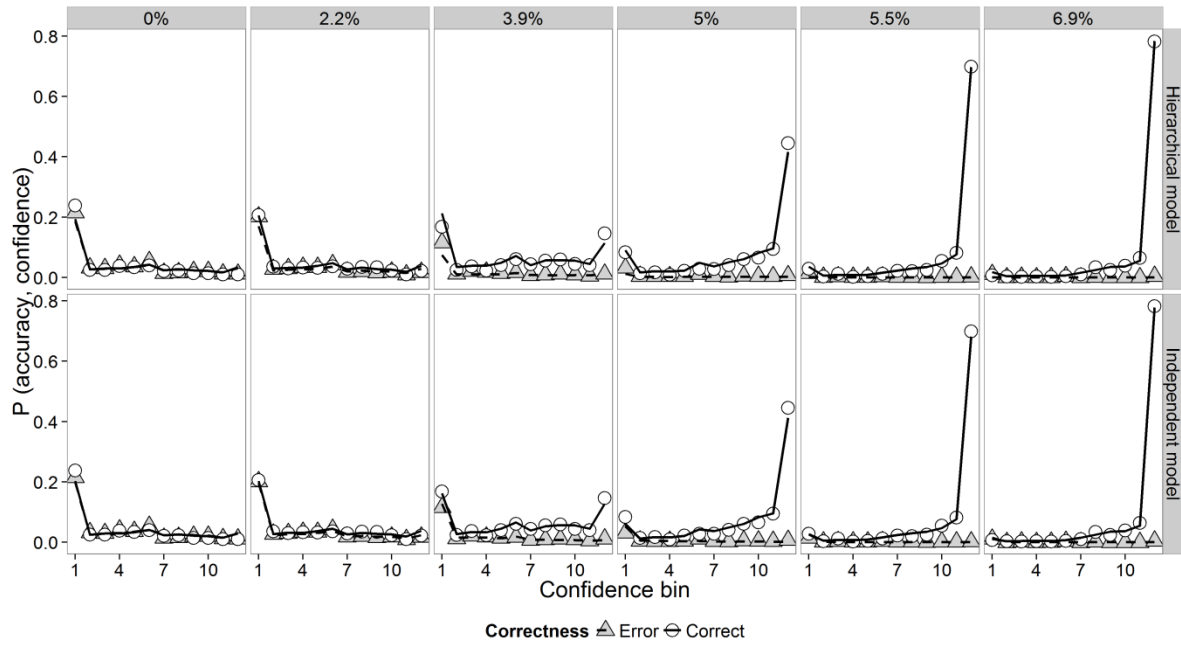
1030



1031

1032 *Fig. 8.* Generalized linear regression according to the independent model of metacognition as
 1033 a function of rating criteria spread τ (X-Axis), rating sensitivity d_m (different lines),
 1034 distributions of evidence (different columns) and link function (different rows). The other
 1035 parameters were fixed: $d = 1$, $\theta = 0$.

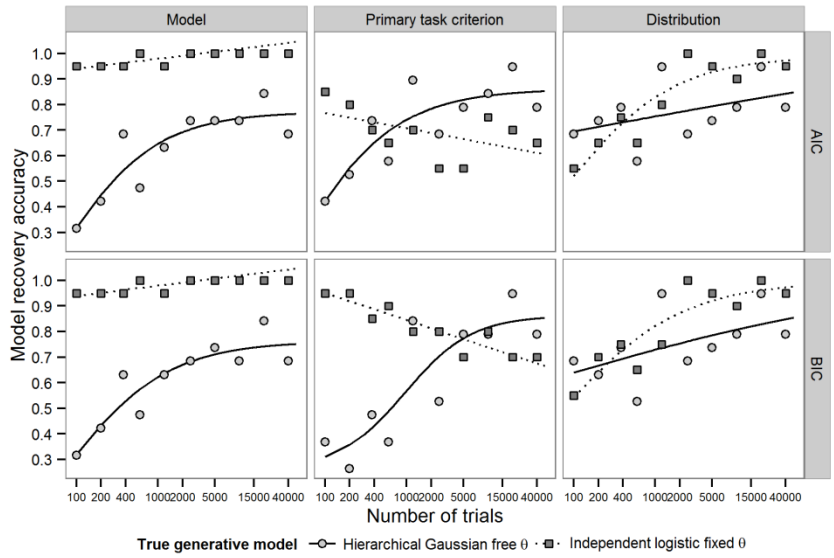
1036



1037

1038 *Fig. 9.* Joint probability of primary task accuracy (symbols) and confidence bins (on the X-
 1039 Axis) as a function of contrast levels (separate columns). Different rows indicate the
 1040 prediction of the two best performing models, the hierarchical logistic model without task
 1041 bias (upper row) and the independent logistic model with task bias (lower row).

1042



1043

1044 *Fig. 10.* Accuracy of model recovery as a function of model feature (model: hierarchical vs.
 1045 independent, bias: θ free vs. fixed, distribution: logistic vs. Gaussian; in different columns),
 1046 the true generative model (different symbols), number of simulated trials (X-Axis) as well as
 1047 the goodness-of-fit measure (AIC_c vs BIC).

1048

1049

1050 **9 Appendix A**

1051 According to SDT theory, when participants are instructed to indicate which out of two
 1052 possible alternative stimuli was presented, their perceptual systems provides sensory
 1053 evidence, with is randomly sampled out of a distribution depending on the stimulus in the
 1054 current trial. This distribution f and the corresponding cumulative density function F can be
 1055 the normal distribution or the logistic distribution, which we both parametrize here by mean
 1056 and standard deviation. When $S = 1$, then

$$x \sim f_{\frac{d}{2},1} \quad (\text{A1})$$

1057 and correspondingly, when $S = 0$, then

$$x \sim f_{-\frac{d}{2},1} \quad (\text{A2})$$

1058 In the hierarchical model, the decision variable z evaluated when selecting a subjective report
 1059 depends on x overlaid by additive noise. Consequently, z will be distributed with a mean of x
 1060 and a standard deviation of σ :

$$z \sim f_{x,\sigma} \quad (\text{A3})$$

1061 Participants are assumed to report confidence about a response $R = 1$ when $z > c_1$, and
 1062 about a response $R = 0$ when $z < c_0$. The probability of giving a correct response and being
 1063 correct given $S = 1$ can be computed as:

$$P(T = 1 \wedge C = 1|S = 1) = P(x > \theta|S) \times P(z > c_1|x > \theta, S = 1) \quad (\text{A4})$$

$$= \int_{\theta}^{\infty} P(x) \times P(z > c_1|x) dx \quad (\text{A5})$$

$$= \int_{\theta}^{\infty} f_{\frac{d}{2},1}(x) \times (1 - F_{x,\sigma}(c_1)) dx \quad (\text{A6})$$

1064 $P(T = 0 \wedge C = 0|S = 0)$, $P(T = 0 \wedge C = 0|S = 1)$, $P(T = 0 \wedge C = 1|S = 0)$, $P(T = 0 \wedge$
 1065 $C = 1|S = 1)$, $P(T = 1 \wedge C = 0|S = 0)$, $P(T = 1 \wedge C = 0|S = 1)$, and $P(T = 1 \wedge C =$
 1066 $1|S = 0)$ are computed analogously:

$$P(T = 0 \wedge C = 0|S = 0) = \int_{\theta}^{\infty} f_{-\frac{d}{2},1}(x) \times F_{x,\sigma}(c_1) dx \quad (\text{A7})$$

$$P(T = 0 \wedge C = 0|S = 1) = \int_{-\infty}^{\theta} f_{\frac{d}{2},1}(x) \times (1 - F_{x,\sigma}(c_0)) dx \quad (\text{A8})$$

$$P(T = 0 \wedge C = 1|S = 0) = \int_{\theta}^{\infty} f_{-\frac{d}{2},1}(x) \times (1 - F_{x,\sigma}(c_1)) dx \quad (\text{A9})$$

$$P(T = 0 \wedge C = 1|S = 1) = \int_{-\infty}^{\theta} f_{\frac{d}{2},1}(x) \times F_{x,\sigma}(c_0) dx \quad (\text{A10})$$

$$P(T = 1 \wedge C = 0 | S = 0) = \int_{-\infty}^{\theta} f_{-\frac{d}{2},1}(x) \times (1 - F_{x,\sigma}(c_0)) dx \quad (\text{A11})$$

$$P(T = 1 \wedge C = 0 | S = 1) = \int_{\theta}^{\infty} f_{\frac{d}{2},1}(x) \times F_{x,\sigma}(c_1) dx \quad (\text{A12})$$

$$P(T = 1 \wedge C = 1 | S = 0) = \int_{-\infty}^{\theta} f_{-\frac{d}{2},1}(x) \times F_{x,\sigma}(c_0) dx \quad (\text{A13})$$

1067 For simplicity, we assume that $c_1 - \theta = \theta - c_0 = \tau$, leaving only one parameter of rating
 1068 criteria τ , which represents the conservativeness of rating criteria.

1069 Eq. (A6) – (A12) can be used to compute the probability of being correct given a subjective
 1070 report of $C = 0$ and $C = 1$, respectively.

$$P(T = 1 | C = 0) = \frac{P(T = 1 \wedge C = 0)}{P(C = 0)} \quad (\text{A14})$$

$$P(T = 1 \wedge C = 0) = \frac{1}{2} \times (P(T = 1 \wedge C = 0 | S = 0) + P(T = 1 \wedge C = 0 | S = 1)) \quad (\text{A15})$$

$$P(C = 0) = \frac{1}{2} \times (P(T = 1 \wedge C = 0 | S = 0) + P(T = 1 \wedge C = 0 | S = 1) + P(T = 0 \wedge C = 0 | S = 0) + P(T = 0 \wedge C = 0 | S = 1)) \quad (\text{A16})$$

$$P(T = 1 | C = 1) = \frac{P(T = 1 \wedge C = 1)}{P(C = 1)} \quad (\text{A17})$$

$$P(T = 1 \wedge C = 1) = \frac{1}{2} \times (P(T = 1 \wedge C = 1 | S = 0) + P(T = 1 \wedge C = 1 | S = 1)) \quad (\text{A18})$$

$$P(C = 1) = \frac{1}{2} \times (P(T = 1 \wedge C = 1 | S = 0) + P(T = 1 \wedge C = 1 | S = 1) + P(T = 0 \wedge C = 1 | S = 0) + P(T = 0 \wedge C = 1 | S = 1)) \quad (\text{A19})$$

1071

1072

1073 **10 Appendix B**

1074 According to the independent model, participants select the response R based on a
 1075 comparison between the sample of sensory evidence x and the primary task criterion θ , just as
 1076 in the standard SDT model. However, subjective reports are assumed to be based on a second
 1077 independent sample of sensory evidence y . The distribution f and the corresponding
 1078 cumulative distribution function F are assumed to be the same as those from which x is
 1079 sampled, except for the mean, which is described by the metacognitive access parameter
 1080 d_m . When $S = 1$, then

$$y \sim f_{\frac{1}{2}d_m,1} \quad (\text{B1})$$

1081 and correspondingly, when $S = 0$, then

$$y \sim f_{-\frac{1}{2}d_m,1} \quad (\text{B2})$$

1082 Participants are assumed to report confidence about a response $R = 1$ when $y > c_1$, and
 1083 about a response $R = 0$ when $y < c_0$. The probability of giving a correct response and being
 1084 correct given $S = 1$ in the dual evidence model is obtained by:

$$P(T = 1 \wedge C = 1|S = 1) = P(x > \theta|S) \times P(y > c_1|S = 1) \quad (\text{B4})$$

$$= \left(1 - F_{\frac{d}{2},1}(\theta)\right) \times \left(1 - F_{\frac{1}{2}d_m,1}(c_1)\right) \quad (\text{B5})$$

1085 $P(T = 0 \wedge C = 0|S = 0)$, $P(T = 0 \wedge C = 0|S = 1)$, $P(T = 0 \wedge C = 1|S = 0)$, $P(T = 0 \wedge$
 1086 $C = 1|S = 1)$, $P(T = 1 \wedge C = 0|S = 0)$, $P(T = 1 \wedge C = 0|S = 1)$, and $P(T = 1 \wedge C =$
 1087 $1|S = 0)$ are computed analogously:

$$P(T = 0 \wedge C = 0|S = 0) = \left(1 - F_{-\frac{d}{2},1}(\theta)\right) \times F_{-\frac{1}{2}d_m,1}(c_1) \quad (\text{B6})$$

$$P(T = 0 \wedge C = 0|S = 1) = F_{\frac{d}{2},1}(\theta) \times \left(1 - F_{\frac{1}{2}d_m,1}(c_0)\right) \quad (\text{B7})$$

$$P(T = 0 \wedge C = 1|S = 0) = \left(1 - F_{-\frac{d}{2},1}(\theta)\right) \times \left(1 - F_{-\frac{1}{2}d_m,1}(c_1)\right) \quad (\text{B8})$$

$$P(T = 0 \wedge C = 1|S = 1) = F_{\frac{d}{2},1}(\theta) \times F_{\frac{1}{2}d_m,1}(c_0) \quad (\text{B9})$$

$$P(T = 1 \wedge C = 0|S = 0) = F_{-\frac{1}{2}d_m,1}(\theta) \times \left(1 - F_{-\frac{1}{2}d_m,1}(c_0)\right) \quad (\text{B10})$$

$$P(T = 1 \wedge C = 0|S = 1) = \left(1 - F_{\frac{1}{2}d_m,1}(\theta)\right) \times F_{\frac{1}{2}d_m,1}(c_1) \quad (\text{B11})$$

$$P(T = 1 \wedge C = 1|S = 0) = F_{-\frac{1}{2}d_m,1}(\theta) \times F_{-\frac{1}{2}d_m,1}(c_0) \quad (\text{B12})$$

1088

1089

1090

1091

Supplementary Material

1092

1093

to

1094

1095

Should metacognition be measured by logistic regression?

1096

1097 Manuel Rausch ^{1,2} and Michael Zehetleitner ^{1,2}

1098

1099 ¹ Katholische Universität Eichstätt-Ingolstadt, Eichstätt, Germany

1100 ² Ludwig-Maximilians-Universität München, Munich, Germany

1101

1102

1103 Correspondence should be addressed at:

1104 Manuel Rausch

1105 Katholische Univerität Eichstätt-Ingolstadt

1106 Psychologie II

1107 Ostenstraße 25, "Waisenhaus"

1108 85072 Eichstätt

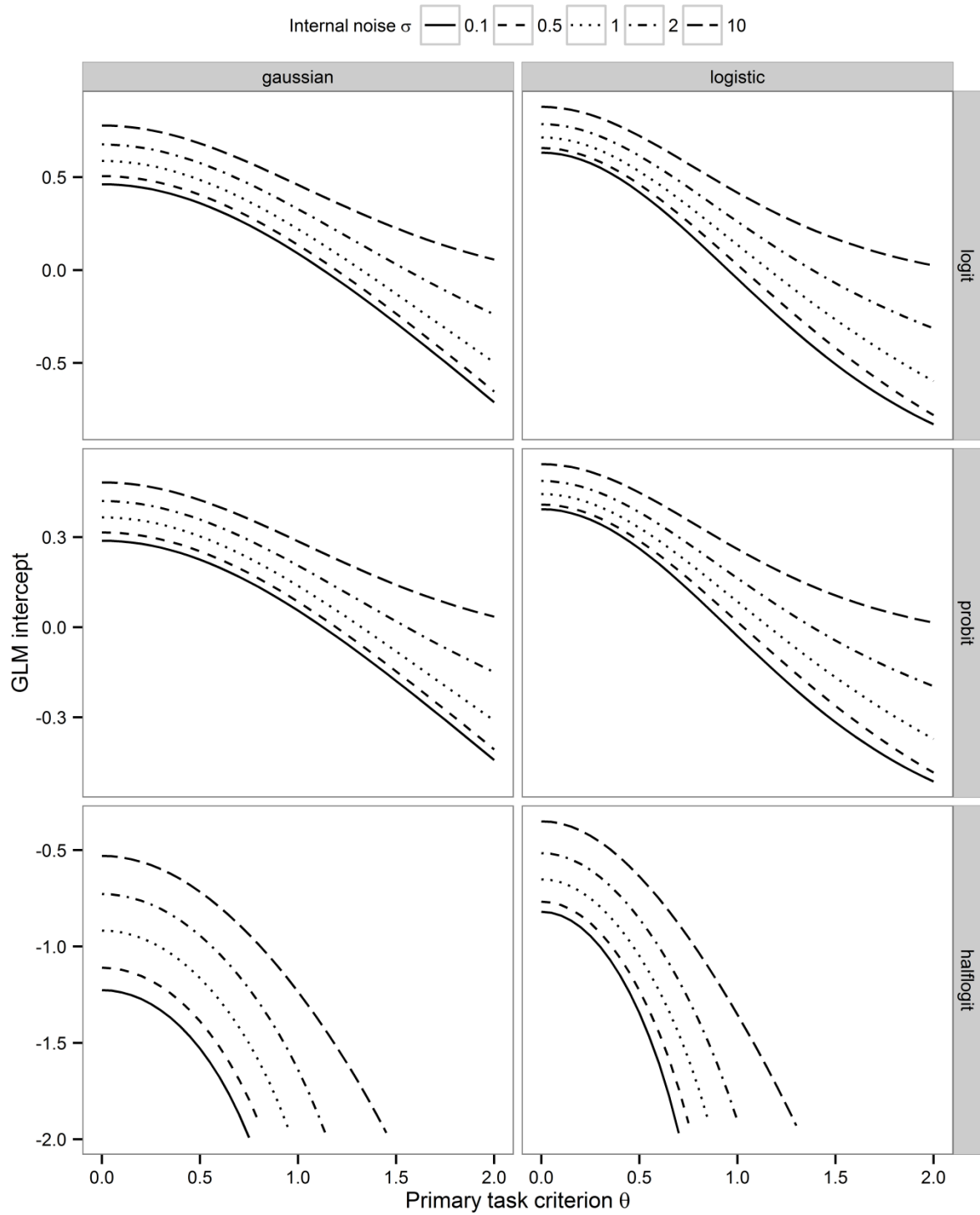
1109 Germany

1110 Phone: +49 8421 93 21639

1111 Email: manuel.rausch@ku.de

1112 <http://www.ku.de/ppf/psychologie/psych2/mitarbeiter/m-rausch/>

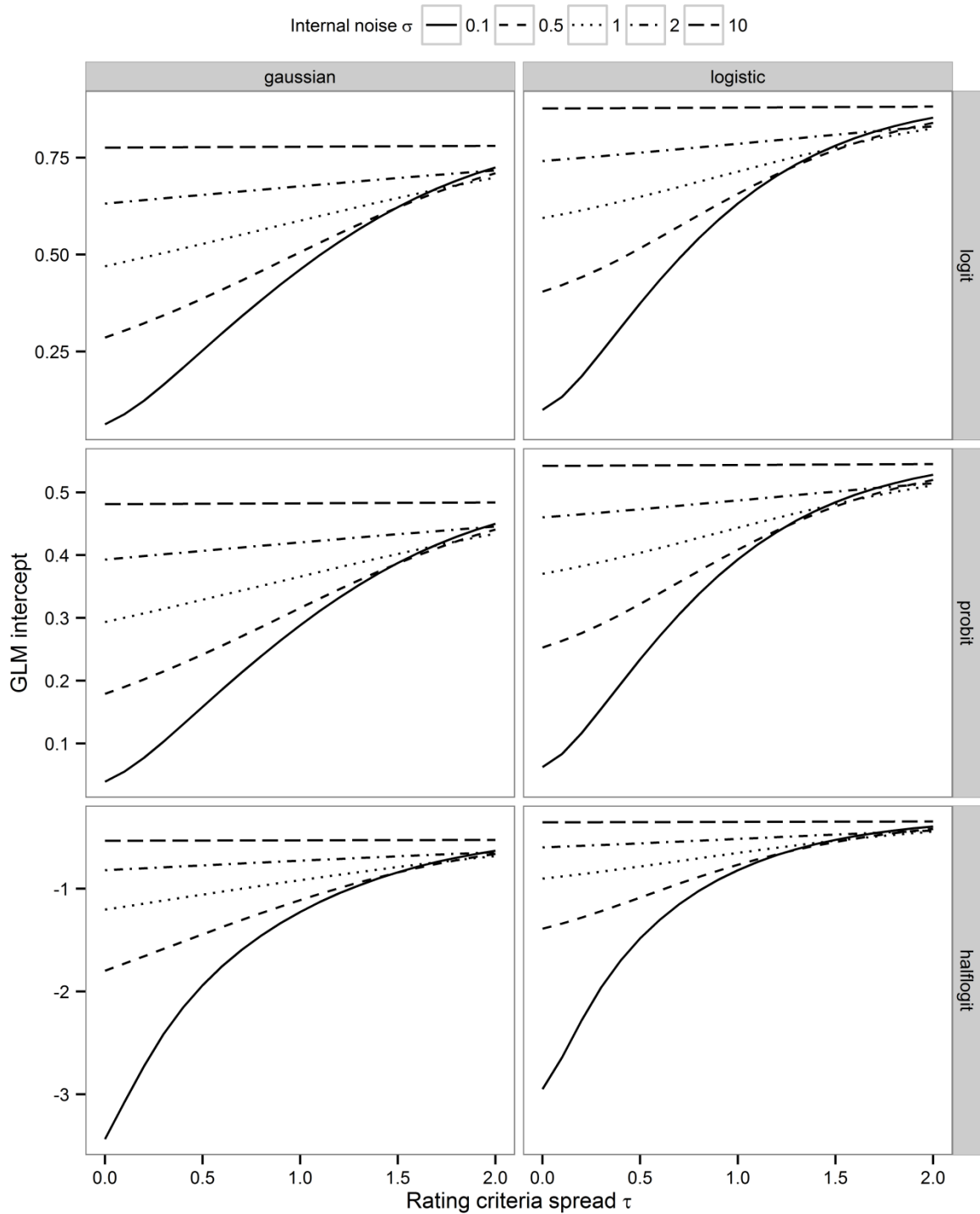
1113



1114

1115 *Figure S1.* Generalized linear regression intercepts according to the hierarchical model of
 1116 metacognition as a function of primary task criterion θ (X-Axis), internal noise σ (different
 1117 lines), distributions of evidence (different columns) and link function (different rows). The
 1118 other parameters were fixed: $d = 1$, $\tau = 1$.

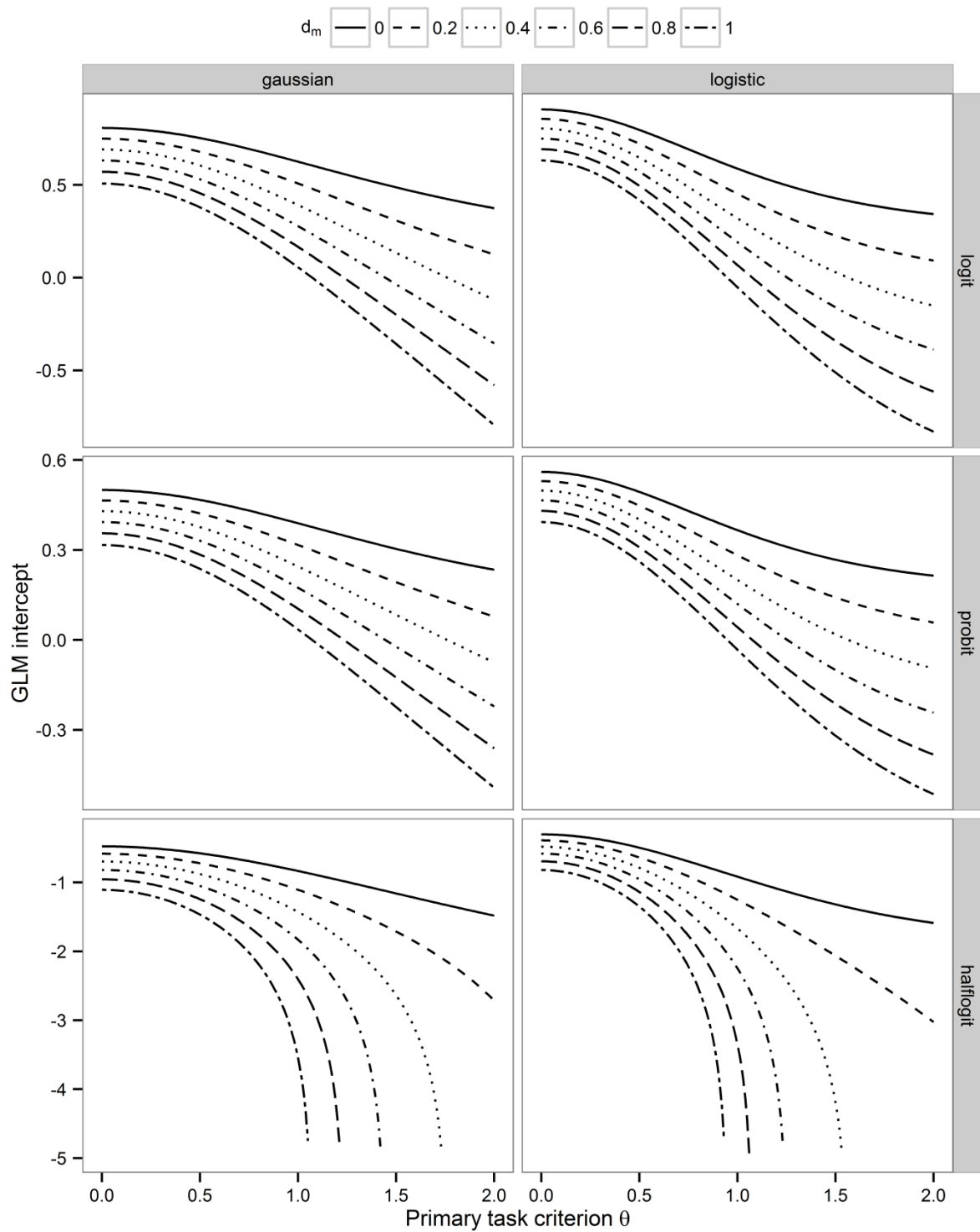
1119



1120

1121 *Figure S2.* Generalized linear regression intercepts according to the hierarchical model of
 1122 metacognition as a function of rating criteria spread τ (X-Axis), internal noise σ (different
 1123 lines), distributions of evidence (different columns) and link function (different rows). The
 1124 other parameters were fixed: $d = 1$, $\theta = 0$.

1125

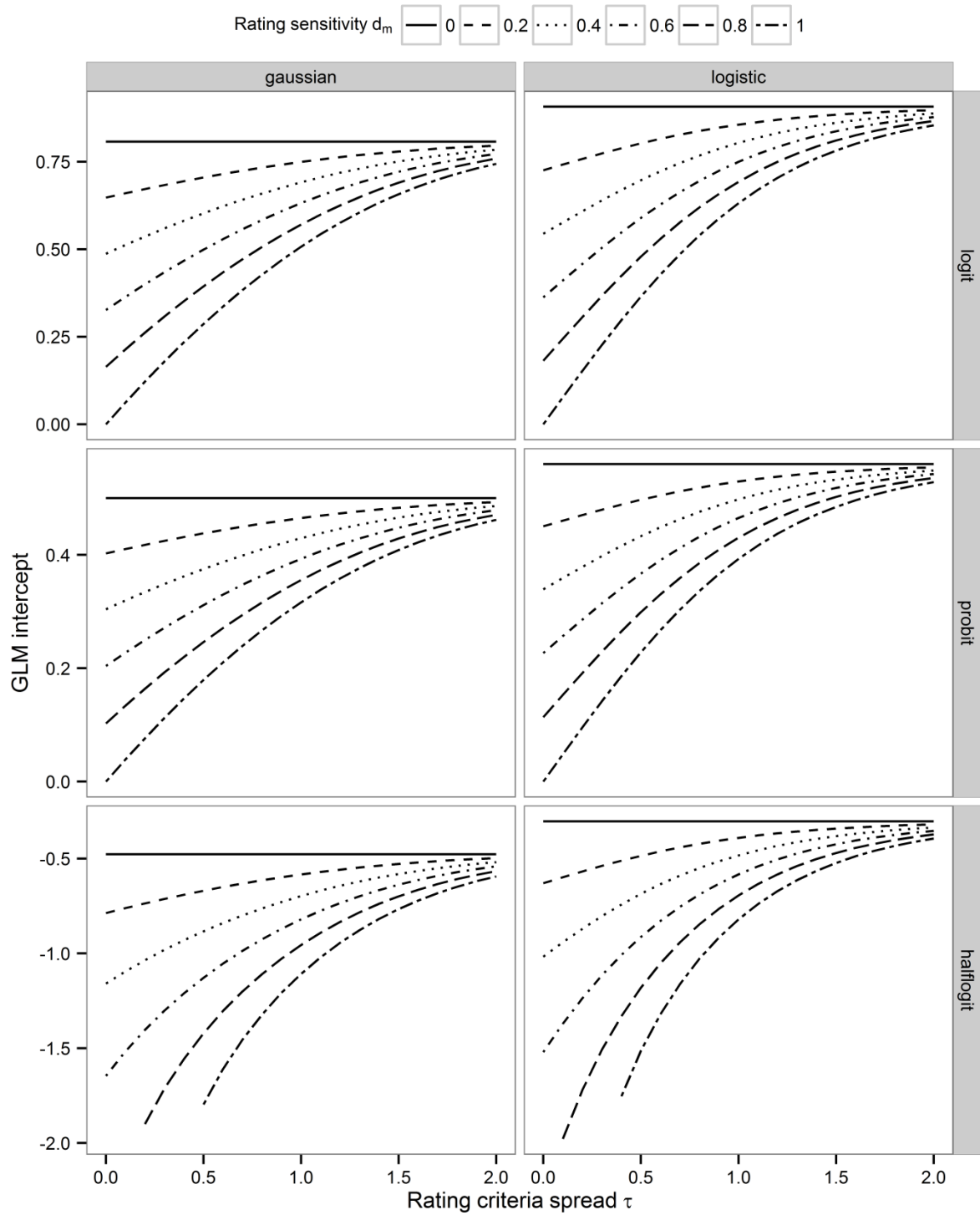


1126

1127 *Figure S3.* Generalized linear regression intercept according to the independent model of
 1128 metacognition as a function of primary task criterion θ (X-Axis), rating sensitivity d_m
 1129 (different lines), distributions of evidence (different columns) and link function (different
 1130 rows). The other parameters were fixed: $d = 1$, $\tau = 1$.

1131

1132



1133

1134 *Figure S4.* Generalized linear regression intercept according to the independent model of
 1135 metacognition as a function of rating criteria spread τ (X-Axis),, rating sensitivity d_m
 1136 (different lines), distributions of evidence (different columns) and link function (different
 1137 rows). The other parameters were fixed: $d = 1, \theta = 0$

1138