

Do modeling choices matter for the reliability of individual difference measures in conflict tasks?

Michelle C. Donzallaz¹, Udo Boehm², Andrew Heathcote^{1,3}, Chris Donkin⁴,
Dora Matzke¹, & Julia M. Haaf⁵

¹ Psychological Methods, University of Amsterdam

² Centrum Wiskunde & Informatica, Amsterdam

³ University of Newcastle, Australia

⁴ LMU Munich

⁵ Department of Psychology, University of Potsdam

Abstract

There is a growing realization that experimental tasks that produce reliable effects in group comparisons can simultaneously provide unreliable assessments of individual differences. Proposed solutions to this “reliability paradox” range from collecting more test trials to modifying the tasks and/or the way in which effects are measured from these tasks. Here we systematically compare two proposed modeling solutions in a cognitive conflict task. Using the ratio of individual variability of the conflict effect (i.e., signal) and the trial-by-trial variation in the data (i.e., noise) obtained from Bayesian hierarchical modeling, we examine whether improving statistical modeling may improve the reliability of individual differences assessment in four Stroop datasets. The proposed improvements are 1) increasing the descriptive adequacy of the statistical models from which conflict effects are derived, and 2) using psychologically-motivated measures from cognitive models. Our results show that modeling choices do not have a consistent effect on the signal-to-noise ratio: the proposed solutions improved reliability in only one of the four datasets. We provide analytical and simulation-based approaches to compute the signal-to-noise ratio for a range of models of varying sophistication and discuss their potential to aid in developing and comparing new measurement solutions to the reliability paradox.

Keywords: reliability, cognitive conflict, Bayesian hierarchical modeling

Cognitive conflict is commonly assessed using experimental paradigms by comparing conditions requiring higher versus lower levels of cognitive control. One such paradigm is the classic color-Stroop task (Stroop, 1935), in which participants indicate the color in which a color-word is printed. The Stroop effect is conventionally measured by the response time (RT) difference between trials in which the printed color and the meaning of the word are different (i.e., incongruent; e.g., “red” written in blue) and trials in which they match (i.e., congruent; e.g., “red” written in

19 red). Similar cognitive conflict paradigms with a congruency manipulation are the Flanker (Erik-
 20 sen & Eriksen, 1974) and the Simon task (Simon & Rudell, 1967). In all of them, the incongruent
 21 condition requires more control to inhibit automatic associations between the wrong response and
 22 irrelevant information (e.g., color). Positive conflict effects, that is the RT difference between the
 23 incongruent and the congruent conditions, provide an index of cognitive control necessary to man-
 24 age the conflict (Hommel, 2011; MacLeod, 1991; Ridderinkhof, Wylie, van den Wildenberg, & et
 25 al., 2021).

26 Conflict effects are highly robust at the population-level – in fact, the Stroop effect is even
 27 considered universal (Haaf & Rouder, 2017; MacLeod, 1991). At the individual-level, however, the
 28 picture is more complicated. While almost everyone seems to show a conflict effect, true individual
 29 differences are assessed with a lot of uncertainty and are masked by measurement error. This sit-
 30 uation has been labeled the "reliability paradox" (Hedge, Powell, & Sumner, 2018): experimental
 31 control, which has the desirable effect of increasing validity, also tends to reduce individual differ-
 32 ences, and hence decreases reliability (Keye, Wilhelm, Oberauer, & et al., 2009; Paap & Greenberg,
 33 2013; Pettigrew & Martin, 2014; Rey-Mermet, Gade, & Oberauer, 2018).

34 To better understand reliability in experimental tasks, consider typical Stroop task data
 35 from I participants with two conditions (i.e., congruent and incongruent trials), and K trials
 36 per condition. For illustrative purposes, we use the simplest model possible, the normal-normal
 37 model. Individual i 's mean RT differences between the conditions follow a normal population-level
 38 distribution,

$$\bar{Y}_{i\text{inc}} - \bar{Y}_{i\text{con}} \sim \text{Normal}(\mu_\theta, \sigma_\theta^2 + \frac{2\sigma^2}{K}),$$

39 with mean μ_θ and variance $\sigma_\theta^2 + \frac{2\sigma^2}{K}$, where μ_θ is the population mean of the conflict effect, σ_θ the
 40 true individual variation and σ the measurement error¹. By "true individual variation", we mean
 41 the extent to which people differ quantitatively in their conflict effect if we had an infinite number
 42 of observations. By measurement error, we mean the confounding variability that exists because
 43 participants' RTs vary across trials, for example due to motor and perceptual processes. Using
 44 hierarchical modeling, one can estimate both σ_θ and σ as well as the population mean μ_θ .

45 Reliability is commonly defined as the ratio between true between-subjects variance σ_θ^2 and
 46 total variance, where the latter is the sum of true between-subjects variance in the effect of interest
 47 and an error term:

$$r = \frac{\sigma_\theta^2}{\sigma_\theta^2 + \sigma_E^2}.$$

¹In this paper, we use the terms "variability" and "variation" interchangeably to refer to both variance and standard deviation (i.e., the square root of the variance)

MCD, AH, and DM are supported by a Vidi grant to DM (VI.Vidi.191.091) from the Dutch Research Council (NWO). UB is supported by a Veni grant (VI.Veni.201G.045) from NWO. JMH is supported by a Veni grant (VI.Veni.201G.019) from NWO. AH is supported by the Australia-US Multidisciplinary University Research Initiative (AUSMURIV000003). This work used the Dutch national e-infrastructure with the support of the SURF Cooperative using grant no. EINF-5776.

The authors have no relevant financial or non-financial interests to disclose.

We would like to thank Suzanne Hoogeveen, Michael Nunez, Jeff Rouder, and Niek Stevenson for valuable discussions.

Correspondence concerning this article: m.c.donzallaz@uva.nl.

48 In experimental tasks with two main conditions of interest, such as the Stroop task, the error
 49 term is twice the squared standard error (see Rouder & Haaf, 2019, for more details),

$$\sigma_E^2 = \frac{2\sigma^2}{K},$$

50 where K is the number of trials performed per condition. Therefore, the reliability term becomes

$$r = \frac{\sigma_\theta^2}{\sigma_\theta^2 + \frac{2\sigma^2}{K}}.$$

51 Rouder and Haaf (2019) pointed out that to understand the reliability paradox it is important
 52 to acknowledge that reliability is not "portable": even for otherwise identical tasks, reliability differs
 53 with K . This is due the fact that as K grows, the standard error decreases and the reliability
 54 increases. Rouder, Kumar, and Haaf (2019) proposed that the suitability of a conflict task for
 55 individual-differences research instead be evaluated in terms of the signal-to-noise ratio (SNR)

$$\gamma = \frac{\sigma_\theta}{\sigma},$$

56 where the signal is the individual variability of the conflict effect, σ_θ , and the noise is the trial-to-
 57 trial variability in RT, σ . Using hierarchical modeling Rouder et al. (2019) found that the ratio, γ ,
 58 was around 1/7 across a large set of conflict tasks, leading them to conclude that for conventional
 59 conflict tasks, $K > 500$ (i.e., over 1000 trials per participant) is required to obtain sufficiently
 60 reliable measurement.

61 Low reliability, and the impracticality of collecting that many trials, has spurred research
 62 to develop more reliable ways of measuring cognitive control in general and cognitive conflict in
 63 particular. One direction has been to abandon conflict measures based on RT differences, as taking
 64 a difference doubles the measurement error (Draheim, Mashburn, Martin, & Engle, 2019). Instead,
 65 it has been suggested that cognitive control should be assessed using accuracy-based, adaptive
 66 tasks or the number of correct (in)congruent trials performed within a fixed time frame, with
 67 confounding from overall speed partialled out. However, some of these approaches sacrifice the
 68 validity of RT difference scores, which is evidenced by a large body of experimental literature, and
 69 risk confounding from individual differences in processing speed (Kucina et al., 2023; Rouder et al.,
 70 2019).

71 Kucina et al. (2023) used an alternative approach to improve reliability. They used the SNR
 72 γ to compare, combine, and refine traditional conflict tasks. They also refined the RT difference
 73 measure itself and took into account that RT distributions are well described by a shifted log-
 74 normal distribution (Heathcote & Love, 2012), improving the descriptive adequacy of the model
 75 used. They reported SNRs for conflict effects in the range of 1/3 to 1/2, corresponding to clas-
 76 sic reliability values of 0.8 or higher, with only 100 trials. In related work, Haines et al. (2020)
 77 found that modeling RT as lognormally distributed improved test-retest reliability in several con-
 78 flict datasets (Hedge, Powell, & Sumner, 2018), a delay-discounting task (Gawronski, Morrison,
 79 Phills, & Galdi, 2017), and an Implicit Association Test (Ahn et al., 2020). More generally, Haines
 80 et al. (2020) argued that simple atheoretical summaries, such as conflict effects based on mean
 81 RT, are unable to validly characterize the underlying psychological processes and recommended a
 82 "generative modeling approach" instead to overcome a "theory-description gap". Consistent with
 83 this recommendation, Hedge, Powell, Bompas, Vivian-Griffiths, and Sumner (2018) demonstrated
 84 that it can be important to simultaneously account for accuracy and RT because inconsistent re-
 85 lationships between accuracy-based and RT-based effects are widespread. Individual differences in
 86 this speed-accuracy trade-off may potentially contribute to the unreliability of RT-based conflict

87 effects. Evidence-accumulation models can disentangle these confounding processes (Hedge, Powell,
88 Bompas, et al., 2018).

89 In summary, the reliability paradox has triggered a wide range of responses in the commu-
90 nity. The proposed solutions range from refining and reinventing traditional tasks and measures,
91 to advocating for theoretically-motivated modeling approaches that provide psychologically mean-
92 ingful measures of cognitive conflict and a more complete and descriptively more accurate account
93 of choice RT performance. Here we investigate the degree to which different modeling choices
94 influence the SNR.

95 **Comparing the Reliability of Model-Based Conflict Effects**

96 Does improving the descriptive adequacy of the model from which conflict effects are derived
97 improve the SNR? Can a cognitive model that provides a psychologically meaningful characteriza-
98 tion of performance detect relatively more signal in a given dataset than a purely statistical model?
99 And to what extent does this depend on the dataset? We will use the SNR of conflict effects as
100 measured by variants of the Stroop task to answer these questions.

101 Rouder and Mehrvarz (2024) suggested that the SNR can be used as an indicator for reliability
102 that is independent of the trial size. To see the relationship between reliability and SNR, note that
103 reliability can be re-expressed in terms of the ratio γ ,

$$r = \frac{\gamma^2}{\gamma^2 + \frac{2}{K}}.$$

104 In theory, the independence of the number of trials makes γ a portable indicator of reliability
105 (Rouder & Haaf, 2019). In practice, however, the SNR does not always prove to be portable.
106 Kucina et al. (2023) modeled the possibility that individual differences in conflict effect and trial-
107 to-trial variability changes with K and computed the SNR as a function of the number of trials.
108 After de-trending to account for decreased mean RT with practice (Evans, Brown, Mewhort, &
109 Heathcote, 2018; Heathcote, Brown, & Mewhort, 2000), their results indicated that the between-
110 subjects variability σ_θ decreased as more trials were performed (i.e., individual differences decreased
111 with practice). Although this result raises doubts about the portability of γ , they showed that the
112 SNR can be useful for comparing reliability among tasks for a fixed K . Whether the SNR is portable
113 or not, however, does not affect its usefulness for our current purpose: comparing reliability among
114 models of the same task for a fixed K .

115 We investigate the effects of closing the “theory-description gap” (Haines et al., 2020) on
116 the SNR by progressively improving the descriptive adequacy and theoretical underpinning of the
117 models used to quantify the conflict effect. We examine how reliably various models can detect
118 individual differences in four Stroop datasets: we move from purely statistical models, such as
119 the normal or (shifted) lognormal distribution that provide a descriptive characterization of RT
120 distributions, to cognitive models, such as the lognormal race and the racing diffusion models that
121 simultaneously account for choice accuracy and RTs and that provide a psychologically meaningful
122 account of performance. We show how the SNR can be computed analytically for these five models
123 and outline a simulation-based approach that is more broadly applicable to any model².

124 The main argument for more complex statistical models of the conflict effect is that de-
125 scriptive adequacy may increase the SNR (Haines et al., 2020). A better fit means more accurate
126 measurement of individual differences and the noise term, the components of the SNR. Additionally,
127 cognitive models may further increase the SNR. By adding psychologically meaningful parameters,
128 some variability that is attributed to noise in statistical models may be explained by cognitive

²Any model with unbiased (e.g., maximum likelihood) estimators

129 models. Therefore, cognitive modeling may increase the SNR by decreasing its denominator, the
130 noise term, which is critical for reliable assessment of individual differences in Stroop effects.

131

Method

132 From the Normal Model to the Racing Diffusion Model

133 The five models employed are illustrated in Figure 1. The first three are purely RT-based
134 Bayesian hierarchical models that were also used by Haines et al. (2020). They are shown in in-
135 creasing order of descriptive adequacy: a normal model, which makes the unrealistic assumption
136 that RTs are normally distributed and assigns probability to negative RTs, followed by two models
137 that account for the characteristic right skew of empirical RT distribution by assuming a lognormal
138 parametrization, one with and one without an estimated shift parameter. The shift parameter en-
139 sures that the lower bound of the RT distribution is shifted away from zero, excluding unrealistically
140 fast RTs.

141 We then use two psychologically-grounded evidence-accumulation models, the lognormal race
142 model (LNR; Heathcote & Love, 2012) and the racing diffusion model (RDM; Tillman, Van Zandt,
143 & Logan, 2020), both providing a comprehensive account of performance by simultaneously ac-
144 counting for response choices and the corresponding RTs (see also Matzke, Logan, & Heathcote,
145 2020). Racing evidence-accumulation models have a long history in psychology because they pro-
146 vide a principled yet flexible approach to describe and explain performance in a broad range of
147 tasks (Heathcote & Matzke, 2022). Both the LNR and the RDM assume that decisions are made
148 by a process of accumulating evidence until a threshold amount is reached, which then triggers a
149 response. In particular, the models assume an independent race between a set of evidence accumu-
150 lators, corresponding to the different response options. In terms of the color-Stroop task, this means
151 separate racers for each response option (e.g., left button for "green" and right button for "red")
152 that race against each other. The racer that wins determines the response. Individuals may start
153 accumulating evidence at different start points and/or set different thresholds to vary the trade-off
154 between the speed of their responses and the accuracy of their choices (Donkin & Brown, 2018).
155 For example, a smaller distance from the starting point of evidence accumulation to the response
156 threshold leads to faster but less accurate responses. Conversely, larger distances lead to slower
157 but more accurate responses. Similar to the shifted lognormal model, the cognitive models also
158 feature a "non-decision time" parameter that accounts for time spent to encode evidence from the
159 choice stimulus and to produce a response. This parameter shifts the finishing time distributions
160 and ensures that the lower bound of the RT distributions is greater than zero.

161 The LNR (Heathcote & Love, 2012) was chosen because it forms a natural extension of the
162 RT-only models by assuming a lognormal parameterization but also accounting for choice accuracy.
163 The model assumes that the rate of evidence accumulation and the distance between the start
164 point and threshold follow independent lognormal distributions. As a result, the time for a single
165 accumulator to reach its threshold (i.e., the finishing time) also follows a lognormal distribution.
166 When accuracy is perfect, the observed RT distribution follows a shifted lognormal distribution,
167 consistent with the use of this distribution by Haines et al. (2020). However, when performance is
168 not perfect, as is typically the case for Stroop data (see Table 1 for average error rates in the four
169 data sets we examine), the only RT-based shifted lognormal model is mis-specified. Here we use
170 the LNR model to assess whether the SNR improves when accuracy is also taken into account.

171 A limitation of the LNR model is that the rate of evidence accumulation and the distance
172 from starting point to threshold are not separately identified because they combine linearly. As
173 a result, the conflict effect of the LNR confounds individual differences in response cuation and
174 the effects of congruency on the rate of evidence accumulation. The RDM (Tillman et al., 2020)

175 addresses this challenge to validity by allowing us to separately estimate evidence-accumulation
 176 rates and response thresholds.

177 The RDM (Tillman et al., 2020) assumes that each of the accumulators is a Wiener diffusion
 178 process with an evidence-accumulation rate v , a starting point of 0, and a response threshold β .
 179 Trial-to-trial variability in behavior is caused by stochastic accumulation, that is, during accu-
 180 mulation, the evidence-accumulation rate changes from moment to moment due to the addition
 181 of fractional Gaussian noise. The finishing time distribution of each accumulator is an inverse
 182 Gaussian distribution (i.e., Wald distribution; Wald, 1947). Here we use the RDM to investigate
 183 whether using psychologically-grounded measures of cognitive conflict as captured in the evidence-
 184 accumulation rate parameter improves the SNR.

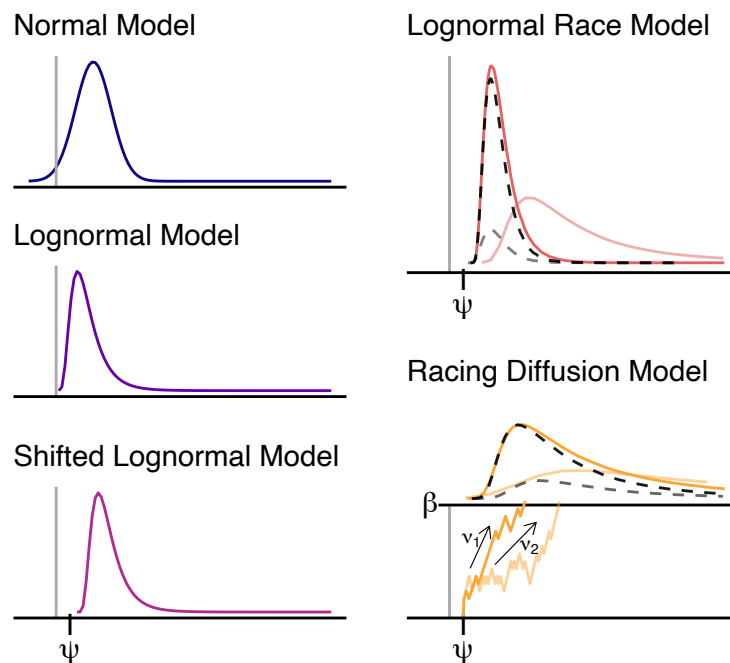


Figure 1. Model overview. The plots show exemplary shapes of the corresponding response and finishing time distributions as assumed by the respective models. The lightgrey vertical line indicates response time (RT) = 0 on the x-axis. The **normal model** crosses this line, demonstrating that the normal model also allows for negative RTs. The **lognormal model** accounts for the fact that RTs can only be positive and that their distribution is typically skewed. The **shifted lognormal model** additionally accounts for the fact that RT distributions are shifted away from zero, whereby ψ indicates the shift parameter. The **lognormal race model** (LNR) additionally accounts for choice errors. The model conceptualizes decision-making as an independent race between evidence accumulators, one for each response option, until a threshold is reached which then triggers a response. Depicted are the finishing time distributions of the matching and mismatching accumulators (i.e., stimulus and response (mis)match). The dashed lines show observed distributions of the respective racers (i.e., those that won the race), scaled by their winning proportions. ψ indicates the non-decision time parameter. The **racing diffusion model** (RDM) separately estimates the evidence-accumulation rate and the threshold parameters. The yellow paths depict an exemplary race between the matching and mismatching accumulators. The matching accumulator (darker yellow) reaches the threshold β first, resulting in a correct response on that trial. ν_1 and ν_2 refer to the matching and mismatching evidence-accumulation rate parameters.

185 We deliberately chose to explore the performance of relatively standard evidence-
 186 accumulation models as opposed to more complex models specifically developed for conflict tasks.
 187 Although conflict models such as the spotlight diffusion model (White, Ratcliff, & Starns, 2011)
 188 align more closely with the cognitive processes assumed to underlie performance in these tasks,
 189 they are limited in terms of tractability: the models grapple with pronounced parameter trade-offs
 190 and hierarchical estimation is not straightforward. As a result, despite their theoretical merits,
 191 they are less suitable as measurement models to examine individual differences in conflict tasks.

192 The full Bayesian hierarchical model specification can be found in the supplementary mate-
 193 rials. The three purely RT-based models, the normal, the lognormal, and the shifted lognormal
 194 models, are parameterized such that the (log) means are decomposed into an intercept α_i and a
 195 conflict or congruency effect θ_i for individual i . We allow the (log) standard deviation to vary
 196 across individuals, σ_i , but not across conditions (i.e., congruent vs. incongruent). For the shifted
 197 lognormal model, we also estimate individual shift parameters ψ_i which remain the same across
 198 conditions. The LNR is parameterized such that the log mean of the accumulator that *matches*
 199 the stimulus is decomposed into an intercept α_{1i} and a conflict effect θ_i . The corresponding log
 200 standard deviation, σ_{1i} , is allowed to vary across individuals but again not across conditions. Due
 201 to high accuracy rates in our analyzed datasets (see Table 1), we fix the log mean and log standard
 202 deviation of the accumulator that mismatches the stimulus across individuals and estimate an in-
 203 tercept only and no conflict effect (i.e., we estimate α_2 and σ_2). Just like in the shifted lognormal
 204 model, we also estimate a non-decision time parameter for each individual, ψ_i . Lastly, for the RDM,
 205 we decompose the matching evidence-accumulation rate into an intercept α_{1i} and a conflict effect
 206 θ_i . Note that we only model a conflict effect on the rate because we expect that the cognitive con-
 207 flict will impact the evidence accumulation rate at which participants respond to the (in)congruent
 208 stimuli. In contrast, the individual response thresholds are set at the beginning of a trial and we
 209 do not expect stimulus encoding or motor control processed to differ across conditions. Similar to
 210 the LNR specification, we estimate an intercept only and do not allow for individual differences in
 211 the mismatching evidence-accumulation rate, α_2 (see Lüken, Heathcote, Haaf, & Matzke, 2023).
 212 Furthermore, we estimate individual response thresholds, β_i , and non-decision time parameters, ψ_i ,
 213 which were fixed across congruency conditions and accumulators.

214 Estimating Signal-to-Noise Ratios

215 All models have a designated parameter reflecting the numerator of the SNR, σ_θ (i.e., in-
 216 dividual differences in the conflict effect). However, this is not the case for its denominator. To
 217 determine the trial-by-trial variation, we estimate the standard error of the conflict effect. Note
 218 that in the case of the normal model, the standard error is simply $\frac{2\sigma^2}{K}$ and that σ can be directly
 219 estimated using hierarchical modeling and interpreted on the real line (see Rouder & Haaf, 2019;
 220 Rouder et al., 2019).

221 For the normal model, we could use the formula provided by Rouder and Mehrvarz (2024)
 222 to obtain a point estimate of γ from RT data. However, we prefer to use Bayesian hierarchical
 223 modeling (Lee, 2011) for all models for several reasons. First, hierarchical models partition the
 224 variability in data into between and within-participant components. Second, Bayesian estimation
 225 naturally provides measures of uncertainty in γ estimates, providing a basis for inference about
 226 the SNR. Third, Bayesian hierarchical models are well suited for estimating the parameters of
 227 evidence-accumulation models such as the LNR and RDM, as well as the shift parameter of the
 228 lognormal distribution when it is estimated (e.g., Heathcote et al., 2019; Stevenson et al., 2024).

229 Hierarchical modeling typically directly provides an estimate of the between-subjects variabil-
 230 ity of interest, σ_θ for all models. The models just need to have parameters reflecting the difference

231 between the congruent and incongruent conditions (i.e., θ_i). For all but the RDM, the individ-
 232 ual noise parameters are also directly provided (i.e., σ_i for the normal and the three lognormal
 233 models): σ_i , the shape parameter of the lognormal distribution and standard deviation of the log-
 234 RTs, directly reflects the individual trial-by-trial variation. Moreover, within models, σ_θ and σ_i
 235 are on the same scale (i.e., on the real line in the normal model and on the log scale in the three
 236 lognormal models). Therefore, the units cancel out, making the ratios comparable to each other.
 237 Since we assume individual noise parameters, we need to average estimates of σ_i across individuals,
 238 $i = 1, \dots, I$:

$$\gamma = \frac{\hat{\sigma}_\theta}{\sqrt{\frac{1}{I} \sum_{i=1}^I \hat{\sigma}_i^2}} \quad (1)$$

239 Finding an expression for the noise term σ_i in the RDM is more complicated. Contrary to the
 240 other models, trial-by-trial variation is not explicitly represented as a model parameter. Besides the
 241 hierarchical structure that we place on top of the model parameters, our specification of the RDM
 242 has one inherent source of variability: within-trial variation of the drift rate, which is described
 243 by a Wiener diffusion process (Tillman et al., 2020). However, this parameter reflects variation
 244 *within* trials and is conventionally set to 1 to make the model identifiable. To get an estimate of
 245 between-trial variation or measurement error, the standard error of the conflict effect needs to be
 246 determined. Let $\hat{\nu}_{i_{con}}$ and $\hat{\nu}_{i_{inc}}$ be individual i 's matching evidence-accumulation rate estimates in
 247 the congruent and incongruent condition such that $\Delta_{\hat{\nu}_i} = \hat{\nu}_{i_{con}} - \hat{\nu}_{i_{inc}}$ represents the conflict effect.
 248 To get an expression for the noise, we need to compute the variance $V(\Delta_{\hat{\nu}_i})$. An expression for this
 249 variance can be derived as follows:

$$\begin{aligned} V(\hat{\nu}_{i_{con}} - \hat{\nu}_{i_{inc}}) &= V(\hat{\nu}_{i_{con}}) + V(\hat{\nu}_{i_{inc}}) - 2\text{Cov}(\hat{\nu}_{i_{con}}, \hat{\nu}_{i_{inc}}) \\ V(\hat{\nu}) &= (\text{E}(\hat{\nu}))^2 - \text{E}(\hat{\nu}^2) \\ \text{E}(\hat{\nu}) &= \sqrt{K} \frac{\Gamma(K - 3/2)}{\Gamma(K - 1)} \left(\frac{1}{K\beta} + v \right) \\ \text{E}(\hat{\nu}^2) &= \frac{K}{K - 2} \left(v^2 + \frac{3}{K^2\beta^2} + \frac{3v}{K\beta} \right) \\ \text{E}(\hat{\nu}_1\hat{\nu}_2) &= \frac{K}{K - 2} \left(\frac{1}{K\beta} + \nu_1 \right) \left(\frac{1}{K\beta} + \nu_2 \right) \\ \text{Cov}(\hat{\nu}_1, \hat{\nu}_2) &= \text{E}(\hat{\nu}_1\hat{\nu}_2) - \text{E}(\hat{\nu}_1)\text{E}(\hat{\nu}_2) \end{aligned}$$

250 Note that Γ is the gamma function. The proof is provided in the supplementary materials.
 251 This effect variance reflects the standard error $\frac{\sqrt{2}\sigma_i}{\sqrt{K_i}}$, which can be transformed to get an estimate
 252 of the individual noise terms σ_i : $\sqrt{\frac{V(\Delta_{\hat{\nu}_i})K_i}{2}}$ where K_i is the number of trials per condition that
 253 individual i performed. Averaging across individuals, the SNR for the RDM is:

$$\gamma = \frac{\hat{\sigma}_\theta}{\sqrt{\frac{1}{I} \sum_{i=1}^I \frac{V(\Delta_{\hat{\nu}_i})K_i}{2}}} \quad (2)$$

254 In order to capture uncertainty in γ estimates, we calculate the ratio for all available posterior
 255 samples of the individual RDM parameters.

256 The standard error of the effect of interest within a model can also be computed using a
 257 simulation-based approach. The details are provided in the supplementary materials. The ad-
 258 vantage of the simulation-based approach is that it is more broadly applicable than an analytical
 259 expression which is specific to the particular model. However, the simulation-based approach

Table 1

Descriptive statistics (mean RT and accuracy), number of participants (I), number of trials per congruency condition (K), and task type.

	<i>RT (SD)</i>	<i>Accuracy</i>	<i>I</i>	<i>K</i>	<i>Task type</i>
Enkavi et al. (2019)	.73 (.21)	.96	522	48	Color Stroop
Pratte et al. (2010)	.77 (.37)	.96	38	164	Color Stroop
Rey-Mermet et al. (2018)	.59 (.15)	.96	128	95	Number Stroop
Von Bastian et al. (2016)	.75 (.25)	.97	121	48	Number Stroop

260 requires an unbiased maximum-likelihood estimator for the quantity of interest, which may not
 261 always be available. For the RDM, this quantity is the evidence-accumulation rate, which has such
 262 a maximum-likelihood estimator (i.e., the natural estimator; see supplementary materials).³

263

Results

264 We applied the five models to four datasets: number-Stroop tasks by (Rey-Mermet et al.,
 265 2018, younger age group) and by Von Bastian, Souza, and Gade (2016), color-Stroop tasks by
 266 Pratte, Rouder, Morey, and Feng (2010) and by Enkavi et al. (2019). We only analyzed the
 267 (in)congruent conditions and retrieved all datasets from the preprocessed database made available
 268 by Haaf, Hoffstadt, and Lesche (2024).

269 Before estimation, we removed the first three trials of every block, neutral trials, and RTs
 270 ≤ 0.25 seconds. Two studies (Enkavi et al., 2019; Rey-Mermet et al., 2018), used response
 271 windows to ensure fast responding. This approach led to somewhat bimodal RT distributions
 272 and convergence problems for some models. We therefore removed very slow responses, specifically
 273 responses that were at the upper limit of the response window (Rey-Mermet et al., 2018: RT
 274 ≥ 1.9 s; Enkavi et al., 2019: RT > 1.5 s). Pratte et al. (2010) did not have a response
 275 window, but the very long RTs caused some convergence problems, so we excluded them (RT
 276 ≥ 3 s). We also excluded participants with fewer than 50% remaining trials. Overall, only
 277 few trials and participants were excluded (von Bastian et al., 2016: 0.9%, Rey-Mermet et al.,
 278 2018: 5% of the observations including all observations from one participant due to having
 279 less than 50% trials after data cleaning; Pratte et al., 2010: 5%; Enkavi et al., 2019: 4% of the
 280 observations including all observations from one participant due to having less than 50% trials after
 281 data cleaning). For the fitting of the (shifted log)normal models, we also removed all incorrect trials.

282

283 We estimated all models using stan(Stan Development Team, 2023)⁴ using weakly informa-
 284 tive priors based on prior predictive checks and known plausible ranges of RTs in seconds⁵. The
 285 descriptive adequacy of the normal model was consistently worse than for the other four models

³There exists also a straightforward method-of-moments estimator for the evidence-accumulation rate, but this is not a maximum-likelihood estimator and is biased (see the supplementary materials).

⁴We used the following R packages: R (Version 4.2.1; R Core Team, 2022) and the R-packages *DBI* (Version 1.1.3; R Special Interest Group on Databases (R-SIG-DB), Wickham, & Müller, 2022), *dplyr* (Version 1.1.3; Wickham, François, Henry, Müller, & Vaughan, 2023), *extraDistr* (Version 1.9.1; Wolodzko, 2020), *ggplot2* (Version 3.4.3; Wickham, 2016), *papaja* (Version 0.1.1; Aust & Barth, 2022), *patchwork* (Version 1.1.3; Pedersen, 2023), *RSQLite* (Version 2.3.1; Müller, Wickham, James, & Falcon, 2023), *rstan* (Version 2.21.5; Stan Development Team, 2022), *scales* (Version 1.2.1; Wickham & Seidel, 2022), *StanHeaders* (Version 2.21.0.7; Stan Development Team, 2020), *stringr* (Version 1.5.0; Wickham, 2022), *tidybayes* (Version 3.0.2; Kay, 2022), *tidyr* (Version 1.3.0; Wickham, Vaughan, & Girlich, 2023), and *tinylabels* (Version 0.2.3; Barth, 2022).

⁵See the supplementary materials for the prior specification, estimation results, and the detailed posterior predictive assessment of the absolute descriptive adequacy of the models (i.e., the match between the observed response time distributions and the posterior predictions).

286 because the normal model cannot account for the characteristic skewness of RT data (see Figure
 287 2). The shifted lognormal, the LNR, and the RDM showed good descriptive adequacy, however the
 288 RDM and the LNR both had some difficulties capturing slow errors in the incongruent condition,
 289 particularly for the dataset by (Rey-Mermet et al., 2018, see the Appendix and supplementary
 290 materials). To quantify descriptive adequacy, we computed the root mean squared error (RMSE)
 291 for the correct RTs and in the case of the LNR and RDM, additionally for the proportion of correct
 292 responses⁶. We simulated 500 datasets from the posterior predictive distributions and computed
 293 the RMSE for each of the simulated datasets and several quantiles (see Table 2 and 3 in the Ap-
 294 pendix). For most datasets, RMSE steadily decreased for the first three models and the RMSE for
 295 the LNR tended to be smaller than for the RDM.

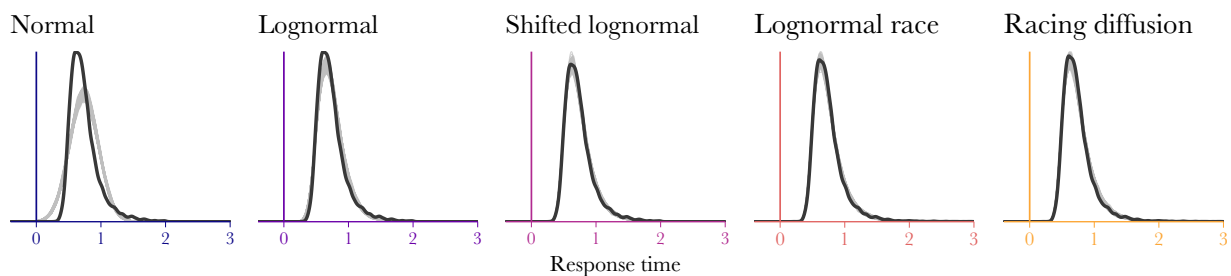


Figure 2. Observed (black) and predicted (grey) correct response time distributions collapsed across participants and conditions for the dataset by Von Bastian et al. (2016). The grey lines represent 500 samples from the posterior predictive distribution. The lognormal race and the racing diffusion model distributions also contain the incorrect responses.

296 We then computed the SNRs by plugging the parameter posterior samples into Equations (1)
 297 and (2) (see Figure 3). Note that the simulation-based approach leads to approximately the same
 298 result⁷. Within datasets, the credible intervals overlap substantially, suggesting that the choice of
 299 the model does not substantially influence the ratio in the analyzed datasets. The only exception
 300 is the Von Bastian et al. (2016) dataset, where the posterior median increases from below $\frac{1}{10}$ for
 301 the normal model to around $\frac{1}{5}$ for the RDM. Notably, the SNRs of the dataset by Enkavi et al.
 302 (2019) are considerably higher than those of the other datasets.

⁶We did not use model selection criteria such as the deviance information criterion (DIC) due to the fact that the models are not comparable as they were fitted to different data (i.e., the (shifted log)normal models were fitted to correct responses only).

⁷see the supplementary materials for the comparisons

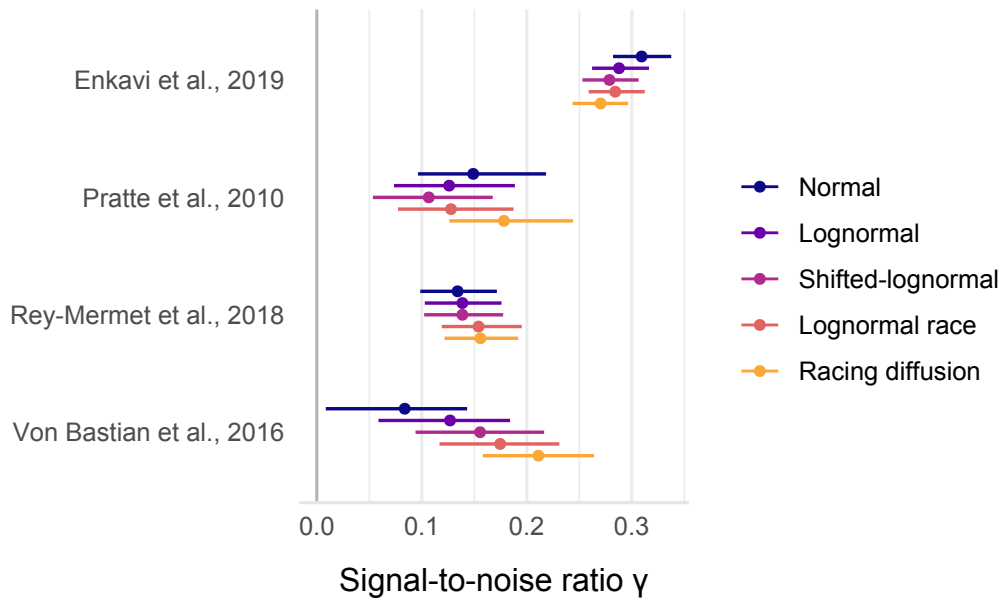


Figure 3. Posterior medians and 95% credible interval of the signal-to-noise ratios γ as computed by the analytical method.

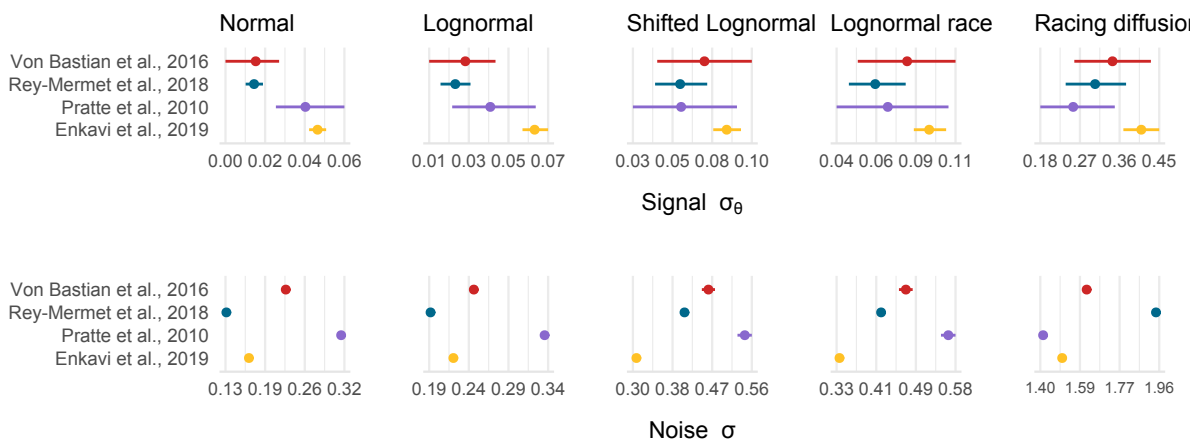


Figure 4. Posterior medians and 95% credible intervals of the numerator (i.e., signal) and denominator (i.e., noise) of the signal-to-noise ratios.

303 Looking at the components of the SNRs separately (Figure 4), the ratio differences across
 304 datasets and within models seem to be driven by both the extent of individual differences in the
 305 Stroop effect (i.e., signal) and by trial-by-trial variation (i.e., noise). Note that the credible interval
 306 (CrI) of the noise term are narrower than those of the signal term because the former is estimated
 307 from all observations. In contrast, the CrI of σ_θ differ in width across datasets as more participants
 308 leads to more precise estimation (see Table 1). Across models, the order of the terms stays largely
 309 the same – with two exceptions. (1) Moving from the normal model to the RDM for the data by
 310 Von Bastian et al. (2016), we gain relatively more signal, though note that the scales of the signal
 311 and the noise terms are not directly comparable. Within the components of the other datasets,

312 there does not appear to be a dominating signal or noise term across models, which is not surprising
313 given that the CrI of the SNRs are overlapping.

314

Discussion

315 Experimental tasks that produce robust group-level differences can simultaneously produce
316 unreliable individual differences. Proposed solutions to this “reliability paradox” (Hedge, Powell,
317 & Sumner, 2018) include collecting more trials per participant, modifying the cognitive tasks,
318 using more descriptively adequate models, and using psychologically-motivated measures derived
319 from cognitive models to quantify individual differences in the construct of interest. Here we
320 focused on the last two solutions, and built on the work by Rouder and Mehrvarz (2024) and
321 Rouder et al. (2019) by using signal-to-noise ratios (SNRs) as a tool for comparing models of
322 varying sophistication. The SNR is a useful measure of reliability quantifying how well individual
323 differences can be detected relative to measurement error. We provided analytical solutions to
324 compute the ratio for five models from the output of Bayesian hierarchical modeling: the normal,
325 lognormal, shifted lognormal, lognormal race (LNR), and racing diffusion models (RDM). Moreover,
326 we provided a general algorithm that - whenever an unbiased estimator of the quantity of interest
327 is available - can be used to compute the SNR for any model. We then applied the five models to
328 four Stroop datasets.

329 Our analysis showed that modeling choices do not have a consistent effect on the SNR: the
330 proposed solutions improved reliability in only one of the four datasets, and we found no consistent
331 ordering or pattern as to which model yielded the highest SNR across datasets. Notably, improved
332 descriptive accuracy did not correspond to a higher SNR. The normal model is clearly unable
333 to accommodate the slow tail of typical RT distributions, yet the corresponding SNRs are not
334 consistently worse than those of the RDM or lognormal models that can account for the skewness
335 in observed RTs. These differences did not appear to be explained by an effect of trial number or
336 sample size. Perhaps a more systematic assessment of more datasets might help identify dataset
337 or task characteristics that are predictive of higher SNRs for one model over another. As of now,
338 the preferred model can only be established after fitting all models using hierarchical modeling and
339 computing the SNRs.

340 In our analyses, we only looked at a subset of potential models. For example, other relevant
341 models are the diffusion decision model (Ratcliff & McKoon, 2008; Ratcliff & Smith, 2004) or
342 the linear ballistic accumulator model (Heathcote & Love, 2012). However, our approach can be
343 extended to these models. In both cases the trial-by-trial variation is not directly reflected in
344 a model parameter. Yet, there are two ways of computing the measurement error term of the
345 SNR: (1) one can try to derive it analytically by computing the variance of the effect of interest
346 or standard error, as we have done for the RDM. (2) one can use a simulation-based approach if
347 unbiased maximum likelihood estimators are available.

348 The focus of this paper is the reliability of measures. However, the validity of a measure
349 (i.e., its ability to measure the intended construct) is also important. Cognitive model parameters
350 have a clear psychological interpretation and so intrinsically support valid inferences as long as the
351 model accurately represents the data and the process that generates it. So, even if a particular
352 model parameter might not lead to a higher SNR, it might still be a better measure. For instance,
353 suppose we analyze a dataset in which speed-accuracy trade-offs are present. Because the LNR or
354 RDM can account for such trade-offs by combining information about accuracy and speed they are
355 likely to provide more valid measures. However, estimating a model providing a valid psychological
356 account may also require particular qualities in the data and the design from which it comes. For
357 example, Lüken et al. (2023) showed that low error rates compromised the quality of estimates of

358 the parameters of the diffusion decision model and the linear ballistic accumulator. It is possible
359 that the low error rates in the datasets we analyzed here may have had the same effect on the
360 LNR and RDM, and that is the reason why they did not consistently perform better than the rest.

361 Clearly, both reliability and validity need to be taken into account whenever researchers
362 attempt to answer substantive questions about the nature of individual differences in cognitive
363 control. Descriptive adequacy is also important, as a model that clearly misfits the data is unlikely
364 to be valid. However, good fit alone does not ensure validity, the model must also provide a sensible
365 account of psychological processes that could plausibly generate the data, preferably one that is
366 backed up by converging evidence from prior literature (i.e., the same type of model has provided
367 accurate and coherent accounts of data from related tasks and manipulations).

368 In sum, the SNR can be used as a tool to identify the statistical or cognitive model that
369 is best suited to examine individual differences in conflict data. We have provided analytical and
370 simulation-based approaches to compute the SNR for a range of models of varying sophistication
371 and showed that models that provide a better, and potentially more valid, description do not
372 necessarily consistently improve the reliability with which individual differences in cognitive conflict
373 are measured. Hence, we recommend that our methodology be deployed on a case-by-case basis to
374 assess the effects of model choices on reliability.

375

Declarations

376 **Funding:** MCD, AH, and DM are supported by a Vidi grant to DM (VI.Vidi.191.091) from
377 the Dutch Research Council (NWO). UB is supported by a Veni grant (VI.Veni.201G.045) from
378 NWO. JMH is supported by a Veni grant (VI.Veni.201G.019) from NWO. AH is supported by
379 the Australia-US Multidisciplinary University Research Initiative (AUSMURIV000003). This work
380 used the Dutch national e-infrastructure with the support of the SURF Cooperative using grant
381 no. EINF-5776.

382 **Conflicts of interest/Competing interests:** The authors have no relevant financial or non-
383 financial interests to disclose.

384 **Ethics approval:** Not applicable

385 **Consent to participate:** Not applicable

386 **Consent for publication:** Not applicable

387 **Availability of data and materials:** The archival data used in this article is available at [https://](https://osf.io/fq8ep)
388 osf.io/fq8ep

389 **Code availability (software application or custom code):** The code is available at [https://](https://osf.io/fq8ep)
390 osf.io/fq8ep

391 **Authors' contributions:** MCD: conceptualization, methodology, formal analysis, writing and
392 editing, visualization; UB: methodology, writing and editing, visualization; AH: conceptualization,
393 methodology, writing and editing, supervision CD: conceptualization, editing, supervision; DM:
394 conceptualization, methodology, writing and editing, supervision, funding acquisition; JMH: con-
395 ceptualization, methodology, writing and editing, supervision.

References

396

- 397 Ahn, W.-Y., Gu, H., Shen, Y., Haines, N., Hahn, H. A., Teater, J. E., . . . Pitt, M. A. (2020). Rapid, precise,
 398 and reliable measurement of delay discounting using a Bayesian learning algorithm. *Scientific Reports*,
 399 *10*(1), 12091. doi: 10.1038/s41598-020-68587-x
- 400 Donkin, C., & Brown, S. D. (2018). Response times and decision-making. In E.-J. Wagenmakers (Ed.),
 401 *Stevens' handbook of experimental psychology and cognitive neuroscience: Vol. 5. methodology* (4th
 402 ed., pp. 349–382). John Wiley & Sons.
- 403 Draheim, C., Mashburn, C. A., Martin, J. D., & Engle, R. W. (2019). Reaction Time in Differential and
 404 Developmental Research: A Review and Commentary on the Problems and Alternatives. *Psychological*
 405 *Bulletin*, *145*(5), 508–535. doi: 10.1037/bul0000192
- 406 Enkavi, A. Z., Eisenberg, I. W., Bissett, P. G., Mazza, G. L., MacKinnon, D. P., Marsch, L. A., & Poldrack,
 407 R. A. (2019). Large-scale analysis of test–retest reliabilities of self-regulation measures. *Proceedings*
 408 *of the National Academy of Sciences*, *116*(12), 5472–5477. doi: 10.1073/pnas.1818430116
- 409 Eriksen, B. A., & Eriksen, C. W. (1974). Effects of noise letters upon the identification of a target letter in
 410 a nonsearch task. *Perception & Psychophysics*, *16*, 143–149. doi: 10.3758/BF03203267
- 411 Evans, N. J., Brown, S. D., Mewhort, D. J. K., & Heathcote, A. (2018). Refining the Law of Practice.
 412 *Psychological Review*, *125*(4), 592–605. doi: 10.1037/rev0000105
- 413 Gawronski, B., Morrison, M., Phillips, C. E., & Galdi, S. (2017). Temporal Stability of Implicit and
 414 Explicit Measures. *Personality and Social Psychology Bulletin*, *43*(3), 300–312. doi: 10.1177/
 415 0146167216684131
- 416 Haaf, J. M., Hoffstadt, M., & Lesche, S. (2024). *Attentional control data collection: A resource for efficient*
 417 *data reuse*. <https://doi.org/10.31234/osf.io/4evy6>.
- 418 Haaf, J. M., & Rouder, J. N. (2017). Developing constraint in bayesian mixed models. *Psychological Methods*,
 419 *22*(4), 779–798. doi: 10.1037/met0000156
- 420 Haines, N., Kvam, P. D., Irving, L. H., Smith, C., Beauchaine, T. P., Pitt, M. A., . . . Turner, B. (2020).
 421 Learning from the reliability paradox: How theoretically informed generative models can advance the
 422 social, behavioral, and brain sciences.
 423 doi: 10.31234/osf.io/xr7y3
- 424 Heathcote, A., Brown, S., & Mewhort, D. J. K. (2000). The power law repealed: The case for an exponential
 425 law of practice. *Psychonomic Bulletin & Review*, *7*(2), 185 – 207. doi: 10.3758/BF03212979
- 426 Heathcote, A., Lin, Y.-S., Reynolds, A., Strickland, L., Gretton, M., & Matzke, D. (2019). Dynamic models
 427 of choice. *Behavior Research Methods*, *51*, 961–985. doi: 10.3758/s13428-018-1067-y
- 428 Heathcote, A., & Love, J. (2012). Linear deterministic accumulator models of simple choice. *Frontiers in*
 429 *psychology*, *3*. doi: 10.3389/fpsyg.2012.00292/abstract
- 430 Heathcote, A., & Matzke, D. (2022). Winner takes all! what are race models, and why and how should
 431 psychologists use them? *Current Directions in Psychological Science*, *31*(5), 383–394. doi: 10.1177/
 432 09637214221095852
- 433 Hedge, C., Powell, G., Bompas, A., Vivian-Griffiths, S., & Sumner, P. (2018). Low and Variable Correlation
 434 Between Reaction Time Costs and Accuracy Costs Explained by Accumulation Models: Meta-Analysis
 435 and Simulations. *Psychological Bulletin*, *144*(11), 1200–1227. doi: 10.1037/bul0000164
- 436 Hedge, C., Powell, G., & Sumner, P. (2018). The reliability paradox: Why robust cognitive tasks do not
 437 produce reliable individual differences. *Behavior Research Methods*, *50*(3), 1166–1186. doi: 10.3758/
 438 s13428-017-0935-1
- 439 Hommel, B. (2011). The Simon effect as tool and heuristic. *Acta Psychologica*, *136*(2), 189–202. doi:
 440 10.1016/j.actpsy.2010.04.011
- 441 Keye, D., Wilhelm, O., Oberauer, K., & et al. (2009). Individual differences in conflict-monitoring: testing
 442 means and covariance hypothesis about the simon and the eriksen flanker task. *Psychological Research*,
 443 *73*, 762–776. doi: 10.1007/s00426-008-0188-9
- 444 Kucina, T., Wells, L., Lewis, I., Salas, K. d., Kohl, A., Palmer, M. A., . . . Heathcote, A. (2023). Calibration
 445 of cognitive tests to address the reliability paradox for decision-conflict tasks. *Nature Communications*,
 446 *14*(1), 2234. doi: 10.1038/s41467-023-37777-2
- 447 Lee, M. D. (2011). How cognitive modeling can benefit from hierarchical Bayesian models. *Journal of*
 448 *Mathematical Psychology*, *55*(1), 1–7. doi: 10.1016/j.jmp.2010.08.013

- 449 Lüken, M., Heathcote, A., Haaf, J. M., & Matzke, D. (2023, October 18). *Parameter identifiability in*
450 *evidence-accumulation models: The effect of error rates on the diffusion decision model and the linear*
451 *ballistic accumulator*. Retrieved from <https://doi.org/10.31234/osf.io/wsgnt>
- 452 MacLeod, C. M. (1991). Half a Century of Research on the Stroop Effect: An Integrative Review. *Psycho-*
453 *logical bulletin*, *109*(2), 163 – 203. doi: 10.1037/0033-2909.109.2.163
- 454 Matzke, D., Logan, G. D., & Heathcote, A. (2020). A cautionary note on evidence-accumulation models
455 of response inhibition in the stop-signal paradigm. *Computational Brain Behavior*, *3*, 269–288. doi:
456 10.1007/s42113-020-00075-x
- 457 Paap, K. R., & Greenberg, Z. I. (2013). There is no coherent evidence for a bilingual advantage in executive
458 processing. *Cognitive Psychology*, *66*(2), 232–258. doi: 10.1016/j.cogpsych.2012.12.002
- 459 Pettigrew, C., & Martin, R. C. (2014). Psychology and aging. *Psychology and Aging*, *29*(2), 187–204. doi:
460 10.1037/a0036085
- 461 Pratte, M. S., Rouder, J. N., Morey, R. D., & Feng, C. (2010). Exploring the differences in distributional
462 properties between stroop and simon effects using delta plots. *Attention, Perception, & Psychophysics*,
463 *72*, 2013–2025. doi: 10.3758/APP.72.7.2013
- 464 Ratcliff, R., & McKoon, G. (2008). The diffusion decision model: Theory and data for two-choice decision
465 tasks. *Neural Computation*, *20*(4), 873–922. doi: 10.1162/neco.2008.12-06-420
- 466 Ratcliff, R., & Smith, P. L. (2004). A comparison of sequential sampling models for two-choice reaction
467 time. *Psychological Review*, *111*(2), 333–367. doi: 10.1037/0033-295X.111.2.333
- 468 Rey-Mermet, A., Gade, M., & Oberauer, K. (2018). Should we stop thinking about inhibition? searching
469 for individual and age differences in inhibition ability. *Journal of Experimental Psychology: Learning,*
470 *Memory, and Cognition*, *44*(4), 501–526. doi: 10.1037/xlm0000450
- 471 Ridderinkhof, K. R., Wylie, S. A., van den Wildenberg, W. P. M., & et al. (2021). The arrow of time:
472 Advancing insights into action control from the arrow version of the eriksen flanker task. *Attention,*
473 *Perception, & Psychophysics*, *83*, 700–721. doi: 10.3758/s13414-020-02167-z
- 474 Rouder, J. N., & Haaf, J. M. (2019). A psychometrics of individual differences in experimental tasks.
475 *Psychonomic Bulletin & Review*, *26*(2), 452–467. doi: 10.3758/s13423-018-1558-y
- 476 Rouder, J. N., Kumar, A., & Haaf, J. M. (2019). Why Most Studies of Individual Differences With Inhibition
477 Tasks Are Bound To Fail.
478 doi: 10.31234/osf.io/3cjr5
- 479 Rouder, J. N., & Mehrvarz, M. (2024). Hierarchical-model insights for planning and interpreting individual-
480 difference studies of cognitive abilities. *Current Directions in Psychological Science*, *33*(2), 128–135.
481 doi: 10.1177/09637214231220923
- 482 Simon, J. R., & Rudell, A. P. (1967). Auditory s-r compatibility: The effect of an irrelevant cue on
483 information processing. *Journal of Applied Psychology*, *51*(3), 300–304. doi: 10.1037/h0020586
- 484 Stan Development Team. (2023). Stan modeling language users guide and reference manual [Computer
485 software manual]. Retrieved from <https://mc-stan.org>
- 486 Stevenson, N., Donzallaz, M. C., Innes, R. J., Forstmann, B., Matzke, D., & Heathcote, A. (2024, January
487 30). Emc2: An r package for cognitive models of choice. *PsyArXiv*. Retrieved from [https://doi.org/](https://doi.org/10.31234/osf.io/2e4dq)
488 [10.31234/osf.io/2e4dq](https://doi.org/10.31234/osf.io/2e4dq) doi: 10.31234/osf.io/2e4dq
- 489 Stroop, J. R. (1935). Studies of interference in serial verbal reactions. *Journal of Experimental Psychology*,
490 *18*(6), 643. doi: 10
- 491 Tillman, G., Van Zandt, T., & Logan, G. D. (2020). Sequential sampling models without random between-
492 trial variability: the racing diffusion model of speeded decision making. *Psychonomic Bulletin &*
493 *Review*, *27*, 911–936. doi: 10.3758/s13423-020-01719-6
- 494 Von Bastian, C. C., Souza, A. S., & Gade, M. (2016). No evidence for bilingual cognitive advantages: A test of
495 four hypotheses. *Journal of Experimental Psychology: General*, *145*(2), 246. doi: 10.1037/xge0000120
- 496 Wald, A. (1947). *Sequential analysis*. John Wiley & Sons.
- 497 White, C. N., Ratcliff, R., & Starns, J. J. (2011). Diffusion models of the flanker task: Discrete versus gradual
498 attentional selection. *Cognitive Psychology*, *63*(4), 210–238. doi: 10.1016/j.cogpsych.2011.08.001

499

Appendix

500 Below, we present plots of the model fits (observed and predicted cumulative distribution
 501 functions) and tables displaying the RMSE for all models and datasets. Additional figures assessing
 502 the fits can be found in the supplementary materials.

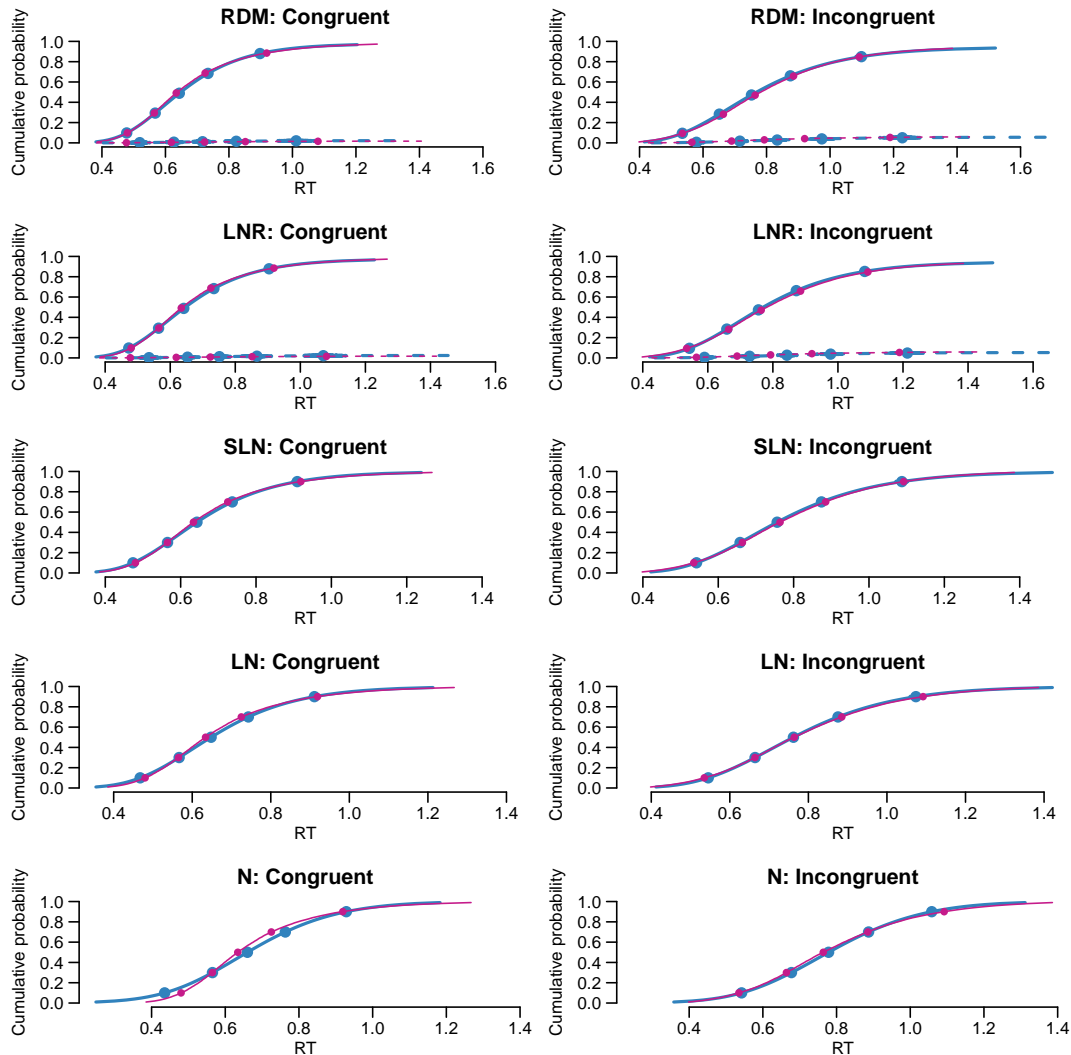
503 **Posterior predictive checks**

Figure 5. Observed and predicted cumulative distribution functions (CDFs) for Enkavi et al. (2019). Observed = pink, blue = predicted, RDM = racing diffusion model, LNR = lognormal race model, SLN = shifted lognormal model, LN = lognormal model, N = normal model. The predictions are based on 500 sampled datasets from the posterior predictive distribution and plotted is the median including 95% credible interval. The points show the 10%, 30%, 50%, 70%, and 90% quantiles averaged across participants, separately for the two congruency conditions. Note that for the RDM and the LNR, both the (defective) CDFs of correct and incorrect responses are shown, whereas for the other models, only the CDFs of correct responses are depicted.

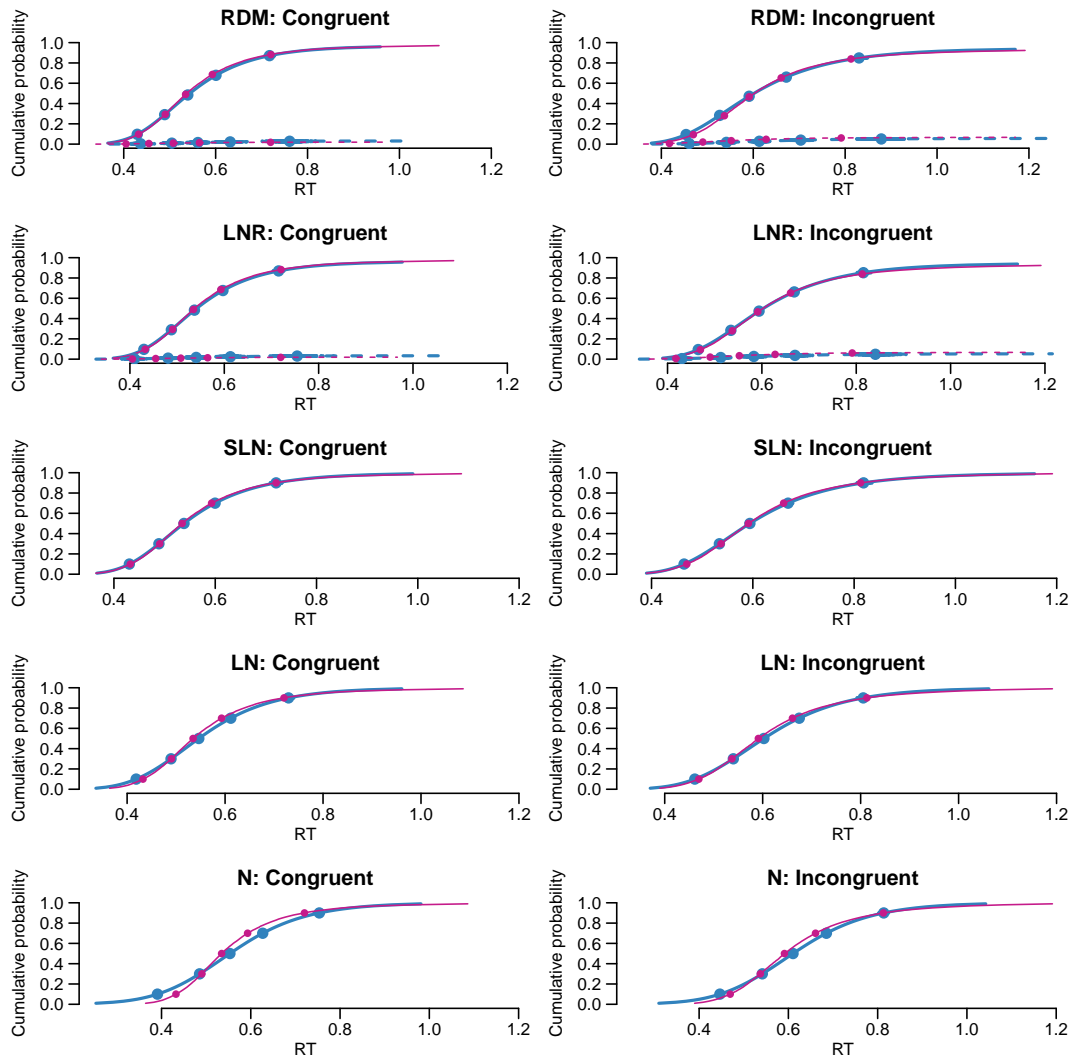


Figure 6. Observed and predicted cumulative distribution functions (CDFs) for Rey-Mermet et al. (2018). Observed = pink, blue = predicted, RDM = racing diffusion model, LNR = lognormal race model, SLN = shifted lognormal model, LN = lognormal model, N = normal model. The predictions are based on 500 sampled datasets from the posterior predictive distribution and plotted is the median including 95% credible interval. The points show the 10%, 30%, 50%, 70%, and 90% quantiles averaged across participants, separately for the two congruency conditions. Note that for the RDM and the LNR, both the (defective) CDFs of correct and incorrect responses are shown, whereas for the other models, only the CDFs of correct responses are depicted.

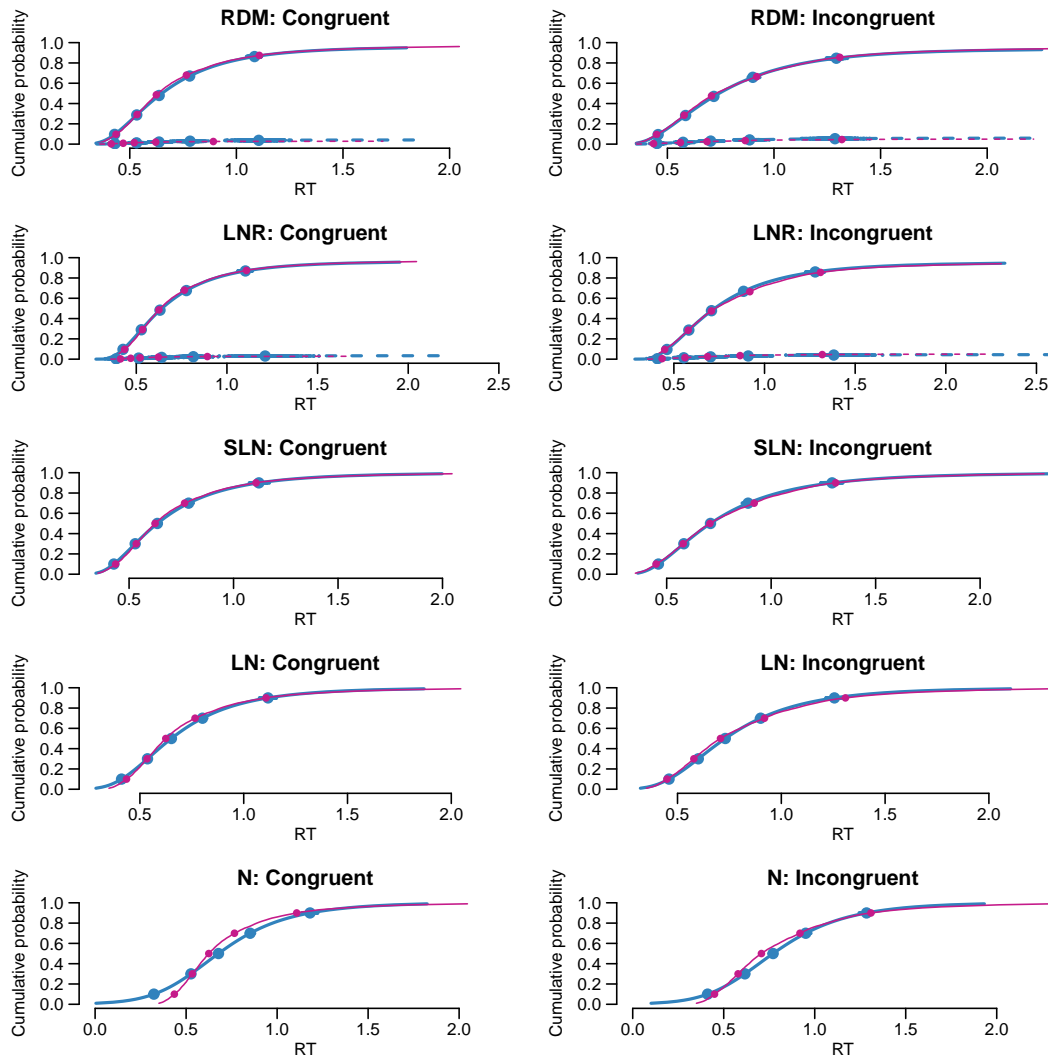


Figure 7. Observed and predicted cumulative distribution functions (CDFs) for Pratte et al. (2010). Observed = pink, blue = predicted, RDM = racing diffusion model, LNR = lognormal race model, SLN = shifted lognormal model, LN = lognormal model, N = normal model. The predictions are based on 500 sampled datasets from the posterior predictive distribution and plotted is the median including 95% credible interval. The points show the 10%, 30%, 50%, 70%, and 90% quantiles averaged across participants, separately for the two congruency conditions. Note that for the RDM and the LNR, both the (defective) CDFs of correct and incorrect responses are shown, whereas for the other models, only the CDFs of correct responses are depicted.

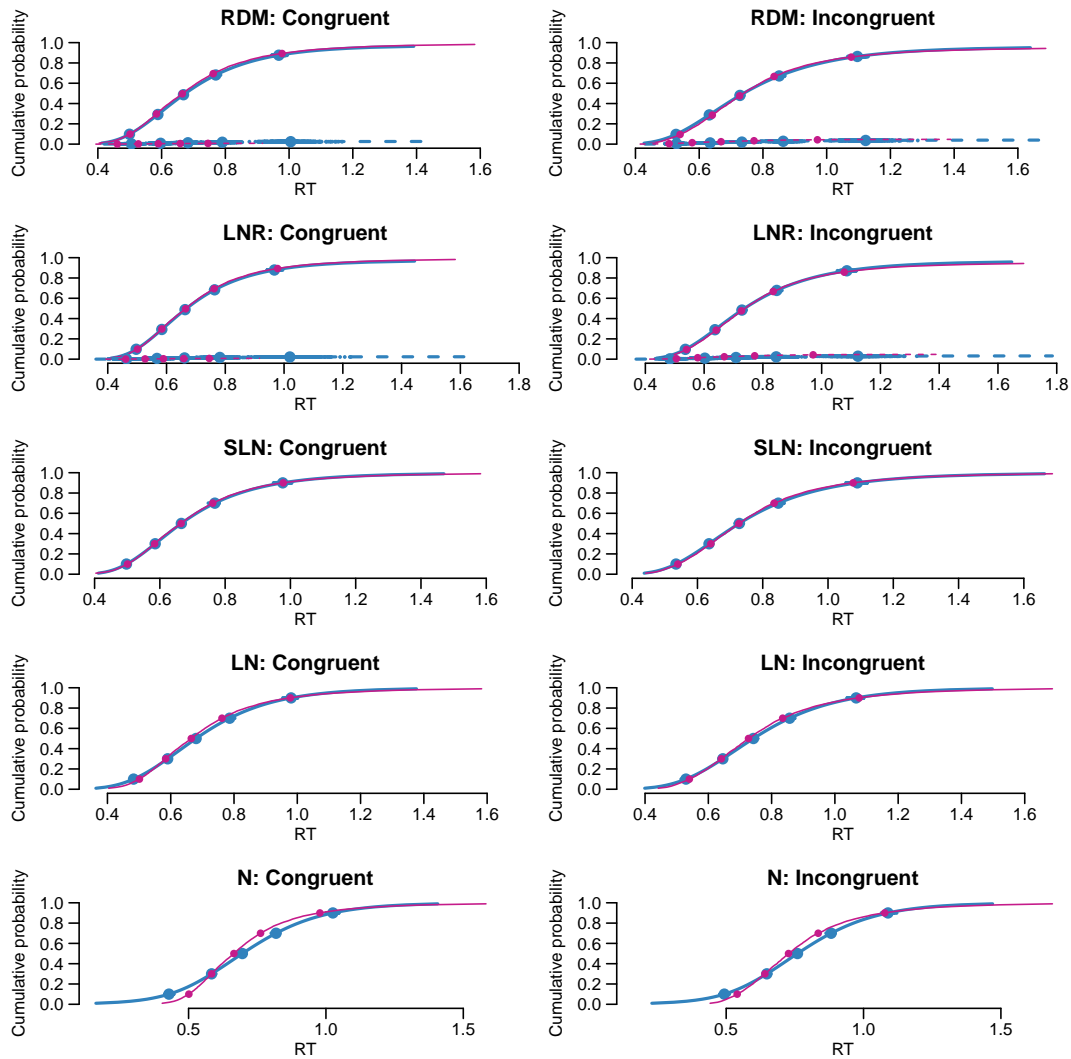


Figure 8. Observed and predicted cumulative distribution functions (CDFs) for Von Bastian et al. (2016). Observed = pink, blue = predicted, RDM = racing diffusion model, LNR = lognormal race model, SLN = shifted lognormal model, LN = lognormal model, N = normal model. The predictions are based on 500 sampled datasets from the posterior predictive distribution and plotted is the median including 95% credible interval. The points show the 10%, 30%, 50%, 70%, and 90% quantiles averaged across participants, separately for the two congruency conditions. Note that for the RDM and the LNR, both the (defective) CDFs of correct and incorrect responses are shown, whereas for the other models, only the CDFs of correct responses are depicted.

504 **Goodness of fit (RMSE)**

Table 2

Median and 95 % credible interval of the root mean squared error based on the 1%, 10%, 30%, 50%, 70%, 90%, and 99% quantiles of the correct response times and 500 samples from the posterior predictive distribution.

	Enkavi et al. (2019)	Pratte et al. (2010)	Rey-Mermet et al. (2018)	Von Bastian et al. (2016)
N	0.054 [0.051, 0.057]	0.168 [0.152, 0.184]	0.069 [0.064, 0.073]	0.126 [0.116, 0.136]
LN	0.011 [0.009, 0.013]	0.073 [0.045, 0.100]	0.054 [0.047, 0.059]	0.086 [0.071, 0.100]
SLN	0.019 [0.013, 0.024]	0.016 [0.006, 0.049]	0.027 [0.018, 0.035]	0.033 [0.013, 0.051]
LNR	0.016 [0.011, 0.021]	0.018 [0.007, 0.048]	0.031 [0.023, 0.038]	0.031 [0.023, 0.038]
RDM	0.025 [0.018, 0.031]	0.049 [0.018, 0.077]	0.027 [0.019, 0.034]	0.027 [0.019, 0.034]

Table 3

Median and 95 % credible interval of the root mean squared error computed on the proportion of correct responses using 500 samples from the posterior predictive distribution.

	Enkavi et al. (2019)	Pratte et al. (2010)	Rey-Mermet et al. (2018)	Von Bastian et al. (2016)
LNR	0.001 [0.000, 0.003]	0.002 [0.000, 0.005]	0.001 [0.000, 0.004]	0.002 [0.000, 0.005]
RDM	0.001 [0.000, 0.003]	0.011 [0.006, 0.016]	0.001 [0.000, 0.004]	0.006 [0.002, 0.010]