

Computations of confidence are modulated by mentalizing ability

E van der Plas^{1,2,3}, D Mason³, LA Livingston^{3,4}, J Craigie⁵, F Happé³ and SM Fleming^{1,2,6}

¹Wellcome Centre for Human Neuroimaging, University College London, WC1N 3BG, London UK

²Department of Experimental Psychology, University College London, WCH1H 0AP London, UK

³Social, Genetic and Developmental Psychiatry Centre, King's College London, SE5 8AF London, UK

⁴School of Psychology, Cardiff University, CF10 3AT Cardiff, UK

⁵Centre of Medical Law and Ethics, Dickson Poon School of Law, King's College London, WC2R 2LS London, UK

⁶Max Planck Centre for Computational Psychiatry and Ageing Research, University College London, WC1B 5EH London, UK

Email: elisa.plas.18@ucl.ac.uk

Code and data availability: All anonymized behavioural data and analysis scripts supporting the findings of this work are available on GitHub (github.com/metacoglab/MetaMenta-project).

Preregistration availability: <https://osf.io/u6ecx/>

Author Contributions: All authors conceptualized the project. EvdP and DM collected the data, EvdP analysed the data, EvdP and SMF wrote the paper with revisions from DM, JC, LAL, and FH. All authors approved the manuscript.

Competing Interest Statement: The authors declare no competing interest.

Keywords: Metacognition, mentalizing, autism spectrum condition, perceptual decision-making

This PDF file includes:

Main Text (8,180 words)

Figures 1 to 4

Supplementary Materials (1,295 words)

Supplementary Figures 1 to 6

ABSTRACT

1 Do people have privileged and direct access to their own minds, or do we infer our own
2 thoughts and feelings indirectly, as we would infer the mental states of others? In this study
3 we shed light on this question by examining how *mentalizing ability*—the set of processes
4 involved in understanding other people’s thoughts and feelings—relates to *metacognitive*
5 *efficiency*—the ability to reflect on one’s own performance. In a general population sample
6 ($N = 477$) we showed that mentalizing ability and self-reported socio-communicative skills
7 are positively correlated with perceptual metacognitive efficiency, even after controlling for
8 choice accuracy. By modelling the trial-by-trial formation of confidence we showed that
9 mentalizing ability predicted the association between response times and confidence,
10 suggesting those with better mentalizing ability were more sensitive to inferential cues to
11 self-performance. In a second study we showed that both mentalizing and metacognitive
12 efficiency were lower in autistic participants ($N = 40$) when compared with age, gender, IQ,
13 and education-matched non-autistic participants. Together, our results suggest that the ability
14 to understand other people’s minds predicts self-directed metacognition.

“The sorts of things that I can find out about myself are the same as the sorts of things that I can find out about other people, and the methods of finding them out are much the same.”

– G. Ryle in *The Concept of Mind* (1949)

INTRODUCTION

15 In 1949, Ryle proposed that the cognitive mechanisms employed to understand
16 ourselves are similar to those involved in understanding the feelings and experiences of other
17 people (Ryle, 1949). Since then, various proposals have echoed Ryle in suggesting that
18 *explicit metacognition*—the capacity for conscious evaluation of one’s own mental states
19 (Fleming et al., 2010; Fleming & Lau, 2014; Frith, 2012; Yeung & Summerfield, 2012) and
20 *mentalizing*—the capacity to evaluate and understand other people’s mental states (Abell et
21 al., 2000; David et al., 2008; Rosenblau et al., 2015; White et al., 2009; White et al., 2011)
22 have a common neurocomputational basis (Carruthers, 2009; Dimaggio et al., 2008; Fleming
23 & Daw, 2017; Frith, 2012; Vaccaro & Fleming, 2018).

24 According to recent perspectives on the developmental trajectory of metacognition,
25 while “core” or implicit mechanisms for self-monitoring and tracking uncertainty may be in
26 place early in infancy (Goupil & Kouider, 2016), explicit metacognition emerges around the
27 ages of 2-3 (e.g. Hembacher & Ghetti, 2014; see Goupil & Kouider, 2019 for a review), and
28 continues to be shaped in childhood and adolescence (Fandakova et al., 2017; Weil et al.,
29 2013). One potential driver of this continued development of explicit metacognition is that a
30 growing understanding of other people’s mental states may be used to refine awareness of
31 ourselves (Carruthers, 2009). For example, repeatedly perceiving a parent expressing
32 uncertainty together with their hesitation may allow a child to recognize and express
33 uncertainty when they themselves are hesitating. This hypothesis predicts that introspection is
34 not a distinct natural kind, but is instead grounded in the same processes used to understand

35 the mental states of others (Carruthers, 2009; Gazzaniga, 1995, 2000; Gopnik, 1993; Wegner,
36 2002; Wilson, 2002). This view makes several testable predictions, for example, that people
37 with a good mentalizing ability should also have good metacognitive ability; and that if
38 children have problems with inferring the mental states of others (e.g., because of a
39 neurodevelopmental condition such as autism), they may also develop difficulties with
40 understanding their own minds.

41 The second prediction can be directly studied in the context of Autism Spectrum
42 Condition (ASD)—a neurodevelopmental condition that is, in part, characterised by
43 nonverbal and verbal communicative problems, deficits in socio-emotional reciprocity
44 (American Psychiatric Association, 2013) and mentalizing difficulties (Happé, 2015;
45 Livingston & Happé, in press). If our view is correct, difficulties with understanding other
46 people’s thoughts and social communication (as is typical in autism) should also affect the
47 development of metacognition in this condition.

48 Metacognition is often quantified in laboratory tasks as the ability to provide accurate
49 confidence ratings about self-performance in a range of cognitive domains. “Good”
50 metacognitive ability is indicated by reporting lower confidence when wrong, and higher
51 confidence when right (Fleming et al., 2010; Fleming & Lau, 2014; Frith, 2012; Yeung &
52 Summerfield, 2012). This is known as metacognitive “sensitivity” and is distinct from
53 metacognitive “bias”, the tendency to be more or less confident overall (Fleming & Lau,
54 2014). Mentalizing, on the other hand, is often assessed as participants’ ability to understand
55 what agents are thinking or intending from observations of their actions and expressions
56 (Abell et al., 2000; Baron-Cohen et al., 2001; White et al., 2011). “Good” mentalizing ability
57 is indicated by correct assessment of others’ mental states. To date, six studies have examined
58 associations between metacognition and mentalizing in children or adults with autism

59 (Carpenter et al., 2019; Grainger et al., 2016; Nicholson et al., 2019; 2020; Wojcik et al.,
60 2013; Williams et al., 2018). Three of the six papers suggest, in line with the idea that
61 mentalizing and metacognition have a similar neuro-computational mechanism, that autistic
62 individuals have metacognitive difficulties that are commensurate with their mentalizing
63 capacity (Grainger et al., 2016; Nicholson et al., 2020; Williams et al., 2018). However, the
64 remaining three studies did not find deficits in metacognition in autistic compared with non-
65 autistic participants despite finding deficits in mentalizing ability (Wojcik et al., 2013;
66 Carpenter et al., 2019). Taken together, the existing data indicate a link between
67 metacognition and mentalizing, but not unequivocally so.

68 One difficulty with interpreting findings on metacognition is that its measurement is
69 often confounded by other aspects of task performance, which itself may vary across
70 individuals and clinical groups. For example, many of the studies reviewed above computed
71 people's metacognitive sensitivity as the Goodman-Kruskall gamma correlation between
72 trial-by-trial accuracy and confidence (Nelson, 1984), a measure known to be confounded by
73 *type I sensitivity* (task performance) and *metacognitive bias* (people's average confidence
74 scores) (Fleming & Lau, 2014; Maniscalco & Lau, 2012, 2014; Masson & Rotello, 2009;
75 Rahnev & Fleming, 2019; **Figure 1a**). The impact of this confound may be particularly
76 pertinent in studies comparing autistic and non-autistic people, as sensory (hyper-) sensitivity
77 (Ewbank et al., 2016; Lieder et al., 2019; Pirrone et al., 2017) and over-confidence
78 (McMahon et al., 2016; Milne et al., 2002; Zalla et al., 2015) are sometimes found to be
79 higher in autistic compared to non-autistic groups. In other words, previously reported
80 measures of *metacognitive* sensitivity may have been confounded by higher *sensory*
81 sensitivity in autistic participants.

82 A powerful approach to control for task performance confounds in studies of
83 metacognition is to use model-based metrics derived from signal detection theory, that allow
84 metacognitive sensitivity to be expressed in the same units as task performance while also
85 controlling for metacognitive bias (meta- d' ; Maniscalco & Lau, 2012, 2014). Notably, a
86 recent study identifying a positive correlation between metacognitive and mentalizing ability
87 when using this meta- d' metric to quantify metacognitive sensitivity (Nicholson et al., 2020).
88 Nicholson and colleagues (2020) measured both implicit (behavioural) and explicit (verbal)
89 metrics of choice uncertainty (defined as ‘opting-out’ from choosing or verbally reporting
90 lower confidence, respectively) and measured mentalizing ability from participants’
91 descriptions of short animations of abstract figures that vary in their level of intentionality
92 (Abell et al., 2000). The authors found that explicit, but not implicit, metacognitive sensitivity
93 was positively correlated with mentalizing ability, and significantly lower among autistic
94 children. In a second study on neurotypical adults, the authors leveraged a dual-task condition
95 in which participants completed a mentalizing or non-mentalizing-related cognitive task
96 alongside a metacognition task and found that the dual mentalizing task significantly lowered
97 metacognitive sensitivity compared to conditions in which the dual task did not require
98 mentalizing (Nicholson et al., 2020). Together these findings suggested that mentalizing and
99 metacognitive ability share a common neurocognitive basis which is commensurately
100 impaired in autistic individuals.

101 However, despite this promising result, further limitations in the measurement of both
102 mentalizing and metacognition in Nicholson et al (2020) are worth considering. First,
103 mentalizing ability was scored from participants’ written descriptions of the triangles’ mental
104 states. It has been proposed that this type of question is more prone to confounds of verbal
105 fluency than, for example, multiple-choice assessments of mentalizing (White et al., 2011).
106 This may be particularly problematic in studies of autism given that differences in verbal

107 fluency are commonly observed in this condition (Livingston, Carr, et al., 2019; Livingston et
108 al., in press; Spek et al., 2009). Second, in the metacognition task, decisions were of varying
109 choice difficulty, with some perceptual discriminations (of colour, or dot density) being
110 easier than others. When task difficulty is varying between trials and subjects, it may affect
111 measures of metacognitive ability, even when d' is controlled for (Rahnev & Fleming, 2019).
112 Finally, participants received trial-by-trial feedback on their confidence ratings, where they
113 were rewarded for reporting higher confidence on correct trials and lower confidence on error
114 trials (i.e., better metacognition was incentivized). This may have created a disadvantage for
115 autistic participants who may have difficulties with interpreting and learning from ambiguous
116 or implicit feedback (Broadbent & Stokes, 2013; Greene et al., 2019; Reed, 2019; Robic et
117 al., 2015; Sapey-Triomphe et al., 2018; Zwart et al., 2018). In other words, it could be that
118 the lower metacognitive ability in the autistic group was a consequence of failing to
119 maximize rewards on the basis of the ambiguous feedback.

120 Across two studies, we set out to control for some of the factors that might have
121 influenced the results of these previous studies by adopting experimental and computational
122 methods that are considered optimal for the assessment of metacognitive sensitivity (Rahnev
123 & Fleming, 2019; Fleming, 2017). Specifically, we measured metacognition using a
124 psychophysical task on which participants make repeated perceptual judgements and rated
125 their confidence in being correct. In order to match sensory sensitivity across participants and
126 over the course of the experiment within the same participant, we employed a staircase
127 procedure that continually adjusted sensory evidence strength on the basis of people's
128 responses. In addition, we measured the same participants' *mentalizing ability* on a separate
129 task in which they watched short animations of abstract figures that moved across the screen
130 according to distinct types of interaction (Abell et al., 2000), similar to that used by
131 Nicholson et al. (2020). Instead of providing a verbal description of each interaction,

132 participants indicated their answer using multiple choice selection (White et al., 2011;
133 Livingston et al., in press). We controlled for type 1 performance in the measurement of
134 metacognition by computing *metacognitive efficiency* ($\text{meta-}d'/d'$), which controls for type 1
135 sensitivity and metacognitive bias using the $\text{meta-}d'$ model (Maniscalco & Lau, 2012, 2014).
136 Moreover, we estimated metacognitive efficiency within a Bayesian hierarchical model that
137 allows optimal estimation of the relationship between metacognitive efficiency and individual
138 differences in mentalizing ability, while also taking into account uncertainty surrounding
139 each individual subject's parameter estimates (Fleming, 2017; Harrison et al., 2020).

140 Having confirmed a link between metacognition and mentalizing, in a second set of
141 analyses we investigated *how* the computation of confidence is modulated by mentalizing
142 ability by building hierarchical regression models of trial-by-trial confidence ratings. We
143 reasoned that, if metacognition and mentalizing rely on similar inferential processes and cues,
144 mentalizing ability should facilitate the use of behavioural cues that are similarly predictive
145 of the mental states of others. Work in cognitive psychology has often shown that people
146 have poor access to the reasons for their actions but instead infer these from contextual cues
147 (even if these cues are experimentally decoupled from the true underlying intention;
148 Gazzaniga, 1995, 2000; Nisbett & Wilson, 1977; Wegner, 2002; Wilson, 2002). For example,
149 when asked to rate their confidence in a previous decision, people's confidence reports may
150 be affected by various (behavioural) cues that are more or less related to the decision, such as
151 response times (Kiani et al., 2014; Patel et al., 2012), social context (Bang et al., 2017, 2020;
152 Van der Plas et al., 2021), as well as the quantity and reliability of evidence (Campbell-
153 Meiklejohn et al., 2010; De Martino et al., 2017; Kiani & Shadlen, 2009; Pleskac &
154 Busemeyer, 2010). Intriguingly, response times have been shown to have a causal impact on
155 the confidence levels people ascribe not only to themselves, but also to others (Palser et al.,
156 2018; Patel et al., 2012). In a series of exploratory analyses, we therefore asked whether

157 confidence was more tightly coupled to response times among participants with better
158 mentalizing ability.

159 In two independent behavioural experiments, we tested three [pre-registered](#) premises
160 of the hypothesis that metacognition and mentalizing are inter-related, namely that: (1)
161 metacognition and mentalizing ability are positively correlated, even after controlling for
162 first-order performance; (2) metacognitive efficiency is lower in people with autism, and in
163 participants with greater autistic traits; and (3) especially in those with greater difficulties
164 with social communication and understanding but not non-social autistic traits. We also
165 assessed the extent to which response times predict confidence on a trial-by-trial level by
166 conducting exploratory hierarchical regression models, asking whether the predictions of
167 confidence interacted with mentalizing ability.

METHODS

Experiment 1.

168 **Participants.** We recruited N = 501 proficient English speaking participants via Prolific
169 (<https://www.prolific.com>), a recruitment platform more representative of real populations
170 than standard student samples (Palan & Schitter, 2018). All participants accessed the
171 experiment with a desktop computer or laptop (no tablets or smartphones). Exclusion criteria
172 were responding incorrectly to a “catch” question (e.g., “If you are still paying attention,
173 please select x as your answer”); performing below or above pre-defined accuracy cut-offs
174 (60% and 90% respectively) on the metacognition task; or rating the same confidence on
175 more than 90% of the trials on the metacognition task. This resulted in the exclusion of N =
176 23 participants (5% of the total sample), leaving N = 477 participants for further analysis
177 (168 female, mean age: 28.73, SEM = 0.52 years). All participants gave informed consent
178 before the experiment, which was approved by the University College London Ethics
179 Committee (1260/003).

180 **Metacognition task.** Stimuli were programmed in JavaScript using JSPsych (version 5.0.3)
181 and hosted on the online research platform Gorilla (<https://gorilla.sc/>). Participants made 168
182 decisions across four blocks concerning which box was filled with a higher density of dots
183 (left or right, indicated by pressing the “W” or “E” key, respectively without a time limit).
184 The boxes were two black squares (each 250 x 250 pixels) which were each subdivided into
185 grids of 625 cells that were filled with 313 dots. Choice difficulty was manipulated by
186 adjusting the dot difference between boxes according to a “2-down-1-up” staircase
187 procedure: dot difference increased after every error and decreased after two consecutive
188 correct answers. Dots seemed to flicker, an effect created by replotting five different
189 configurations of the same dot difference level for 150 ms each, for a full stimulus duration
190 of 750 ms (Rollwage et al., 2018). On 26 practice trials participants received immediate

191 feedback. During the remaining trials, participants did not receive feedback but had to rate
192 their confidence that their decision was correct (on a scale from 1 “*Guessing*” to 6 “*Certainly*
193 *correct*”, without a time limit; Rouault et al., 2018).

194 ***Mentalizing task.*** We used a validated online version of the Frith-Happé Triangle Task
195 (Abell et al., 2000; Livingston et al., in press). Participants were shown twelve short (34-35
196 sec.) animations of one large red and one small blue triangle. The way in which the triangles
197 moved was manipulated across three conditions: in random animations they moved
198 purposelessly around; in Goal-Directed animations they interacted behaviourally; and in four
199 Theory of Mind (ToM) animations they interacted in a way that involves responding to the
200 other’s mental states. Participants were scored on their accuracy in classifying which
201 category the interaction pertained to (mentalizing classification) giving a score ranging
202 between 0-12 (i.e., participants could score one point after each animation). In addition, we
203 computed participants’ accuracy in categorizing the feelings of the triangles (mentalizing
204 ability; White et al., 2011). Mentalizing ability was scored as the number of correctly
205 identified mental states of each of the two triangles, after each ToM animation that had been
206 correctly identified in the mentalizing classification question. This type of mental state
207 attribution requires tracking the triangle’s intentions throughout the animation and cannot
208 simply be deduced from the general kinematics of the triangle, therefore making it less
209 susceptible to compensatory strategies. Participants had to watch the complete animation
210 before the questions appeared, after which they were allowed to decide without a time limit.
211 All animations were presented in pseudo-randomized order and after three practice
212 animations on which participants received immediate feedback.

213 ***Additional measures.*** After the two computer tasks, which were presented in
214 counterbalanced order, the following questionnaires were administered: (1) the Autism

215 Quotient-10 (AQ-10) a brief assessment of autistic traits (a higher score indicates more
216 autistic traits; Allison et al., 2012); (2) the RAADS-14, a screening tool for autistic traits in
217 adult populations which asks whether each trait was present either in childhood, currently,
218 both or neither (with a higher score indicating more autistic traits; Eriksson et al., 2013); (3)
219 the Beck Cognitive Insight Scale (BCIS) an assessment of people's ability to distinguish
220 between objective reality and subjective experience (Beck et al., 2004); and (4) the
221 International Cognitive Ability Resource (ICAR) a brief assessment of fluid intelligence
222 (Condon & Revelle, 2016). More details on these questionnaires are provided in
223 **Supplementary Materials.**

224 **Statistics.** The hypotheses and analyses for this study were pre-registered
225 (<https://osf.io/vgy7a/>). Validation checks are reported in the **Supplementary Material** and
226 consisted of Spearman's rho correlations (which are recommended for ordinal data) to assess
227 relationships between main composite survey scores. Equal variances were assumed if not
228 otherwise specified. We report P values at a 0.05 alpha level and the 95% confidence interval
229 (95% CI) of the test statistic. Type-1 cognitive and type-2 metacognitive parameters were
230 estimated using the open source HMeta-d toolbox (<https://github.com/metacoglab/Hmeta-d>)
231 implemented in MATLAB (version 9.7.0). Type-2 $meta - d'$, the ability to determine one's
232 accuracy with confidence ratings, was inferred using Markov chain Monte Carlo (MCMC)
233 Bayesian sampling procedures using JAGS (<http://mcmc-jags.sourceforge.net>) across 30,000
234 samples after a burn-in of 1,000 samples distributed across three chains. Our parameter of
235 interest was M_{ratio} ($meta-d'/d'$), or metacognitive efficiency, which expresses metacognitive
236 sensitivity ($meta-d'$) relative to task performance (d' ; in other words, an M_{ratio} of 1 implies
237 participants have optimal metacognitive efficiency; Fleming, 2017).

238 We assessed model convergence for each HMeta-d model by ensuring that the
239 consistency of the posteriors within and between chains, the Gelman-Rubin (G-R) \hat{R} statistic,

240 was below 1.1 (Gelman & Rubin, 1992) and by visually inspecting the chains
 241 **(Supplementary Materials)**. In addition, each reported model was checked for reliability by
 242 conducting posterior predictive checks which are summarized in the **Supplementary**
 243 **Materials**.

244 To test the first pre-registered hypothesis of a positive association between
 245 metacognitive efficiency and mentalizing ability, we incorporated a simultaneous hierarchical
 246 estimation of the beta coefficient (β) of the impact of our standardized mentalizing ability
 247 score, *meta*, on the log of metacognitive efficiency, *log(Mratio)*:

$$248 \quad \log(Mratio)_s \sim \log(Mratio)_0 + \beta \textit{meta}_s + \varepsilon_s \quad (1.1)$$

249 *log(Mratio)*₀ denotes baseline group-level metacognitive efficiency; *meta*_s is the
 250 mentalizing score for subject *s*; and ε_s refers to noise that is drawn from a T-distribution with
 251 variance σ_δ and 5 degrees of freedom, multiplied by a noise parameter ζ . We used priors
 252 found to provide the most efficient regression parameter recovery (Harrison et al., 2020),
 253 which were drawn from Gaussians $N(\mu, \sigma)$, half-Gaussians $HN(\mu, \sigma)$ and T-distributions
 254 $T(\mu, \sigma, df)$:

$$\textit{meta}_0 \sim N(0,1)$$

$$\beta \sim N(0,1)$$

$$\sigma_\delta \sim HN(1)$$

$$\zeta \sim Beta(1,1)$$

$$\delta_s = T(0, \sigma_\delta, 5)$$

$$\varepsilon_s = \zeta * \sigma_\delta$$

255 The highest density interval (HDI) represents the ‘credible’ posterior range within
 256 which 95% of the estimated regression coefficient falls. We plotted the HDI for the
 257 regression coefficient and assessed significance by computing the probability that it differed

258 from zero: $P_{\theta}(HDI < 0 | HDI > 0)$, where a higher probability suggests a stronger effect
 259 (Kruschke, 2010).

260 We also calculated $\log(Mratio)_s$ at the individual level for use in post-fit frequentist
 261 analyses. We used a linear model with $\log(Mratio)_s$ as the dependent variable and $menta_s$
 262 and covariates (standardized age, IQ, gender [-1: female, 1: male] and education (edu) [1: no
 263 education, 2: high school or equivalent, 3: some college, 4: BSc, 5: MSc, 6: doctoral]) as
 264 independent variables:

$$\log(Mratio)_s \sim \log(Mratio)_0 + \beta_1 menta_s + \beta_2 age_s + \beta_3 IQ_s + \beta_4 gender_s + \beta_5 edu_s + \varepsilon_s \quad (1.2)$$

265 To test the effect of autistic traits on $\log(Mratio)_s$ we ran the same models specified in
 266 Equations 1.1 and 1.2 but now replacing $menta_s$ with the RAADS-14 main composite
 267 autistic trait scores (Eriksson et al., 2013). In preliminary analyses we failed to replicate
 268 previous findings of a negative correlation between mentalizing ability and AQ-10 scores
 269 (Allison et al., 2012), and therefore (deviating from our pre-registration plan) we decided to
 270 conduct all further analysis of questionnaire data using RAADS-14 scores alone (Bertrams,
 271 2021).

272 To assess the effects of trial-by-trial standardized (log) response times $\log RT$ and
 273 accuracy on confidence, we conducted hierarchical mixed-effect regression models using the
 274 ‘lme4’ package in R (version 3.3.3) and plotted the standardized fixed-effect beta coefficients
 275 of the model fits. We obtained the P -values of the regression coefficients using the *car*
 276 package. All models include a random effect at the participant level and all statistics are
 277 computed at the group level. We report type III Wald chi-square tests (χ^2), degrees of

278 freedom (*df*) for fixed effects, and estimated beta-coefficients (β) together with their standard
 279 errors of the mean (\pm SEM) and *P*-values of the associated contrasts.

280 To investigate if *logRT* informs confidence differently as a function of individual
 281 differences in autistic traits, we tested whether a hierarchical mixed-effect regression model
 282 better predicts trial-by-trial confidence (*conf*) when the predictor variables accuracy (*acc*) [-1:
 283 error, 1: correct], z-score of the log response time (*RT*) and their interactions (Equation 2.1)
 284 were allowed to vary as a function of individual differences in standardized autistic trait
 285 scores (ASD; Equation 2.2.):

$$conf \sim acc + logRT + acc:logRT + (1 + acc + logRT + acc:logRT|subj) \quad (2.1)$$

$$conf \sim ASD: (acc + logRT + acc:logRT) + (1 + acc + logRT + acc:logRT|subj) \quad (2.2)$$

286 The results of the Likelihood Ratio Test are expressed in terms of the *Akaike Information*
 287 *Criterion (AIC)*: $\Delta AIC = AIC_{\text{Equation 2.1}} - AIC_{\text{Equation 2.2}}$, and the *Log Likelihood (LL)*: $\Delta LL =$
 288 $LL_{\text{Equation 2.1}} - LL_{\text{Equation 2.2}}$ with associated *P* values extracted from a type III Wald chi-square
 289 tests (χ^2).

Experiment 2

290 **Participants.** We recruited a sample of *N* = 43 autistic participants via the research charity
 291 Autistica (www.autistica.org.uk). Interested participants first completed an online pre-
 292 screening questionnaire that included questions about mental health and demographics.
 293 Participants that met the inclusion criteria (i.e., aged between 18 and 50 years old and a self-
 294 reported diagnosis of autism spectrum condition by a health professional) were sent a link to
 295 the online experiment that could be accessed with a desktop computer or laptop (no tablets or
 296 smartphones). Exclusion criteria were the same as in Experiment 1. Three participants were

297 excluded: one participant performed below the *a priori* accuracy cut-off and two participants
298 performed above the *a priori* accuracy cut-off. This resulted in the exclusion of N = 3
299 participants (7.5% of the total sample, which is consistent with Experiment 1), leaving data
300 from N=40 participants for analysis (37 female, mean age: 37.90, SEM = 1.59 years). All
301 participants gave informed consent before experiment onset which was approved by the
302 Research Ethics Office at King's College London (HR-19/20-17704).

303 To obtain an equal number of comparison participants we re-analysed a subset of the
304 dataset from Experiment 1 which used the same experimental paradigms and questionnaire
305 battery. The dataset from Experiment 1 consisted of N = 477 English speaking participants
306 from the general population (198 female, mean age: 28.73, SEM = 0.52). Data on mental
307 health conditions was not collected. To ensure that the participants from this dataset provided
308 a comparison group with low autistic traits, we first reduced the number to N = 97
309 participants scoring in the lowest 50% quantile of RAADS-14 and AQ-10 responses (a score
310 lower than 16 and 5, respectively, which is more stringent than the clinical cut-off score;
311 Ashwood et al., 2016; Eriksson et al., 2013). Next, to ensure the groups were well-matched
312 on other characteristics, for each included autistic participant we manually selected a
313 comparison participant of similar gender (a high proportion of females in the autism group
314 meant that it was not possible to find a 1:1 gender match for three participants); who was
315 within ± 5 years from the target age; ± 2 levels from the target education; and ± 5 ICAR points
316 from the target fluid intelligence level. These criteria were identified after initial exploration
317 indicated they provided sufficient flexibility to provide a reasonable match between the two
318 groups on all relevant dimensions. Importantly, participant selection was carried out prior to
319 hypothesis testing.

320 ***Experimental paradigm.*** The experimental procedure was the same as in Experiment 1.

321 **Statistics.** Statistical inference was conducted similarly to analysis of Experiment 1.
 322 Validation checks are reported in the **Supplementary Material**. To investigate if
 323 metacognitive efficiency was different between the autism and comparison group, we fitted a
 324 linear model with *Mratio* from a single-subject fit as dependent variable, clinical group
 325 [autism: -0.5, comparison: 0.5] and covariates (standardized age, IQ, gender [-1: female, 1:
 326 male] and education (edu) [1: no education, 2: high school or equivalent, 3: some college, 4:
 327 BSc, 5: MSc, 6: doctoral]) as independent variables:

$$\begin{aligned} \log(Mratio)_s \sim & \beta_0 + \beta_1 \text{group}_s + \beta_2 \text{age}_s + \beta_3 \text{IQ}_s + \beta_4 \text{gender}_s \\ & + \beta_5 \text{edu}_s + \varepsilon_s \end{aligned} \quad (3.1)$$

328 We also conducted hierarchical regressions using the HMeta-d toolbox in which *Mratio* in
 329 the autism and comparison groups were estimated in separate models that controlled for the
 330 following covariates:

$$meta_s \sim \beta_0 + \beta_1 \text{age}_s + \beta_2 \text{IQ}_s + \beta_3 \text{gender}_s + \beta_4 \text{edu}_s + \varepsilon_s \quad (3.2)$$

To assess significance, we computed the probability P_θ of overlap between the HDI posterior distribution of *Mratio* in the autism and comparison group:

$$P_\theta(HDI_{autism} < HDI_{comparison})$$

331 To assess whether the effect of *logRT* on confidence was different for autistic and
 332 comparison participants, we conducted hierarchical mixed-effect regression models using the
 333 “lme4” package in R (version 3.3.3), similar to the method used in Experiment 1, but now
 334 using a dummy variable denoting clinical group (group [autism: -0.5, comparison: 0.5])
 335 instead of continuous autistic trait scores. To visualize the direction of significant effects we

336 obtained the beta-coefficients of $\log RT$ on confidence for each clinical group and on error
337 and correct trials, separately:

$$conf_{acc/group} \sim \beta_0 + \beta_1 \log RT_s + \beta_2 \text{gender}_s + \beta_3 \text{edu}_s + \varepsilon_s \quad (3.3)$$

RESULTS

338 **Experiment 1**

339 The staircase converged to a stable performance level within and between participants
 340 (choice accuracy: $M = 75\%$, $SEM = 0.23$). Given that staircase variability can affect
 341 estimates of metacognitive sensitivity (Rahnev & Fleming, 2019), we also computed each
 342 individual's experienced stimulus variability (the ratio between the standard deviation of
 343 stimulus difficulty and average stimulus difficulty) and established that stimulus variability
 344 was not correlated with metacognitive efficiency ($r_{S475} = -0.068$, $P = 0.137$; **Supplementary**
 345 **Figure 1.2b**).

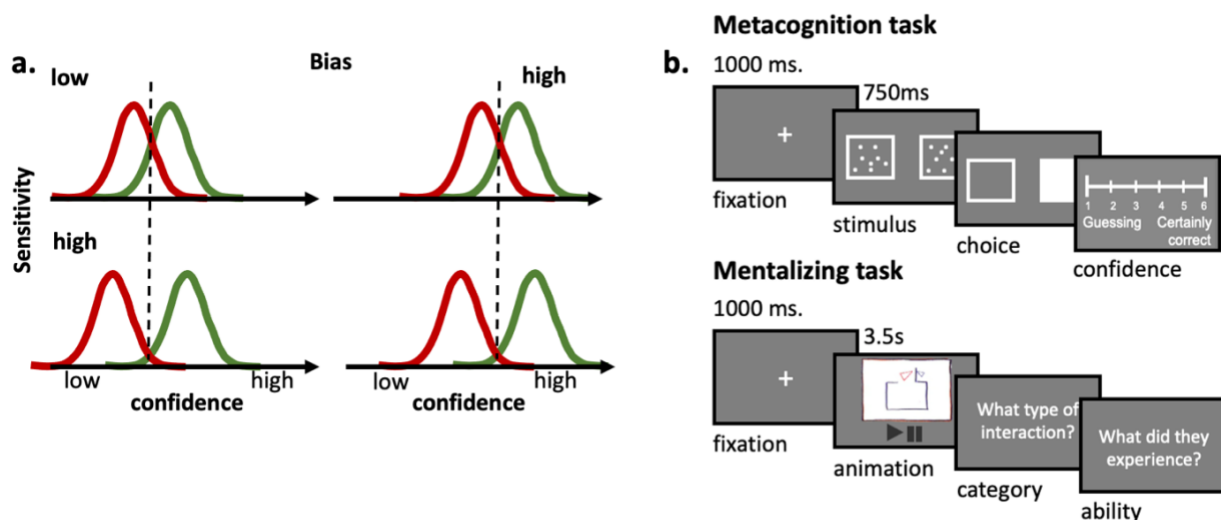


Figure 1. Task design and dissociation between metacognitive sensitivity and bias. *a.* Hypothetical Gaussian distributions of confidence for correct (green) and incorrect (red) decisions. The left panel represents a decider with low confidence; the right panel represents a decider with high confidence. Metacognitive sensitivity is defined as the separation in confidence between correct and incorrect decisions; metacognitive bias is the overall confidence expressed. *b.* On the metacognition task, participants made judgments about which patch with dots had a higher density (left or right). After this, they were asked to rate their confidence on a scale from 1 “Guessing” to 6 “Certainly correct”. On the mentalizing task, participants watched animations of moving triangles and were asked to categorize and interpret the interaction of the triangles.

346 We next investigated the hypothesis of a positive association between metacognitive
 347 efficiency and mentalizing ability within the hierarchical meta- d' model. When we examined
 348 the beta coefficient representing the impact of mentalizing ability on metacognitive

349 efficiency, the HDI was positive and did not encompass zero (hierarchical estimation: 95%
350 HDI [0.01, 0.09]), with 99% of the sampled beta values being higher than zero
351 ($P_{\theta}(\text{HDI mentalizing ability} > 0) = 0.99$; **Figure 2a**) indicating a significant positive relationship.
352 To confirm this effect while controlling for covariates of age, gender, IQ and education, we
353 used a linear regression model with the standardized log metacognitive efficiency from a
354 single-subject model as a dependent variable and mentalizing ability and these covariates as
355 predictor variables. This approach again revealed a positive relationship between mentalizing
356 and metacognition (linear regression model: $\beta_{\text{mentalizing efficiency}} = 0.11$, $SE = 0.05$, $t_{476} =$
357 2.26 , $P = 0.02$) and no effects of the covariates ($P > 0.05$), suggesting that participants who
358 were better at inferring the mental states and interactions on the mentalizing task were also
359 better at tracking their performance on the metacognition task.

360 To investigate how mentalizing was related to metacognition, we next tested the
361 hypothesis that mentalizing is associated with a greater impact of response times on
362 confidence. Specifically, we estimated a hierarchical mixed-effects model predicting trial-by-
363 trial explicit confidence levels on the metacognition task from differences in standardized log
364 response times (logRT) and accuracy [error: -0.5, correct: 0.5] (**Equation 2.1**), and asked
365 whether this model provided a better fit when these predictors were allowed to vary as a
366 function of the participants' mentalizing ability (**Equation 2.2**). A Likelihood Ratio Test
367 indicated that this was the case ($\chi^2(4) = 27.59$, $P = 1.51e-05$) which was also confirmed by
368 several goodness-of-fit indices (log likelihood (LL): $\Delta LL = 13$, Akaike Information Criterion
369 (AIC): $\Delta AIC = -20$, Bayesian Information Criterion (BIC): $\Delta BIC = 17$ and $\Delta \text{Deviance} = -28$),
370 suggesting a significant relationship between mentalizing and the computations underpinning
371 confidence formation.

372 We next asked how mentalizing modulated the construction of confidence by
 373 investigating which predictor variables interacted with mentalizing ability. We found that
 374 participants with better mentalizing ability reported lower overall confidence in their own
 375 responses than participants with lower mentalizing ability (hierarchical linear regression,
 376 main effect of mentalizing ability: $\chi^2(1) = 6.08$, $P = 0.01$, $\beta = -0.04$, $SE = 0.02$). In addition,
 377 participants with higher mentalizing ability scores modulated their confidence ratings more
 378 on the basis of their response times than participants with lower scores of mentalizing ability
 379 (interaction effect of logRT x mentalizing ability: $\chi^2(1) = 21.92$, $P = 2.84e-06$, $\beta = -0.03$, SE
 380 $= 0.006$; **Figure 2b**), consistent with the idea that mentalizing facilitates metacognition by
 381 facilitating self-inference on the basis of externally visible behavioural cues.

382

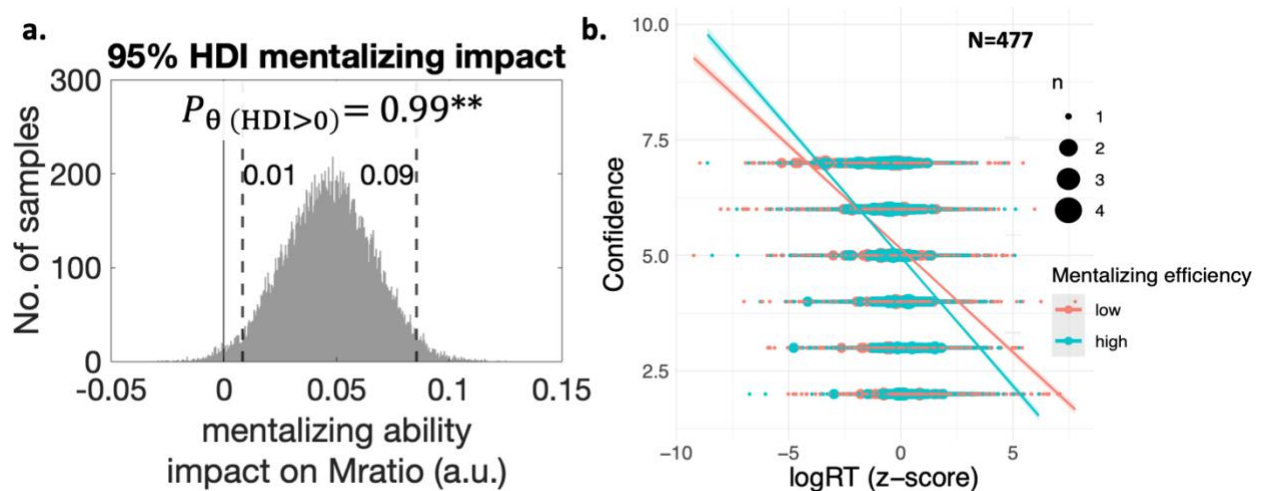


Figure 2. Mentalizing modulates computation of confidence. *a.* Posterior distribution over the regression coefficient relating mentalizing ability to metacognitive ability. The dashed lines represent the 95% highest density interval (HDI), P_{θ} indicates the probability that the posterior samples are greater than zero, $** P < 0.01$ in the frequentist linear model. *b.* Confidence was negatively related to response times (logRT). Trial-by-trial response times have a higher impact on the estimated confidence of participants scoring above the median of mentalizing ability scores (in turquoise) than participants scoring below the median (in pink). Shaded area represents the Standard Deviation from the Mean (\pm SDM).

383 Next, we addressed the second hypothesis of a negative association between
 384 metacognitive efficiency and autistic traits in the general population, as assessed with the

385 AQ-10 (Allison et al., 2012) and the RAADS-14 questionnaires (Eriksson et al., 2013). First,
386 we evaluated whether participants with higher scores of autistic traits had lower mentalizing
387 ability, by conducting a linear regression model with mentalizing ability as the dependent
388 variable and autistic trait scores and the covariates (age, gender, education, IQ) as predictor
389 variables. We found the expected negative relationship between mentalizing ability and
390 RAADS-14 scores (linear regression model: $\beta_{RAADS-14} = -0.002$, $SE = 0.0009$, $t_{476} = -2.21$, P
391 $= 0.03$) but not AQ-10 scores (linear regression model: $\beta_{AQ10} = 0.006$, $SE = 0.004$, $t_{476} =$
392 1.33 , $P = 0.19$). This unexpected finding, together with recent re-evaluations of the reliability
393 of the AQ-10 scale (Bertrams, 2021), and the greater developmental information captured by
394 the RAADS-14, led us to focus on RAADS-14 scores in the remainder of the analyses.

395 Next, we asked whether compromised mentalizing ability in participants with higher
396 scores of autistic traits was associated with lower metacognitive efficiency. To test this, we
397 estimated the correlation between metacognitive efficiency and RAADS-14 scores within a
398 hierarchical regression model. The 95% HDI for the coefficient of RAADS-14 scores was
399 negative on average, ranging from $[-0.057, 0.019]$, but encompassed zero (hierarchical
400 estimation: $P_{\theta (HDI_{RAADS} < 0)} = 0.82$). A frequentist linear model that controlled for the
401 covariates also confirmed that participants with higher scores of autistic traits do not
402 necessarily also have compromised metacognitive efficiency (linear regression model:
403 $\beta_{RAADS14} = -0.05$, $SE = 0.05$, $t_{476} = -1.09$, $P = 0.28$).

404 An alternative explanation hypothesis is that autistic traits as measured by the
405 RAADS-14 do not have a direct impact on the metacognitive efficiency score, but rather
406 affect the construction of confidence. To examine this, we tested if our mixed-effect
407 hierarchical regression model better predicts trial-by-trial confidence levels on the
408 metacognition task when the predictors (accuracy, logRT and their interactions) were allowed

409 to vary as a function of differences in autistic traits. A Likelihood Ratio Test indeed suggests
410 that an interaction term on autistic traits improved the fit of the model ($\chi^2(4) = 14.52, P =$
411 0.006) which was further confirmed by several goodness-of-fit metrics ($\Delta LL: 7, \Delta BIC: -31,$
412 $\Delta AIC: 7$ and $\Delta Deviance: -15$), indicating that the computation of confidence differs as a
413 function of individual differences in autistic traits.

414 We next asked in what way people with higher scores for autistic traits constructed
415 their confidence differently, by testing which predictor variables interacted with RAADS-14
416 scores. We found that participants with higher scores for autistic traits reported lower
417 confidence overall (hierarchical linear regression, main effect of RAADS-14: $\chi^2(1) = 4.86, P =$
418 $0.027, \beta = -0.008, SE = 0.004$). In addition, explicit confidence was more informed by
419 logRT among participants with lower scores for autistic traits than among participants with
420 higher scores for autistic traits (interaction effect of logRT x RAADS-14: $\chi^2(1) = 6.46, P =$
421 $0.011, \beta = 0.004, SE = 0.001$). In **Figure 3a** we plot the extracted beta coefficients of the
422 impact of response times on confidence for participants scoring above and below the median
423 cut-off on autistic traits on error and correct trials separately, which shows that this effect was
424 driven by participants with higher autistic trait scores having a lower impact of response
425 times on error-trials than participants with lower autistic traits (three-way interaction of
426 logRT x RAADS-14 x accuracy: $\chi^2(1) = 4.63, P = 0.031, \beta = -0.003, SE = 0.001$). Together
427 these results suggest that participants with higher autistic traits use response times less to
428 infer they have committed an error than participants with lower autistic trait scores.

429 These results suggest that compromised mentalizing ability may specifically affect the
430 relationship between response times and confidence. We next asked whether specifically
431 social aspects of the autistic phenotype, rather than non-social aspects, negatively impact
432 metacognition. In an exploratory analysis we estimated the correlation between

433 metacognitive efficiency and self-reported social skills with hierarchical regression models.
 434 This analysis revealed that participants with self-reported difficulties in everyday types of
 435 social interaction, measured by the ‘mentalizing’ sub-scale of the RAADS-14, had lower
 436 metacognitive efficiency than participants with better self-reported social skills (hierarchical
 437 estimation: HDI: [-0.07, 0.00], $P_{\theta} (HDI \text{ social skills} < 0) = 0.97$; frequentist linear regression: $\beta =$
 438 -0.09, $SE = 0.05$, $t_{476} = -1.84$, $P = 0.067$; **Figure 3b**). In contrast, the non-social sub-scale of
 439 the RAADS-14 was not associated with metacognitive efficiency (hierarchical estimation:
 440 HDI: [-0.04, 0.04], $P_{\theta} (HDI \text{ nonsocial skills} < 0) = 0.43$; frequentist linear regression: $\beta = -0.007$,
 441 $SE = 0.05$, $t_{476} = -1.14$, $P = 0.89$; **Figure 3c**). Together, these results suggest that self-
 442 reported social, but not non-social, autistic traits are negatively associated with metacognitive
 443 efficiency.

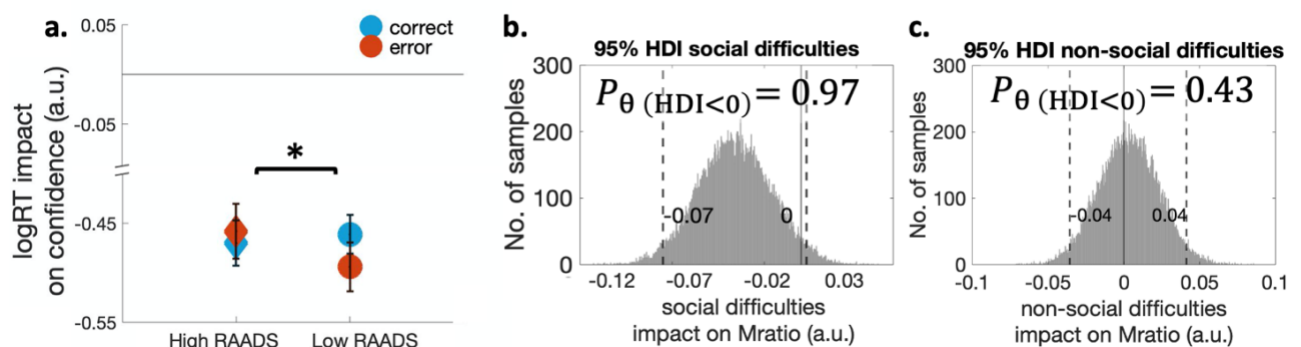


Figure 3. Autistic trait differences modulate metacognitive efficiency. *a.* Standardized beta coefficients of the impact of logRT on confidence from a hierarchical mixed-effect regression model on error trials (red) and correct trials (blue) for participants with high and low RAADS scores (above and below the median cut-off, respectively). *b.* Posterior estimates of the hierarchically estimated beta coefficient relating the social subscale of RAADS-14 to metacognitive efficiency. *c.* Posterior estimates of the hierarchically estimated beta coefficient relating the non-social subscale of RAADS-14 to metacognitive efficiency. The dashed lines represent the 95% highest density intervals (HDI), P_{θ} indicates the probability that the posterior samples are different from zero. Error bars represent group means \pm SEM, * $P < 0.05$ of the interaction effect between RAADS and logRT on confidence.

445 In summary, in Experiment 1 we found a metacognitive benefit for participants with
 446 better mentalizing ability. We further disentangled the mechanism of this effect by showing
 447 that mentalizing ability is associated with a tighter coupling between response times and

448 confidence in errors. Metacognition was less efficient in participants with higher scores for
449 autistic traits, in particular, among participants who report greater difficulties with self-
450 reported social difficulties. Together these results provide initial evidence that metacognitive
451 processes are related to mentalizing capacity.

452

453 **Experiment 2**

454 In Experiment 1 we found that metacognitive and mentalizing abilities are related, potentially
455 by affecting the extent to which response times modulate confidence. Against our
456 expectation, we did not find a statistically significant negative correlation between autistic
457 traits and metacognitive efficiency. One explanation of this null result is that the variation in
458 autistic traits was not pronounced enough in our general population sample to allow
459 estimation of this relationship. In Experiment 2 we sought to compare data from $N = 40$
460 autistic participants recruited via the charity organization Autistica to a matched comparison
461 group of $N = 40$ participants subsampled from the dataset of Experiment 1. As a result of the
462 selection procedure described in **Methods**, both groups had similar age (independent samples
463 t-test, $t_{78} = 0.90$, $P = 0.37$), gender (independent samples t-test, $t_{78} = 1.07$, $P = 0.29$),
464 education ($M_{\text{autism}} = 4.00$, $SE = 0.06$; $M_{\text{comparison}} = 3.92$, $SE = 0.19$; independent samples t-test,
465 $t_{78} = 0.25$, $P = 0.80$) and IQ scores ($M_{\text{autism}} = 9$, $SE = 0.54$; $M_{\text{comparison}} = 7.90$, $SE = 0.52$;
466 independent samples t-test, $t_{76} = 1.45$, $P = 0.15$). In addition, as a result of the calibration
467 procedure, first-order performance on the metacognition task was not statistically different
468 between groups ($M_{\text{autism}} = 0.75$, $SE = 0.01$; $M_{\text{comparison}} = 0.74$, $SE = 0.008$; independent
469 samples t-test, $t_{78} = 0.52$, $P = 0.60$; see **Supplementary Material** for other reliability
470 checks).

471 Having shown that the two groups were matched in terms of demographics and
472 general cognitive ability, we next asked if autistic participants had lower mentalizing ability

473 than comparison participants by testing a linear regression model with mentalizing ability as
474 independent variable and clinical group [autism: -0.5, comparison: 0.5] and the covariates
475 (age, gender, IQ, and education) as predictor variables. When we do this, we find that
476 mentalizing ability was indeed lower for autistic participants than comparison participants,
477 but not significantly so (linear regression: $\beta_{group} = -0.43 (0.25)$, $t_{68} = -1.72$, $P = 0.089$).

478 Next, we use a similar linear regression model to test if the autism group had lower
479 metacognitive efficiency than the comparison group. In line with our pre-registered
480 hypotheses, this indeed revealed significantly lower metacognitive efficiency in autistic
481 participants than in comparison participants (linear regression model: $\beta_{group} = -0.60 (0.25)$,
482 $t_{63} = -2.46$, $P = 0.016$; **Figure 4a**) with no effects of the covariates. We next estimated
483 metacognitive efficiency within a hierarchical model fitted to each group separately, while
484 accounting for the effects of IQ, age, gender and education. The HDI of metacognitive
485 efficiency in the autism group (HDI [0.92, 0.55]) was quantitatively lower than that of the
486 comparison group (HDI [0.84, 0.52]) in 78% of the samples $P_{\theta} (HDI_{ASD} < HDI_{comparison}) =$
487 **0.78 (Figure 4b)**, although did not reach significance at the classical 95% threshold. Taken
488 together these analyses provide some evidence in support of our pre-registered hypothesis of
489 lower metacognitive efficiency in autism.

490 Finally, building upon a hierarchical mixed-effect regression model of trial-by-trial
491 predictions of confidence on the metacognition task, we next tested whether the model could
492 better predict confidence levels when the predictors (**Equation 2.1**), were allowed to vary as
493 a function of whether the subject was autistic or not (**Equation 2.2**). A likelihood ratio test
494 indicated that this was the case ($\chi^2(4) = 966.46$, $P < 2.20e^{-16}$) which was further strongly
495 confirmed by goodness-of-fit indices (ΔLL : -484, ΔAIC : 958, ΔBIC : 929 and $\Delta Deviance$:

966), supporting the prediction that confidence formation in autistic participants is
 967 qualitatively distinct to comparison participants.

968 Consistent with the results of Experiment 1 we found that autistic participants report
 969 lower confidence than comparison participants in general (hierarchical regression model,
 970 main effect of group: $\chi^2(1) = 768.50$, $P < 2.0e^{-16}$, $\beta = 0.82$, $SE = 0.03$). Autistic
 971 participants show a marginally lower impact of response times in error trials than comparison
 972 participants (three-way interaction logRT x group x accuracy: $\chi^2(1) = 3.086$, $P = 0.060$,
 973 $\beta = 0.10$, $SE = 0.06$). In **Figure 4c** we plot the impact of response times on confidence on
 974 error and correct trials separately, which shows that the negative impact of RT on confidence
 975 was less negative in autistic participants than in comparison participants, suggesting a weaker
 976 influence on response times on confidence in error trials.

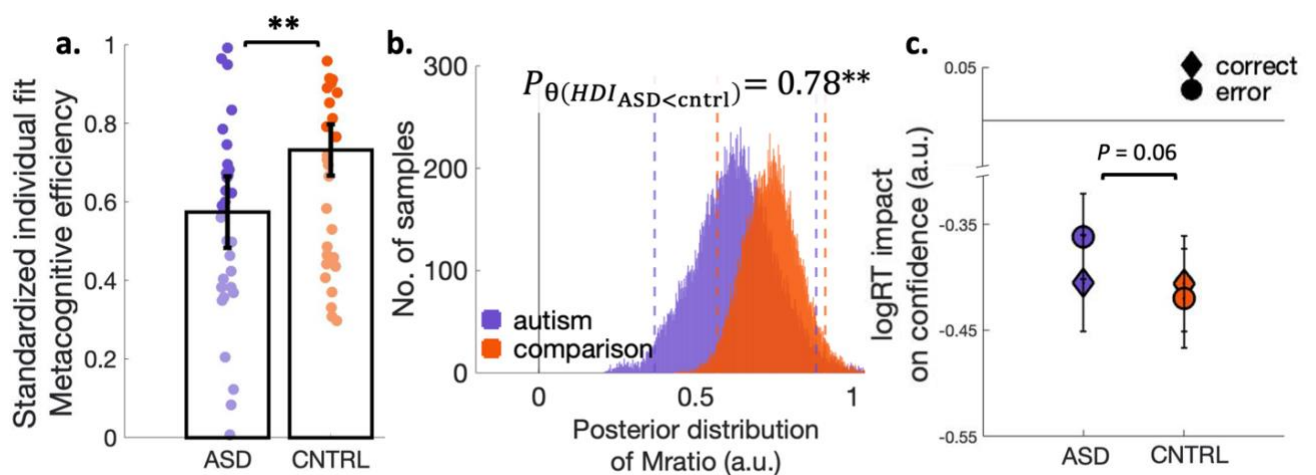


Figure 4. Differences in metacognitive efficiency and confidence formation in autism. a. Metacognitive efficiency estimated from a single-subject Bayesian model fit is significantly lower in the autism group ($N=40$) than in the comparison group ($N=40$). Error bars represent group mean \pm SEM. **b.** Posterior estimates of metacognitive efficiency from independent group model fits (autism in purple, controls in orange) where the dashed lines represent the highest density intervals (HDI) and P_{θ} represents the probability that the HDI of the autism group is lower than the HDI of the comparison group. **c.** Impact of logRT on confidence on error and correct trials for autism and comparison participants. Error bars represent group means \pm SEM.

507 In summary, in Experiment 2 we show that metacognitive efficiency is compromised
508 in autism and reveal a weaker association between response times and confidence in autistic
509 participants in contrast to matched comparison participants.

DISCUSSION

510 Across two behavioural experiments we show that mentalizing ability is positively
511 related to metacognition. In a general population sample of $N = 477$ participants we found
512 that individuals who were better at self-reported social skills and mentalizing could also more
513 reliably track their own accuracy on a perceptual discrimination task. By investigating the
514 trial-by-trial computations of confidence, we were able to investigate precisely how
515 mentalizing relates to metacognition. Notably, mentalizing ability was associated with a
516 tighter coupling between response times and confidence, suggesting that mentalizing ability
517 may facilitate inference on cues to self-performance. In a second dataset with autistic
518 participants, we show that the mentalizing difficulties that characterize this condition are
519 associated both with compromised metacognitive ability and replicate the findings of
520 Experiment 1 that autistic traits are associated with a weaker link between response times and
521 confidence. Together, these findings suggest that processes involved in inferring other
522 people's mental states may also facilitate self-directed metacognition, and vice versa.

523 We quantified metacognition as the ability to reliably separate correct from incorrect
524 decisions with confidence ratings (Flavell, 1979; Fleming et al., 2010; Rollwage et al., 2018;
525 Rouault et al., 2018). Several studies have suggested confidence is 'read out' from how much
526 reliable evidence has been seen, either during the course of the decision itself (Kiani &
527 Shadlen, 2009; Pleskac & Busemeyer, 2010) or after an initial decision has been made (post-
528 decisional evidence processing; Fleming et al., 2018; Resulaj et al., 2009; Talluri et al., 2018;
529 van den Berg et al., 2016). Other studies suggest that response times also provide a
530 behavioural cue to confidence (Kiani et al., 2014; Patel et al., 2012). How, then, might
531 mentalizing play a role in confidence construction? Recent theoretical models suggest that
532 confidence estimates reflect an inference about the state of the decider, informed by

533 behavioural and cognitive cues—suggesting a computational parallel between self- and other-
534 evaluation (Fleming & Daw, 2017). Indeed, evidence strength (Campbell-Meiklejohn et al.,
535 2017) and response times (Patel et al., 2012) appear to be used similarly to infer both one’s
536 own and others’ confidence. However, isolating such metacognitive capacity requires tight
537 control over the evidence going into a decision, to avoid first-order performance and stimulus
538 factors confounding estimates of the confidence-accuracy correlation (Masson & Rotello,
539 2009; Rahnev & Fleming, 2019). Here we used a staircase procedure to control perceptual
540 performance within a narrow range and used a metric of metacognition that is unconfounded
541 by both metacognitive bias and first-order performance. In addition, we used a Bayesian
542 inference approach to estimate the impact of mentalizing ability on metacognitive ability
543 within the same hierarchical model, which ensured that both within- and between-subject
544 variability are appropriately taken into account. These methodological advances may explain
545 why here we found a more robust between-subjects relationship between metacognition and
546 mentalizing than reported previously (Carpenter et al., 2019; Nicholson et al., 2020).

547 Our results are also in line with previous work on autism, suggesting that
548 metacognitive ability may be compromised in autistic individuals to a similar extent to the
549 ability to evaluate other people’s mental states. Autism was characterised as a general “mind-
550 blindness” in 1985 (Simon Baron-Cohen et al., 1985) but, since then, only a handful of
551 studies have extended the study of mentalizing in autism to that of metacognitive ability
552 about one’s own behaviour and mental states (Carpenter et al., 2019; Grainger et al., 2016;
553 Nicholson et al., 2019, 2020; Williams et al., 2018; Wojcik et al., 2013). Some of these
554 studies (Grainger et al., 2016; Nicholson et al., 2020; Williams et al., 2018) but not others
555 (Carpenter et al., 2019; Wojcik et al., 2013), found, in line with our pre-registered hypotheses
556 and findings, that mentalizing and metacognitive ability were commensurately compromised
557 in autism. A notable exception to this general picture is that we unexpectedly found that self-

558 reported autistic traits on the RAADS-14 were not negatively associated with metacognitive
559 efficiency in our general population dataset of Experiment 1. One candidate explanation for
560 this inconsistency is that variation in autistic traits in the general population may not have
561 been pronounced enough to find statistically significant differences in metacognitive
562 efficiency. Another explanation is that metacognitive ability in autism may not be worse on
563 average but rather more extreme (both extremely strong *and* weak; Pariser, 1981; Shields-
564 Wolfe & Gallagher, 1992)—as hinted at by the greater variance in the autistic group
565 estimates (see overlaid dots in **Figure 4a**). Future studies should investigate whether this is
566 the case in larger samples and, if so, whether it can be attributed to autistic people engaging
567 in alternative, perhaps more cognitively demanding, processes to compensate for
568 metacognitive difficulties (Livingston, Colvert, et al., 2019; Livingston, Shah, et al., 2019).
569 Given the range of cues people may use to inform confidence, it will be important for future
570 studies to focus on how the construction of confidence or other mentalizing processes varies
571 across participants. It could be that, in real life, the metacognitive ability of some autistic
572 people is above average but achieved via different routes than those studied in this
573 experiment.

574 Our work goes beyond estimating correlations between metacognition and
575 mentalizing by revealing a potential mechanism through which mentalizing may affect
576 metacognitive processes. Specifically, we show that better mentalizing ability is associated
577 with a tighter coupling between response times and confidence. Previous work has
578 experimentally manipulated response times and found this to have a causal effect on the
579 construction of confidence: when response times are manipulated to be faster, people are
580 subsequently more likely to report being confident (Kiani et al., 2014; Palser et al., 2018).
581 The mentalizing-is-prior theory suggests metacognition consists of a re-application of
582 inferential processes used to understand other people to understand our own mental states

583 (Carruthers, 2009). Our findings are consistent with this view, showing that people with
584 greater proficiency in self-reported social skills and objectively measured mentalizing also
585 had better metacognitive efficiency. In addition, we found that mentalizing ability not only
586 correlated with overall metacognitive efficiency, but specifically with the ability to infer
587 confidence from behavioural cues that would also be visible markers of other people's
588 decision confidence in everyday situations. An important limitation of the current study is
589 that we cannot draw causal conclusions about how mentalizing affects metacognition or vice
590 versa. Future longitudinal work is needed to ask whether exposure to situations requiring
591 mental state inference from behaviour causally affects the development of explicit
592 metacognition. Another limitation of this study is that of domain-generalizability. There is reason
593 to believe that metacognitive efficiency measured from perceptual decision-making is similar
594 to metacognitive efficiency measured in other domains, such as from mnemonic or numerical
595 decision-making tasks (Bronfman et al., 2015; Rouault, McWilliams, et al., 2018; Talluri et
596 al., 2018; van der Plas et al., 2021). However, other studies found selective differences in
597 perceptual metacognition between groups, in the absence of differences in memory
598 metacognition (Fleming et al., 2014). The possibility of dissociations between domains
599 suggests an unlikely, albeit possible, chance that mentalizing ability is only related to
600 metacognitive efficiency when the latter is measured in the context of a perceptual task.
601 Future studies should test the interplay between metacognition and mentalizing across a
602 wider range of cognitive domains.

603 In summary, across two behavioural experiments we demonstrate that mentalizing
604 ability is associated with both greater metacognitive efficiency, and tighter links between
605 response times and confidence. In a general population sample, participants with better social
606 skills were also better at reflecting upon their own performance. In a second dataset we show
607 that autistic participants with generally lower mentalizing ability also had weaker

608 metacognitive ability, in the absence of differences in first-order performance. Together,
609 these results suggest that inferring other people's mental states is related to the ability to
610 evaluate our own decisions.

REFERENCES

- 611 Abell, F., Happé, F., & Frith, U. (2000). Do triangles play tricks? Attribution of mental states to
612 animated shapes in normal and abnormal development. *Cognitive Development, 15*(1), 1–16.
613 [https://doi.org/10.1016/S0885-2014\(00\)00014-9](https://doi.org/10.1016/S0885-2014(00)00014-9)
- 614 Allison, C., Auyeung, B., & Baron-Cohen, S. (2012). Toward Brief “Red Flags” for Autism
615 Screening: The Short Autism Spectrum Quotient and the Short Quantitative Checklist in
616 1,000 Cases and 3,000 Controls. *Journal of the American Academy of Child & Adolescent
617 Psychiatry, 51*(2), 202-212.e7. <https://doi.org/10.1016/j.jaac.2011.11.003>
- 618 Ashwood, K. L., Gillan, N., Horder, J., Hayward, H., Woodhouse, E., McEwen, F. S., Findon, J.,
619 Eklund, H., Spain, D., Wilson, C. E., Cadman, T., Young, S., Stoencheva, V., Murphy, C. M.,
620 Robertson, D., Charman, T., Bolton, P., Glaser, K., Asherson, P., ... Murphy, D. G. (2016).
621 Predicting the diagnosis of autism in adults using the Autism-Spectrum Quotient (AQ)
622 questionnaire. *Psychological Medicine, 46*(12), 2595–2604. PubMed.
623 <https://doi.org/10.1017/S0033291716001082>
- 624 Bang, D., Aitchison, L., Moran, R., Hecce Castanon, S., Rafiee, B., Mahmoodi, A., Lau, J. Y. F.,
625 Latham, P. E., Bahrami, B., & Summerfield, C. (2017). Confidence matching in group
626 decision-making. *Nature Human Behaviour, 1*(6), 0117. [https://doi.org/10.1038/s41562-017-](https://doi.org/10.1038/s41562-017-0117)
627 [0117](https://doi.org/10.1038/s41562-017-0117)
- 628 Bang, D., Ershadmanesh, S., Nili, H., & Fleming, S. M. (2020). Private-public mappings in human
629 prefrontal cortex. *BioRxiv*, 2020.02.21.954305. <https://doi.org/10.1101/2020.02.21.954305>
- 630 Baron-Cohen, S., Wheelwright, S., Hill, J., Raste, Y., & Plumb, I. (2001). The ‘Reading the Mind in
631 the Eyes’ Test revised version: A study with normal adults, and adults with Asperger
632 syndrome or high-functioning autism. *Journal of Child Psychology and Psychiatry, and
633 Allied Disciplines, 42*(2), 241–251.
- 634 Baron-Cohen, Simon, Leslie, A. M., & Frith, U. (1985). Does the autistic child have a “theory of
635 mind”? *Cognition, 21*(1), 37–46. [https://doi.org/10.1016/0010-0277\(85\)90022-8](https://doi.org/10.1016/0010-0277(85)90022-8)

- 636 Beck, A. T., Baruch, E., Balter, J. M., Steer, R. A., & Warman, D. M. (2004). A new instrument for
637 measuring insight: The Beck Cognitive Insight Scale. *Schizophrenia Research*, 68(2), 319–
638 329. [https://doi.org/10.1016/S0920-9964\(03\)00189-0](https://doi.org/10.1016/S0920-9964(03)00189-0)
- 639 Bertrams, A. (2021). Internal reliability, homogeneity, and factor structure of the ten-item Autism-
640 Spectrum Quotient (AQ-10) with two additional response categories. *Experimental Results*, 2,
641 e3. <https://doi.org/10.1017/exp.2020.70>
- 642 Broadbent, J., & Stokes, M. A. (2013). Removal of negative feedback enhances WCST performance
643 for individuals with ASD. *Research in Autism Spectrum Disorders*, 7(6), 785–792.
644 <https://doi.org/10.1016/j.rasd.2013.03.002>
- 645 Bronfman, Z. Z., Brezis, N., Moran, R., Tsetsos, K., Donner, T., & Usher, M. (2015). Decisions
646 reduce sensitivity to subsequent information. *Proceedings of the Royal Society B: Biological*
647 *Sciences*, 282(1810), 20150228. <https://doi.org/10.1098/rspb.2015.0228>
- 648 Campbell-Meiklejohn, D. K., Bach, D. R., Roepstorff, A., Dolan, R. J., & Frith, C. D. (2010). How
649 the Opinion of Others Affects Our Valuation of Objects. *Current Biology*, 20(13), 1165–
650 1170. <https://doi.org/10.1016/j.cub.2010.04.055>
- 651 Campbell-Meiklejohn, D., Simonsen, A., Frith, C. D., & Daw, N. D. (2017). Independent Neural
652 Computation of Value from Other People's Confidence. *The Journal of Neuroscience*, 37(3),
653 673. <https://doi.org/10.1523/JNEUROSCI.4490-15.2016>
- 654 Carpenter, K. L., Williams, D. M., & Nicholson, T. (2019). Putting Your Money Where Your Mouth
655 is: Examining Metacognition in ASD Using Post-decision Wagering. *Journal of Autism and*
656 *Developmental Disorders*, 49(10), 4268–4279. <https://doi.org/10.1007/s10803-019-04118-6>
- 657 Carruthers, P. (2009). How we know our own minds: The relationship between mindreading and
658 metacognition. *The Behavioral and Brain Sciences*, 32(2), 121–138; discussion 138-182.
659 <https://doi.org/10.1017/S0140525X09000545>
- 660 Condon, D. M., & Revelle, W. (2016). Selected ICAR Data from the SAPA-Project: Development
661 and Initial Validation of a Public-Domain Measure. *Journal of Open Psychology Data*, 4(1),
662 1. <https://doi.org/10.5334/jopd.25>

- 663 David, N., Gawronski, A., Santos, N. S., Huff, W., Lehnhardt, F.-G., Newen, A., & Vogeley, K.
664 (2008). Dissociation Between Key Processes of Social Cognition in Autism: Impaired
665 Mentalizing But Intact Sense of Agency. *Journal of Autism and Developmental Disorders*,
666 38(4), 593–605. <https://doi.org/10.1007/s10803-007-0425-x>
- 667 De Martino, B., Bobadilla-Suarez, S., Nouguchi, T., Sharot, T., & Love, B. C. (2017). Social
668 Information Is Integrated into Value and Confidence Judgments According to Its Reliability.
669 *The Journal of Neuroscience*, 37(25), 6066–6074. [https://doi.org/10.1523/JNEUROSCI.3880-](https://doi.org/10.1523/JNEUROSCI.3880-16.2017)
670 16.2017
- 671 Dimaggio, G., Lysaker, P. H., Carcione, A., Nicolò, G., & Semerari, A. (2008). Know yourself and
672 you shall know the other... to a certain extent: Multiple paths of influence of self-reflection
673 on mindreading. *Consciousness and Cognition*, 17(3), 778–789.
674 <https://doi.org/10.1016/j.concog.2008.02.005>
- 675 Eriksson, J. M., Andersen, L. M., & Bejerot, S. (2013). RAADS-14 Screen: Validity of a screening
676 tool for autism spectrum disorder in an adult psychiatric population. *Molecular Autism*, 4(1),
677 49. <https://doi.org/10.1186/2040-2392-4-49>
- 678 Ewbank, M. P., von dem Hagen, E. A. H., Powell, T. E., Henson, R. N., & Calder, A. J. (2016). The
679 effect of perceptual expectation on repetition suppression to faces is not modulated by
680 variation in autistic traits. *Cortex*, 80, 51–60. <https://doi.org/10.1016/j.cortex.2015.10.011>
- 681 Fandakova, Y., Selmecky, D., Leckey, S., Grimm, K. J., Wendelken, C., Bunge, S. A., & Ghetti, S.
682 (2017). Changes in ventromedial prefrontal and insular cortex support the development of
683 metamemory from childhood into adolescence. *Proceedings of the National Academy of*
684 *Sciences*, 114(29), 7582. <https://doi.org/10.1073/pnas.1703079114>
- 685 Flavell, J. H. (1979). Metacognition and cognitive monitoring: A new area of cognitive–
686 developmental inquiry. *American Psychologist*, 34(10), 906–911.
687 <https://doi.org/10.1037/0003-066X.34.10.906>
- 688 Fleming, S. M., Weil, R. S., Nagy, Z., Dolan, R. J., & Rees, G. (2010). Relating Introspective
689 Accuracy to Individual Differences in Brain Structure. *Science*, 329(5998), 1541–1543.
690 <https://doi.org/10.1126/science.1191883>

- 691 Fleming, S. M. (2017a). HMeta-d: Hierarchical Bayesian estimation of metacognitive efficiency from
692 confidence ratings. *Neuroscience of Consciousness*, 2017(1).
693 <https://doi.org/10.1093/nc/nix007>
- 694 Fleming, S. M. (2017b). HMeta-d: Hierarchical Bayesian estimation of metacognitive efficiency from
695 confidence ratings. *Neuroscience of Consciousness*, 2017(1).
696 <https://doi.org/10.1093/nc/nix007>
- 697 Fleming, S. M., & Daw, N. D. (2017). Self-evaluation of decision-making: A general Bayesian
698 framework for metacognitive computation. *Psychological Review*, 124(1), 91–114.
699 <https://doi.org/10.1037/rev0000045>
- 700 Fleming, S. M., & Lau, H. C. (2014). How to measure metacognition. *Frontiers in Human*
701 *Neuroscience*, 8. <https://doi.org/10.3389/fnhum.2014.00443>
- 702 Fleming, S. M, Ryu, J., Golfinos, J. G., & Blackmon, K. E. (2014). Domain-specific impairment in
703 metacognitive accuracy following anterior prefrontal lesions. *Brain : A Journal of Neurology*,
704 137(Pt 10), 2811–2822. PubMed. <https://doi.org/10.1093/brain/awu221>
- 705 Fleming, S. M., van der Putten, E. J., & Daw, N. D. (2018). Neural mediators of changes of mind
706 about perceptual decisions. *Nature Neuroscience*, 21(4), 617–624.
707 <https://doi.org/10.1038/s41593-018-0104-6>
- 708 Frith, C. D. (2012). The role of metacognition in human social interactions. *Philosophical*
709 *Transactions of the Royal Society B: Biological Sciences*, 367(1599), 2213–2223.
710 <https://doi.org/10.1098/rstb.2012.0123>
- 711 Gazzaniga, M. S. (1995). Principles of human brain organization derived from split-brain studies.
712 *Neuron*, 14(2), 217–228. [https://doi.org/10.1016/0896-6273\(95\)90280-5](https://doi.org/10.1016/0896-6273(95)90280-5)
- 713 Gazzaniga, M. S. (2000). Cerebral specialization and interhemispheric communication: Does the
714 corpus callosum enable the human condition? *Brain*, 123(7), 1293–1326.
715 <https://doi.org/10.1093/brain/123.7.1293>
- 716 Gelman, A., & Rubin, D. B. (1992). Inference from Iterative Simulation Using Multiple Sequences.
717 *Statist. Sci.*, 7(4), 457–472. <https://doi.org/10.1214/ss/1177011136>

- 718 Gopnik, A. (1993). How we know our minds: The illusion of first-person knowledge of intentionality.
719 *Behavioral and Brain Sciences*, *16*(1), 1–14, 29–113.
720 <https://doi.org/10.1017/S0140525X00028636>
- 721 Goupil, L., & Kouider, S. (2016). Behavioral and Neural Indices of Metacognitive Sensitivity in
722 Preverbal Infants. *Current Biology : CB*, *26*(22), 3038–3045.
723 <https://doi.org/10.1016/j.cub.2016.09.004>
- 724 Goupil, L., & Kouider, S. (2019). Developing a Reflective Mind: From Core Metacognition to
725 Explicit Self-Reflection. *Current Directions in Psychological Science*, *28*(4), 403–408.
726 <https://doi.org/10.1177/0963721419848672>
- 727 Grainger, C., Williams, D. M., & Lind, S. E. (2016). Metacognitive monitoring and control processes
728 in children with autism spectrum disorder: Diminished judgement of confidence accuracy.
729 *Consciousness and Cognition*, *42*, 65–74. <https://doi.org/10.1016/j.concog.2016.03.003>
- 730 Greene, R. K., Zheng, S., Kinard, J. L., Mosner, M. G., Wiesen, C. A., Kennedy, D. P., & Dichter, G.
731 S. (2019). Social and nonsocial visual prediction errors in autism spectrum disorder. *Autism*
732 *Research*, *12*(6), 878–883. <https://doi.org/10.1002/aur.2090>
- 733 Happe, F. (2015). Autism as a neurodevelopmental disorder of mind-reading. *Journal of the British*
734 *Academy*, *3*. <https://doi.org/10.5871/jba/003.197>
- 735 Harrison, O. K., Garfinkel, S. N., Marlow, L., Finnegan, S., Marino, S., Nanz, L., Allen, M.,
736 Finnemann, J., Keur-Huizinga, L., Harrison, S. J., Stephan, K. E., Pattinson, K., & Fleming,
737 S. M. (2020). The Filter Detection Task for measurement of breathing-related interoception
738 and metacognition. *BioRxiv*, 2020.06.29.176941. <https://doi.org/10.1101/2020.06.29.176941>
- 739 Hembacher, E., & Ghetti, S. (2014). Don't look at my answer: Subjective uncertainty underlies
740 preschoolers' exclusion of their least accurate memories. *Psychological Science*, *25*(9),
741 1768–1776. <https://doi.org/10.1177/0956797614542273>
- 742 Kiani, R., Corthell, L., & Shadlen, M. N. (2014). Choice Certainty Is Informed by Both Evidence and
743 Decision Time. *Neuron*, *84*(6), 1329–1342. <https://doi.org/10.1016/j.neuron.2014.12.015>

- 744 Kiani, R., & Shadlen, M. N. (2009). Representation of Confidence Associated with a Decision by
745 Neurons in the Parietal Cortex. *Science*, 324(5928), 759.
746 <https://doi.org/10.1126/science.1169405>
- 747 Kruschke, J. K. (2010). *Doing Bayesian Data Analysis: A Tutorial with R and BUGS* (1st ed.).
748 Academic Press, Inc.
- 749 Lieder, I., Adam, V., Frenkel, O., Jaffe-Dax, S., Sahani, M., & Ahissar, M. (2019). Perceptual bias
750 reveals slow-updating in autism and fast-forgetting in dyslexia. *Nature Neuroscience*, 22(2),
751 256–264. <https://doi.org/10.1038/s41593-018-0308-9>
- 752 Livingston, L. A., & Happé, F. (in press). Understanding atypical development through social
753 cognitive theory: Lessons from autism. In *The cognitive basis of social interaction across the*
754 *lifespan*. H. J. Ferguson, V. E. A. Brunson & E.E.F. Bradford (Eds.). Oxford University
755 Press.
- 756 Livingston, L. A., Carr, B., & Shah, P. (2019). Recent Advances and New Directions in Measuring
757 Theory of Mind in Autistic Adults. *Journal of Autism and Developmental Disorders*, 49(4),
758 1738–1744. <https://doi.org/10.1007/s10803-018-3823-3>
- 759 Livingston, L. A., Colvert, E., the Social Relationships Study Team, Bolton, P., & Happé, F. (2019).
760 Good social skills despite poor theory of mind: Exploring compensation in autism spectrum
761 disorder. *Journal of Child Psychology and Psychiatry*, 60(1), 102–110.
762 <https://doi.org/10.1111/jcpp.12886>
- 763 Livingston, L.A., Shah, P., & Happé, F. (2019). Compensatory strategies below the behavioural
764 surface in autism: A qualitative study. *The Lancet Psychiatry*, 6(9), 766–777.
765 [https://doi.org/10.1016/S2215-0366\(19\)30224-X](https://doi.org/10.1016/S2215-0366(19)30224-X)
- 766 Livingston, L. A., Shah, P., White, S., & Happé, F. (in press). Further developing the Frith-
767 Happé animations: A quicker, more objective, and web-based test of theory of mind
768 for autistic and neurotypical adults. *Autism Research*.
- 769 Maniscalco, B., & Lau, H. (2012). A signal detection theoretic approach for estimating metacognitive
770 sensitivity from confidence ratings. *Consciousness and Cognition*, 21(1), 422–430.
771 <https://doi.org/10.1016/j.concog.2011.09.021>
772

- 773 Maniscalco, B., & Lau, H. (2014). Signal detection theory analysis of type 1 and type 2 data: Meta-d',
774 response-specific meta-d', and the unequal variance SDT model. *The Cognitive Neuroscience*
775 *of Metacognition.*, 25–66. https://doi.org/10.1007/978-3-642-45190-4_3
- 776 Masson, M. E. J., & Rotello, C. M. (2009). Sources of bias in the Goodman-Kruskal gamma
777 coefficient measure of association: Implications for studies of metacognitive processes.
778 *Journal of Experimental Psychology. Learning, Memory, and Cognition*, 35(2), 509–527.
779 <https://doi.org/10.1037/a0014876>
- 780 McMahon, C. M., Henderson, H. A., Newell, L., Jaime, M., & Mundy, P. (2016). Metacognitive
781 Awareness of Facial Affect in Higher-Functioning Children and Adolescents with Autism
782 Spectrum Disorder. *Journal of Autism and Developmental Disorders*, 46(3), 882–898.
783 <https://doi.org/10.1007/s10803-015-2630-3>
- 784 Milne, E., Swettenham, J., Hansen, P., Campbell, R., Jeffries, H., & Plaisted, K. (2002). High motion
785 coherence thresholds in children with autism. *Journal of Child Psychology and Psychiatry*,
786 43(2), 255–263. <https://doi.org/10.1111/1469-7610.00018>
- 787 Nelson, T. O. (1984). A comparison of current measures of the accuracy of feeling-of-knowing
788 predictions. *Psychological Bulletin*, 95(1), 109–133. [https://doi.org/10.1037/0033-](https://doi.org/10.1037/0033-2909.95.1.109)
789 2909.95.1.109
- 790 Nicholson, T., Williams, D., Lind, S., Grainger, C., & Carruthers, P. (2020). Linking metacognition
791 and mindreading: Evidence from autism and dual-task investigations. *Journal of*
792 *Experimental Psychology. General*. <https://doi.org/10.1037/xge0000878>
- 793 Nicholson, T., Williams, D. M., Grainger, C., Lind, S. E., & Carruthers, P. (2019). Relationships
794 between implicit and explicit uncertainty monitoring and mindreading: Evidence from autism
795 spectrum disorder. *Consciousness and Cognition*, 70, 11–24.
796 <https://doi.org/10.1016/j.concog.2019.01.013>
- 797 Nisbett, R. E., & Wilson, T. D. (1977). Telling more than we can know: Verbal reports on mental
798 processes. *Psychological Review*, 84(3), 231–259. [https://doi.org/10.1037/0033-](https://doi.org/10.1037/0033-295X.84.3.231)
799 295X.84.3.231

- 800 Palan, S., & Schitter, C. (2018). Prolific.ac—A subject pool for online experiments. *Journal of*
801 *Behavioral and Experimental Finance*, *17*, 22–27. <https://doi.org/10.1016/j.jbef.2017.12.004>
- 802 Palser, E. R., Fotopoulou, A., & Kilner, J. M. (2018). Altering movement parameters disrupts
803 metacognitive accuracy. *Consciousness and Cognition: An International Journal*, *57*, 33–40.
804 <https://doi.org/10.1016/j.concog.2017.11.005>
- 805 Pariser, D. (1981). Nadia's Drawings: Theorizing about an Autistic Child's Phenomenal Ability.
806 *Studies in Art Education*, *22*(2), 20–31. <https://doi.org/10.1080/00393541.1981.11650279>
- 807 Patel, D., Fleming, S. M., & Kilner, J. M. (2012). Inferring subjective states through the observation
808 of actions. *Proceedings of the Royal Society B: Biological Sciences*, *279*(1748), 4853–4860.
809 <https://doi.org/10.1098/rspb.2012.1847>
- 810 Pirrone, A., Dickinson, A., Gomez, R., Stafford, T., & Milne, E. (2017). Understanding perceptual
811 judgment in autism spectrum disorder using the drift diffusion model. *Neuropsychology*,
812 *31*(2), 173–180. <https://doi.org/10.1037/neu0000320>
- 813 Pleskac, T. J., & Busemeyer, J. R. (2010). Two-stage dynamic signal detection: A theory of choice,
814 decision time, and confidence. *Psychological Review*, *117*(3), 864–901.
815 <https://doi.org/10.1037/a0019737>
- 816 Rahnev, D., & Fleming, S. M. (2019). How experimental procedures influence estimates of
817 metacognitive ability. *Neuroscience of Consciousness*, *2019*(niz009).
818 <https://doi.org/10.1093/nc/niz009>
- 819 Reed, P. (2019). Unpredictability reduces over-selective responding of individuals with ASD who
820 have language impairments. *Research in Autism Spectrum Disorders*, *57*, 35–45.
821 <https://doi.org/10.1016/j.rasd.2018.10.006>
- 822 Resulaj, A., Kiani, R., Wolpert, D. M., & Shadlen, M. N. (2009). Changes of mind in decision-
823 making. *Nature*, *461*(7261), 263–266. <https://doi.org/10.1038/nature08275>
- 824 Robic, S., Sonié, S., Fonlupt, P., Henaff, M.-A., Touil, N., Coricelli, G., Mattout, J., & Schmitz, C.
825 (2015). Decision-Making in a Changing World: A Study in Autism Spectrum Disorders.
826 *Journal of Autism and Developmental Disorders*, *45*(6), 1603–1613.
827 <https://doi.org/10.1007/s10803-014-2311-7>

- 828 Rollwage, M., Dolan, R. J., & Fleming, S. M. (2018). Metacognitive Failure as a Feature of Those
829 Holding Radical Beliefs. *Current Biology*, 28(24), 4014-4021.e8.
830 <https://doi.org/10.1016/j.cub.2018.10.053>
- 831 Rosenblau, G., Kliemann, D., Heekeren, H. R., & Dziobek, I. (2015). Approximating Implicit and
832 Explicit Mentalizing with Two Naturalistic Video-Based Tasks in Typical Development and
833 Autism Spectrum Disorder. *Journal of Autism and Developmental Disorders*, 45(4), 953–965.
834 <https://doi.org/10.1007/s10803-014-2249-9>
- 835 Rouault, M., McWilliams, A., Allen, M. G., & Fleming, S. M. (2018). Human Metacognition Across
836 Domains: Insights from Individual Differences and Neuroimaging. *Personality Neuroscience*,
837 1. <https://doi.org/10.1017/pen.2018.16>
- 838 Rouault, M., Seow, T., Gillan, C. M., & Fleming, S. M. (2018). Psychiatric Symptom Dimensions Are
839 Associated With Dissociable Shifts in Metacognition but Not Task Performance. *Biological*
840 *Psychiatry*, 84(6), 443–451. <https://doi.org/10.1016/j.biopsych.2017.12.017>
- 841 Ryle, G. (1949). Meaning and Necessity. *Philosophy*, 24(88), 69–76. JSTOR.
- 842 Sapey-Triomphe, L.-A., Sonié, S., Hénaff, M.-A., Mattout, J., & Schmitz, C. (2018). Adults with
843 Autism Tend to Undermine the Hidden Environmental Structure: Evidence from a Visual
844 Associative Learning Task. *Journal of Autism and Developmental Disorders*, 48(9), 3061–
845 3074. <https://doi.org/10.1007/s10803-018-3574-1>
- 846 Shields-Wolfe, J., & Gallagher, P. A. (1992). Functional utilization of splinter skills for the
847 employment of a young adult with autism. *Focus on Autistic Behavior*, 7(4), 1–16.
- 848 Spek, A., Schatorjé, T., Scholte, E., & van Berckelaer-Onnes, I. (2009). Verbal fluency in adults with
849 high functioning autism or Asperger syndrome. *Neuropsychologia*, 47(3), 652–656.
850 <https://doi.org/10.1016/j.neuropsychologia.2008.11.015>
- 851 Talluri, B. C., Urai, A. E., Tsetsos, K., Usher, M., & Donner, T. H. (2018). Confirmation Bias through
852 Selective Overweighting of Choice-Consistent Evidence. *Current Biology*, 28(19), 3128-
853 3135.e8. <https://doi.org/10.1016/j.cub.2018.07.052>

- 854 Vaccaro, A. G., & Fleming, S. M. (2018). Thinking about thinking: A coordinate-based meta-analysis
855 of neuroimaging studies of metacognitive judgements. *Brain and Neuroscience Advances*, 2,
856 239821281881059. <https://doi.org/10.1177/2398212818810591>
- 857 van den Berg, R., Zylberberg, A., Kiani, R., Shadlen, M. N., & Wolpert, D. M. (2016). Confidence Is
858 the Bridge between Multi-stage Decisions. *Current Biology*, 26(23), 3157–3168.
859 <https://doi.org/10.1016/j.cub.2016.10.021>
- 860 van der Plas, E. A. A., Shiqi, Z., Keer, D., Bang, D., Nicholas, W., Jian, L., & Stephen, F. (2021).
861 *Isolating cultural contributors to confidence. (2021, March 8).*
862 <https://doi.org/10.31234/osf.io/sjh7d>
- 863 Wegner, D. M. (2002). *The illusion of conscious will.* (pp. xi, 405). MIT Press.
- 864 Weil, L. G., Fleming, S. M., Dumontheil, I., Kilford, E. J., Weil, R. S., Rees, G., Dolan, R. J., &
865 Blakemore, S.-J. (2013). The development of metacognitive ability in adolescence.
866 *Consciousness and Cognition*, 22(1), 264–271. <https://doi.org/10.1016/j.concog.2013.01.004>
- 867 White, S., Hill, E., Happé, F., & Frith, U. (2009). Revisiting the Strange Stories: Revealing
868 Mentalizing Impairments in Autism. *Child Development*, 80(4), 1097–1117.
869 <https://doi.org/10.1111/j.1467-8624.2009.01319.x>
- 870 White, S. J., Coniston, D., Rogers, R., & Frith, U. (2011). Developing the Frith-Happé animations: A
871 quick and objective test of Theory of Mind for adults with autism. *Autism Research*, 4(2),
872 149–154. <https://doi.org/10.1002/aur.174>
- 873 Williams, D. M., Bergström, Z., & Grainger, C. (2018a). Metacognitive monitoring and the
874 hypercorrection effect in autism and the general population: Relation to autism(-like) traits
875 and mindreading. *Autism*, 22(3), 259–270. <https://doi.org/10.1177/1362361316680178>
- 876 Williams, D. M., Bergström, Z., & Grainger, C. (2018b). Metacognitive monitoring and the
877 hypercorrection effect in autism and the general population: Relation to autism(-like) traits
878 and mindreading. *Autism*, 22(3), 259–270. <https://doi.org/10.1177/1362361316680178>
- 879 Wilson, M. (2002). Six views of embodied cognition. *Psychonomic Bulletin & Review*, 9(4), 625–636.
880 <https://doi.org/10.3758/BF03196322>

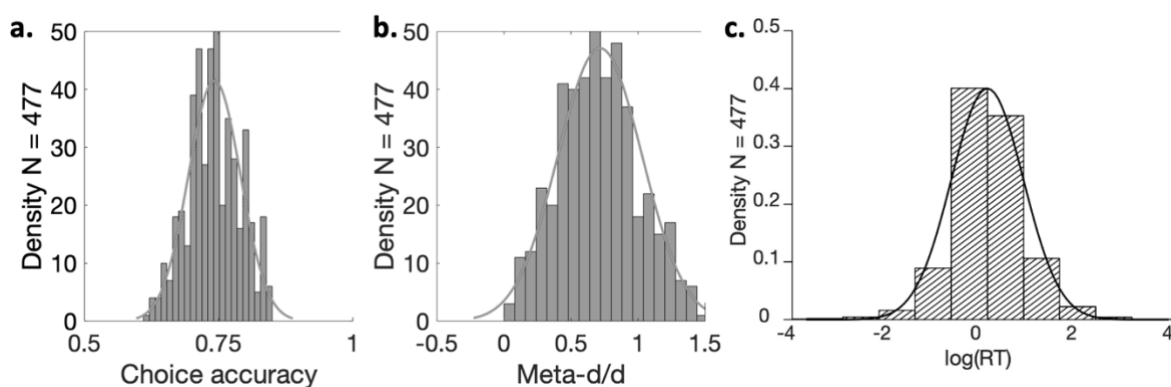
- 881 Wojcik, D. Z., Moulin, C. J. A., & Souchay, C. (2013). Metamemory in children with autism:
882 Exploring “feeling-of-knowing” in episodic and semantic memory. *Neuropsychology*, *27*(1),
883 19–27. <https://doi.org/10.1037/a0030526>
- 884 Yeung, N., & Summerfield, C. (2012). Metacognition in human decision-making: Confidence and
885 error monitoring. *Philosophical Transactions of the Royal Society B: Biological Sciences*,
886 *367*(1594), 1310–1321. <https://doi.org/10.1098/rstb.2011.0416>
- 887 Zalla, T., Miele, D., Leboyer, M., & Metcalfe, J. (2015). Metacognition of agency and theory of mind
888 in adults with high functioning autism. *Consciousness and Cognition*, *31*, 126–138.
889 <https://doi.org/10.1016/j.concog.2014.11.001>
- 890 Zwart, F. S., Vissers, C. Th. W. M., Kessels, R. P. C., & Maes, J. H. R. (2018). Implicit learning
891 seems to come naturally for children with autism, but not for children with specific language
892 impairment: Evidence from behavioral and ERP data: Implicit learning intact in ASD but
893 altered in SLI. *Autism Research*, *11*(7), 1050–1061. <https://doi.org/10.1002/aur.1954>
- 894

SUPPLEMENTARY MATERIAL

Experiment 1

1.1 Performance and validation checks

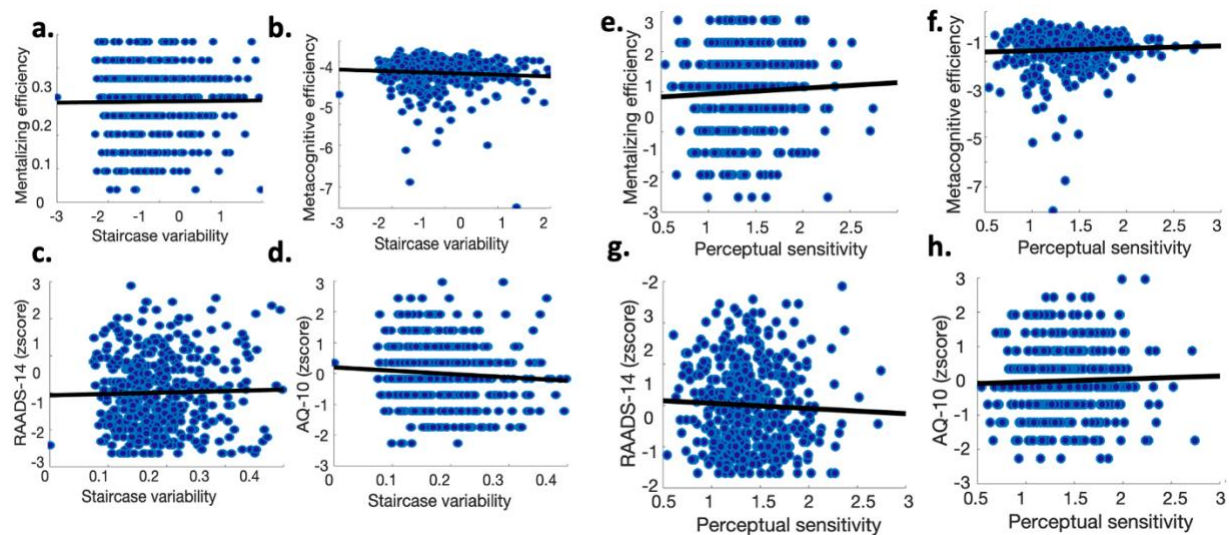
Average choice accuracy on the metacognition task ($M = 74.16\%$, $SEM = 0.002$; **Supplementary Figure 1.1a**), metacognitive efficiency ($M = 0.693$, $SEM = 0.016$; **Supplementary Figure 1.1b**) and the log of response times ($\log RT$; $M = -1.405e^{-07}$, $SEM = 0.046$; **Supplementary Figure 1.1c**) were similar to those of previous studies (Rouault et al., 2018; Rollwage et al., 2018).



Supplementary Figure 1.1. Choice accuracy on the metacognition task. *a. Histogram distribution of choice accuracy. b. Histogram distribution of metacognitive efficiency (meta- d'/d'). c. Histogram distribution of the log of standardized response times ($\log RT$). All variables are derived from the metacognition task and plotted for the group as a whole ($N=477$).*

As an indication of the reliability of mentalizing task variables, we asked whether the two mentalizing measures from the Happé-Frith Triangle Task were measuring a similar mentalizing construct. This was the case, with a positive correlation between the mentalizing feelings and mentalizing category scores: Spearman's $r = 0.37$, $P = 2.73e-16$. In addition, to establish whether the autistic trait surveys and Frith-Happé triangle task were measuring a similar mentalizing construct, we tested whether people with more autistic traits on the mentalizing subscale of the RAADS-14 also had lower mentalizing ability on the Frith-Happé Triangle Task, which was also the case (Spearman's $r = -0.11$, $P = 0.017$).

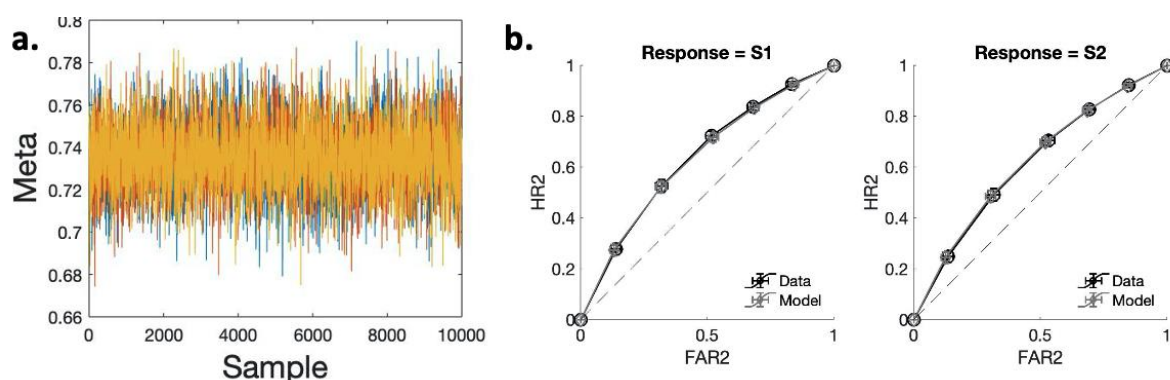
We next sought to ensure key variables related to metacognition and mentalizing were independent of first-order perceptual task performance. We first calculated each individual's experienced stimulus variability (the ratio between the standard deviation of stimulus difficulty and average stimulus difficulty) and correlated this with the main variables of interest. Staircase variability was not correlated with mentalizing ability ($r_{S475} = 0.005$, $P = 0.91$; **Supplementary Figure 1.2a**) metacognitive efficiency ($r_{S475} = -0.068$, $P = 0.137$; **Supplementary Figure 1.2b**), RAADS-14 scores ($r_{S475} = 0.0015$, $P = 0.974$; **Supplementary Figure 1.2c**) or AQ-10 scores ($r_{S475} = -0.066$, $P = 0.149$; **Supplementary Figure 1.2d**). The same validation checks were conducted for perceptual sensitivity, which was not correlated with mentalizing ability ($r_{S475} = 0.0655$, $P = 0.1524$; **Supplementary Figure 1.2e**) metacognitive efficiency ($r_{S475} = -0.0513$, $P = 0.264$; **Supplementary Figure 1.2f**), RAADS-14 scores ($r_{S475} = -0.0536$, $P = 0.2437$; **Supplementary Figure 1.2g**) or AQ-10 scores ($r_{S475} = 0.0359$, $P = 0.435$; **Supplementary Figure 1.2h**).



Supplementary Figure 1.2. Correlations between the main variables of interest. a-d: Staircase variability, the ratio of the standard deviation and the mean dot difference, was not correlated with **a.** mentalizing ability, **b.** metacognitive efficiency ($meta-d'/d'$), **c.** autistic traits as measured by the RAADS-14 **d.** autistic traits as measured with the AQ-10. **e-h:** Perceptual sensitivity (d') was not correlated with **e.** mentalizing ability, **f.** metacognitive efficiency ($meta-d'/d'$), **g.** autistic traits as measured by the RAADS-14, **h.** autistic traits as measured with the AQ-10.

1.2. Posterior predictive checks

Next, we test whether the HMeta-d models used in estimating metacognitive efficiency were reliable by means of convergence checks and posterior predictive checks. The hierarchical regression model predicting metacognition from mentalizing ability scores converged well, indicated by the Gelman-Rubic statistics ($\hat{R} = 0.99997$ and see plotted chains in **Supplementary Figure 1.3a**). In addition, posterior predictive plots captured key patterns of the participants' confidence responses, with model and predicted type ROCs closely overlapping (**Supplementary Figure 1.3b**). The same was true for the hierarchical regression models with RAADS-14 scores ($\hat{R} = 1.0003$ **Supplementary Figure 1.3c, d**) and AQ-10 scores ($\hat{R} = 1.0006$ **Supplementary Figure 1.3e, f**).



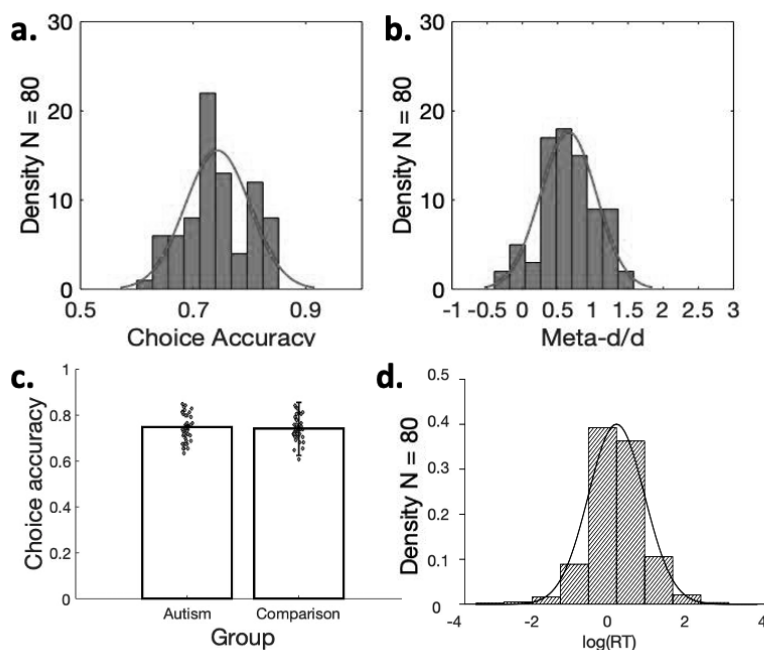
Supplementary Figure 1.3. Posterior predictive checks on HMeta-d fits in Experiment 1. a. MCMC chains for parameter meta-d'/d' (metacognitive efficiency) from the hierarchical regression model. **b.** Observed and model estimates for the Type 2 ROC curves for leftward (S1) and rightward (S2) responses from the regression meta-d model fits. Error bars represent the mean \pm standard error of the mean.

Experiment 2

2.1. Performance and validation checks

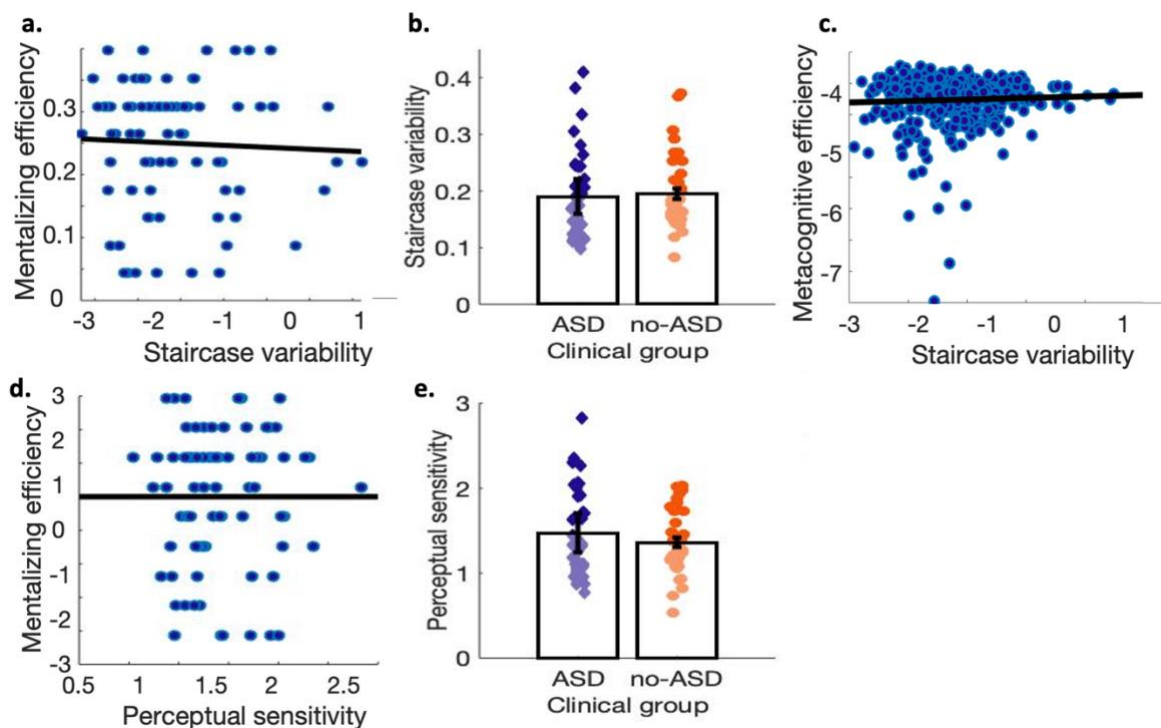
Average choice accuracy on the metacognition task ($M = 74.34\% \pm 0.006$) was normally distributed ($W = 0.98$, $P = 0.12$; **Supplementary Figure 2.1a**) and is visually similar to those of the larger dataset (**Supplementary Figure 1.1**). Metacognitive efficiency or meta - d'/d'

($M=0.653 \pm 0.045$) was also normally distributed ($W = 0.987$, $P = 0.60$; **Supplementary Figure 2.1b**) and similar to that in Experiment 1 (**Supplementary Figure 1.1**). As a result of the calibration procedure, first-order performance was not statistically different between the autism ($M = 0.75 \pm 0.01$) and comparison groups ($M = 0.74 \pm 0.008$; equal variances: $P = 0.73$, $K = 0.15$; independent samples t-test, $t_{78} = 0.519$, 95% CI = [-0.019, 0.032], $P = 0.61$; **Supplementary Figure 2.1c**). Finally, we averaged the log of response times (logRT) across trials of the metacognition task for each subject and plotted the distribution in **Supplementary Figure 2.1d**. Average logRT in the autism group ($M= -7.39e^{-17} \pm 6.05e^{-17}$) and in the comparison group ($M= -2.59e^{-17} \pm 6.17e^{-17}$) were not statistically different ($t_{71} = 0.49$, 95% CI = [-2.43, 1.47], $P = 0.63$).



Supplementary Figure 2.1. Choice accuracy on the metacognition task. *a.* Histogram distribution of choice accuracy on the metacognition task in the group as a whole ($N=80$). *b.* Histogram distribution of metacognitive efficiency (meta- d'/d') on the metacognition task in the group as a whole ($N=80$). *c.* Average choice accuracy was matched for autism ($N=40$) and comparison participants ($N=40$) on the metacognition task. Error bars represent group mean \pm SEM. *d.* Histogram distribution of the log of standardized response times (logRT) on the metacognition task in the group as a whole ($N=80$).

We again sought to ensure key variables related to metacognition and mentalizing were independent of first-order perceptual task performance. Staircase variability was not correlated with mentalizing ability ($r_{s78} = -0.044$, $P = 0.71$; **Supplementary Figure 2.2a**) and was not statistically different between groups (95% CI = [-0.036, 0.026], $t_{78} = -0.31$, $P = 0.756$; **Supplementary Figure 2.2b**). In addition, staircase variability was not correlated with metacognitive efficiency ($r_{s78} = 0.031$, $P = 0.782$; **Supplementary Figure 2.2c**). Perceptual sensitivity (d') was not correlated with mentalizing ability ($r_{s78} = 0.011$, $P = 0.924$; **Figure 3.2d**) and was not statistically different between groups (95% CI = [-0.083, 0.309], $t_{78} = 1.15$, $P = 0.253$; **Supplementary Figure 2.2e**).



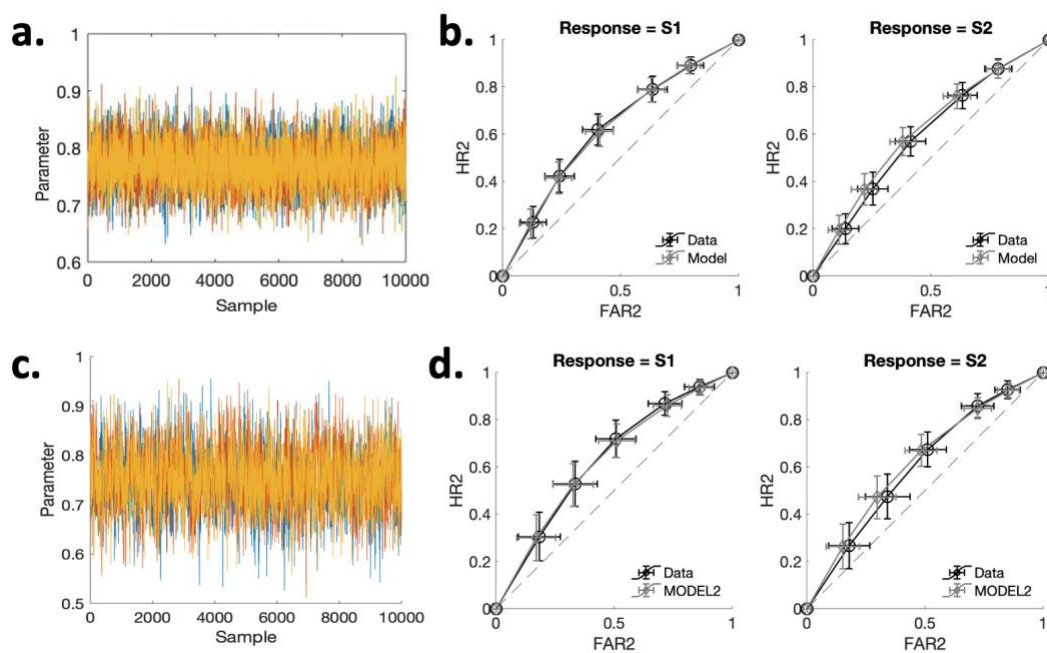
Supplementary Figure 2.2. Correlations between the main variables of interest. a. Mentalizing ability and staircase variability in the sample as a whole ($N=80$) were not correlated. **b.** Staircase variability was not different between the autism ($N=40$) and comparison groups ($N=40$). **c.** Metacognitive efficiency ($meta-d'/d'$) and staircase variability in the sample as a whole ($N=80$) were not correlated. **d.** Mentalizing ability and perceptual ability (d') in the sample as a whole were not correlated. **e.** Perceptual ability (d') was not statistically different between autism ($N=40$) and comparison participants ($N=40$). Error bars represent the group means \pm SEM.

2.2 Posterior predictive checks

Finally, we asked whether the two HMeta-d models fitted to Experiment 2 data were reliable by means of convergence checks and posterior predictive checks. The hierarchical regression model converged well, indicated by the Gelman-Rubic statistics ($\hat{R}_{Mratio}=1.0001$ and plotted chains in **Supplementary Figure 2.3a**). In addition, posterior predictive plots recaptured key patterns of the participants' confidence responses correctly (**Supplementary Figure 2.3b**).

The same was true for separate model fits to the comparison group ($\hat{R}_{Mratio}=1.0014$,

Supplementary Figure 2.3c, d) and autism group ($\hat{R}_{Mratio}=1.002$, **Supplementary Figure 2.3e, f**).



Supplementary Figure 2.3. Posterior predictive checks on HMeta-d fits in Experiment 2. a.

895 MCMC chains for parameter meta-d'/d' (metacognitive efficiency) from the hierarchical
 896 regression model on autistic participants' data ($N = 40$) and **c.** on comparison participants'
 897 data ($N = 40$). **b.** Observed and model estimates for the Type 2 ROC curves for leftward (S1)
 898 and rightward (S2) responses from the hierarchical regression model are plotted for autistic
 899 participants' data ($N=40$.) and **d.** on comparison participants' data ($N = 40$). Error bars
 900 represent the mean \pm standard error of the mean.

SUPPLEMENTARY REFERENCES

- 901 Rahnev, D., & Fleming, S. M. (2019). How experimental procedures influence estimates of
902 metacognitive ability. *Neuroscience of Consciousness*, 2019(niz009).
903 <https://doi.org/10.1093/nc/niz009>
- 904 Rollwage, M., Dolan, R. J., & Fleming, S. M. (2018). Metacognitive Failure as a Feature of Those
905 Holding Radical Beliefs. *Current Biology*, 28(24), 4014-4021.e8.
906 <https://doi.org/10.1016/j.cub.2018.10.053>
- 907 Rouault, M., Seow, T., Gillan, C. M., & Fleming, S. M. (2018). Psychiatric Symptom Dimensions Are
908 Associated With Dissociable Shifts in Metacognition but Not Task Performance. *Biological*
909 *Psychiatry*, 84(6), 443–451. <https://doi.org/10.1016/j.biopsych.2017.12.017>