1    **Egocentric Biases are Determined by the Precision of Self-related Predictions.**

2    Leora Sevi[1*], Mirta Stantic[1], Jennifer Murphy[2], Michel-Pierre Coll[3], Caroline Catmur[4] &
3    Geoffrey Bird[1,5*]

4    [1] Department of Experimental Psychology, University of Oxford, United Kingdom

5    [2] Department of Psychology, Royal Holloway, University of London, United Kingdom

6    [3] Department of Psychology, McGill University, Canada

7    [4] Department of Psychology, Institute of Psychiatry, Psychology & Neuroscience, King's

8    College London, United Kingdom

9    [5] Social, Genetic & Developmental Psychiatry Centre, Institute of Psychiatry, Psychology &

10    Neuroscience, King's College London, United Kingdom

11    *Correspondence: leora.sevi@psy.ox.ac.uk and geoff.bird@psy.ox.ac.uk.

## Abstract

According to predictive processing theories, emotional inference involves simultaneously minimising discrepancies between predictions and sensory data relating to both one's own and others' states, achievable by altering either one's own state (empathy) or perception of another's state (egocentric bias) so they are more congruent. We tested a key hypothesis of these accounts, that predictions are weighted in inference according to their precision (inverse variance). If correct, more precise self-related predictions should bias perception of another's emotional expression to a greater extent than less precise predictions. We manipulated predictions about upcoming own-pain (low or high magnitude) using cues that afforded either precise (a narrow range of possible magnitudes) or imprecise (a wide range) predictions. Participants judged pained facial expressions presented concurrently with own-pain to be more intense when own-pain was greater, and precise cues increased this biasing effect. Implications of conceptualising interpersonal influence in terms of predictive processing are discussed.

*Keywords*: emotion recognition, predictive coding, empathy, generative models, precision, predictive interoceptive coding

29    The notion that the brain is an inferential machine, generating predictions to explain

30    the sensory data it receives in order to test models about the state of the world, is becoming

31    increasingly influential in cognitive neuroscience. Within the predictive processing

32    framework (Clark, 2016; Friston, 2010; Hohwy, 2013), the brain continually tests predictions

33    about the world, generated by models, against incoming sensory data. Discrepancies between

34    predictions and sensory data (prediction errors) are resolved either through, 1) the updating of

35    models generating predictions such that they better fit sensory data, or 2) the performance of

36    'action' (whether cognitive, motoric or interoceptive), in order to minimise the discrepancy

37    between predictions and sensory data. Which of these strategies are enacted in order to

38    reduce prediction errors is a function of the relative expected precision (uncertainty,

39    confidence — or, mathematically, inverse variance) of predictions and prediction errors: if

40    prediction errors are more precise than predictions then models are updated; if not, action

41    occurs.

42    One feature of the models specified by predictive processing theories is that they are

43    hierarchical; at lower levels, they attempt to explain unimodal sensory data, whereas at higher

44    levels they are multimodal, generating exteroceptive (e.g., visual, auditory), proprioceptive,

45    and interoceptive (e.g., hunger, satiety, pain) predictions. These higher levels allow

46    predictions and prediction errors relating to contingent events in different modalities to

47    contextualise each other, allowing for more abstract representations of a cause of sensory

48    data, including the action goals (Kilner, Friston, & Frith, 2007), mental states (Friston &

49    Frith, 2015; Koster-Hale & Saxe, 2013) and affective states (Barrett & Simmons, 2015;

50    Demekas, Parr, & Friston, 2020; Ondobaka, Kilner, & Friston, 2017; Peng, Huang, Liu, &

51    Cui, 2019; Quattrocki & Friston, 2014; Seth, 2013; Seth & Friston, 2016) of ourselves and

52    other people. Importantly, this feature allows action in one modality to resolve prediction

53    error in another (Pezzulo, Rigoli, & Friston, 2015), across individuals. Thus, models can link

54   exteroceptive predictions about states of the other and interoceptive/proprioceptive

55   predictions about the states of the self, and should either of these fail to explain all the

56   sensory data, then prediction error in one domain can be reduced via 'action' in another. This

57   means that exteroceptive predictions concerning states of the other can induce change in the

58   states of the self via interoceptive/proprioceptive 'action', and interoceptive/proprioceptive

59   predictions concerning states of the self can induce change in the perception of another's

60   state via exteroceptive 'action'.

61          As an example, consider the case in which one agent, Derek, observes another

62   agent, Rodney, in pain. In order to estimate the causes of the exteroceptive sensory data

63   before him (i.e., Rodney's pained expression), Derek can use his own model of pain.

64   Providing Derek has experienced a developmental environment in which others (e.g.,

65   caregivers) responded to his pain by displaying pained expressions/vocalisations themselves,

66   Derek's pain model will include predictions relating both to the sight/sound of another in

67   pain and the feeling of pain in himself (Bird & Viding, 2014; Heyes & Bird, 2007).  Thus,

68   activation of Derek's pain model will generate both interoceptive (what pain will feel like)

69   and exteroceptive (e.g., what another's face will look like) predictions. These exteroceptive

70   predictions will provide a good fit to the exteroceptive sensory data (Rodney's pained

71   expression). However, the interoceptive predictions about Derek's own pain would not be

72   fulfilled in this situation and so prediction errors would be generated. As outlined earlier,

73   these prediction errors could be resolved if the predictions cause the instantiation of a pained

74   state in Derek (interoceptive action), i.e., they cause Derek to feel empathy for Rodney.

75          Meanwhile, Rodney's pain model, if it is the same as Derek's, is generating the same

76   interoceptive and exteroceptive predictions. The interoceptive predictions are fulfilled by

77   Rodney's own pain and, if Derek did indeed empathise with Rodney and make a pained

78   expression, there would also be no exteroceptive prediction errors and Rodney's

79    interoceptive data would be fully explained. However, if Derek was not empathic, then

80    exteroceptive prediction errors would be generated, which could be resolved by biasing

81    perception of Derek's expression such that it appears more pained (exteroceptive action), a

82    form of 'emotional projection', or egocentric bias. Under Bayesian theories of perception,

83    this process would be formalised as the exteroceptive predictions acting as a prior, which

84    when combined with sensory evidence to form the percept (i.e., the posterior), act to cause

85    Rodney's expression to be perceived by Derek as more pained than the sensory evidence

86    alone would suggest.

87         While turn-taking in songbirds has been successfully modelled using the predictive

88    processing framework (Friston & Frith, 2015), empirical evidence for interpersonal effects of

89    hierarchical generative models as specified by the predictive processing framework is scarce.

90    Despite plentiful evidence of another's state impacting that of the self (Blakemore, Bristow,

91    Bird, Frith, & Ward, 2005; Chapon, Perchet, Garcia-Larrea, & Frot, 2019; Heyes, 2011;

92    Lamm, Decety, & Singer, 2011; Liu et al., 2019) and several studies demonstrating that one's

93    own state can influence inference of another's state (Edey, Yon, Cook, Dumontheil, & Press,

94    2017; Pezzulo et al., 2018; Rütgen et al., 2015, 2021; Silani, Lamm, Ruff, & Singer, 2013),

95    these empirical studies have not demonstrated, for example, that the degree to which

96    predictions about one's own state influences perception of another's state is determined by

97    their precision (a fundamental tenet of predictive processing). It is this prediction that the

98    present study was designed to test.

99         In brief, an upcoming interoceptive state (pain) was signalled to participants using

100   a cue which afforded a precise or imprecise prediction as to that interoceptive state (i.e., the

101   magnitude of the pain to be experienced). Participants were asked to judge the intensity of a

102   pained facial expression which was presented visually at the same time as the pain was

103   delivered. Crucially, under predictive processing accounts, exteroceptive and interoceptive

104     'hypotheses' about the world outside the brain are biased by expectations. This can be

105     achieved by increasing the precision of units that encode signals the agent expects to

106     encounter (Friston, 2018; Press & Yon, 2019). Accordingly, it was predicted that more

107     precise expectations about participants' upcoming pain would lead to more precise

108     interoceptive predictions (see Hoskin et al., 2019). The precision of the interoceptive

109     predictions should determine the precision of associated exteroceptive predictions and

110     therefore (under Bayesian perception accounts) the degree to which those exteroceptive

111     predictions influence perception of the other's state. Accordingly, it was predicted that

112     precise interoceptive predictions about participants' own pain states should cause a greater

113     influence of this state on perception of the other – specifically that the receipt of painful

114     electrical stimulation should bias perception of another's pain state more when accompanied

115     by precise interoceptive predictions, than when accompanied by imprecise interoceptive

116     predictions.

117                                              **Method**

118     **Participants**

119             In the absence of available data to conduct power calculations, an opportunity sample

120     was collected in which all participants fulfilling the inclusion criteria who responded to the

121     advertisement over six months of data collection were tested. The final sample was composed

122     of 25 females and 24 males between the ages of 18 and 43 years ($M = 23.5$, $SD = 5.86$). All

123     participants had normal or corrected-to-normal vision, rated the maximum electrical

124     stimulation as at least an 8 out of 10 (details below), were not diagnosed with any

125     neurodevelopmental disorder, nor did they meet the criterion for severe alexithymia (20-item

126     Toronto Alexithymia Scale (TAS-20; Bagby, Parker, & Taylor, 1994) score > 60) as

127     alexithymia has been associated with impaired interoception (Brewer, Cook, & Bird, 2016;

128     mean TAS-20 score 41.8, $SD = 9.04$). Participants did not report taking any prescription

129    medications with stimulant, sedative, or analgesic effects. Participants were also asked to

130    have a full night's sleep before the experimental session, and to refrain from caffeine

131    consumption on the day of testing. Participants were excluded from analyses if they deviated

132    more than three standard deviations from the group mean on measures of pain rating

133    consistency (two participants) or habituation (two participants). All participants gave written

134    informed consent, and the study was approved by the Central University Research Ethics

135    Committee, University of Oxford. Participants received a small honorarium for their

136    participation.

**Electrical Stimulation and Thresholding**

138    Pain stimuli consisted of 200 $\mu$s electrical pulses generated by a Digitimer DS7A

139    Constant Current Stimulator (Digitimer Ltd, Hertfordshire, United Kingdom). Stimuli were

140    controlled by a custom MATLAB script and administered via a bar electrode (two disc

141    electrodes with 9 mm diameter and 30 mm spacing) attached to the underside of the forearm

142    of the non-dominant hand.

143    Stimulation levels were calibrated for each participant, creating a personalized '1' to

144    '10' scale of pain. A value of '1' corresponded to a minimally painful pin-prick sensation,

145    while '10' was the most painful stimulation participants were willing to receive up to 30

146    times over the following hour, which did not cause wincing, blinking, or a lapse in focus.

147    Each participant received an ascending series of electrical stimulations, starting at an

148    imperceptible level (1 mA), until they reported first feeling a painful pin-prick sensation.

149    Starting from above this value, a series of stimulations of descending intensity was given

150    until participants reported no longer feeling the pin-prick sensation. The ascending and

151    descending painful thresholds were averaged to give the participant's '1' value. The intensity

152    level was then further increased until the participant reported reaching '10'. Again, starting

153    from above this value, a descending series of stimulations was given until participants

154 reported the intensity dropping below '10' value, and the '10' value was taken as the average

155 of the ascending and descending thresholds. The mean difference between '10' and '1'

156 stimulation intensities was 40.1 mA ($SD = 22.0$). Provisional stimulation levels for values '2'

157 through '9' were calculated as equidistant points between the '1' and '10' values. For each

158 value, the provisional stimulation level was adjusted via further calibration according to

159 participant feedback in increasingly fine intervals until the participant's subjective rating

160 matched the assigned value.

161 **Measures of Pain Reporting and Degree of Habituation**

162 Before the main task, in a pre-test phase, participants received each of their 10

163 individually-calibrated stimulation intensities twice. The order of intensities was random, but

164 held constant across all participants so that any effects of order on pain perception would be

165 equal across participants. Participants were asked to rate each stimulation out of 10, based on

166 the scale used during calibration. From these data, estimates of participant accuracy

167 (correlation between the average of the two pre-test ratings and the actual intensity level) and

168 consistency (correlation between the first and second pre-test rating for each shock level)

169 were calculated. After the main task, in a post-test phase, this procedure was repeated, with

170 each stimulation level being presented only once. Comparison of the pre- and post-test data

171 allowed a measure of habituation to be derived (the mean difference between the post-test

172 and the average of the pre-test ratings across intensity levels) for each participant.

173 **Emotional Facial Expression Stimuli**

174 Stimuli were images of a female actor displaying happy and pained facial expressions

175 of varying intensities, created by morphing each expression with a neutral expression using

176 Morpheus Photo Morpher (Morpheus Development, Howell, Michigan). Original stimuli

177 were obtained and validated by Simon, Craig, Gosselin, Belin, & Rainville (2008). Morphed

178    images were converted to grey-scale and cropped into an oval shape to occlude hair, neck and

179    peripheral information.

180        For both pain and happiness, 18 intermediate images between the neutral (0%) and

181    the emotional expression (100%) were initially produced in 5% increments. A pilot study (n

182    = 50) conducted using these images revealed that participants required 10% more happiness

183    in happy morphs than the amount of pain required in pain morphs to judge the facial image as

184    happy/pained, respectively. Therefore, to equalize perceived intensity of the two emotions,

185    the final happy stimuli consisted of five morphs selected from a range of intensities

186    (minimum 35%, maximum 70% intensity) each of which were 10% more intense than the

187    corresponding pained morphs (minimum 25%, maximum 60% intensity; Figure 1). Stimuli

188    were 222 x 293 pixels in size, presented on a grey background in Psychtoolbox (Brainard,

189    1997) and viewed from a distance of approximately 60 cm. Presentation time was 425 ms.

190    **Figure 1**

191    *Figure removed from preprint.*

192    **Pain Cues**

193        In order to manipulate the precision of pain predictions, participants were presented

194    with a cue prior to receiving each stimulation that informed them, with high or low precision,

195    whether they were going to receive a high- or a low-pain stimulation. Cues were shown as

196    horizontal bars, signifying the range from minimum (1) to maximum (10) pain, with a shaded

197    region indicating the range of possible intensities of the upcoming stimulation. For low

198    precision cues, this shaded region occupied 50% of the bar, indicating that the pain could be

199    anywhere from minimum to mid-way (for low pain) or mid-way to maximum (for high pain).

200    For high precision cues, 10% of the bar was shaded, centred around 25% (for low pain) and

201    75% (for high pain).

**Design**

The design consisted of three variables manipulated on a within-subjects basis: pain stimulation magnitude (Own-Pain: high or low), precision of pain expectation (Precision: high or low) and the type of expressed emotion (Emotion: pain or happiness) and trials representing the factorial combination of these three factors were presented equally over 120 trials in blocks of 24 trials. Blocks consisted of equal numbers of trials from every combination of experimental factors, presented in a random order. In low precision conditions, each facial image was paired once each with a stimulation of level '1', '3', and '5' (in the low own-pain condition) or a stimulation of level '6', '8' and '10' (in the high own-pain condition). In the high precision conditions, the stimulation given was always '3' in the low own-pain condition and '8' in the high own-pain condition. This ensured that the mean stimulation intensity received was equal across high and low precision conditions (i.e., '3' or '8') and also within each facial image.
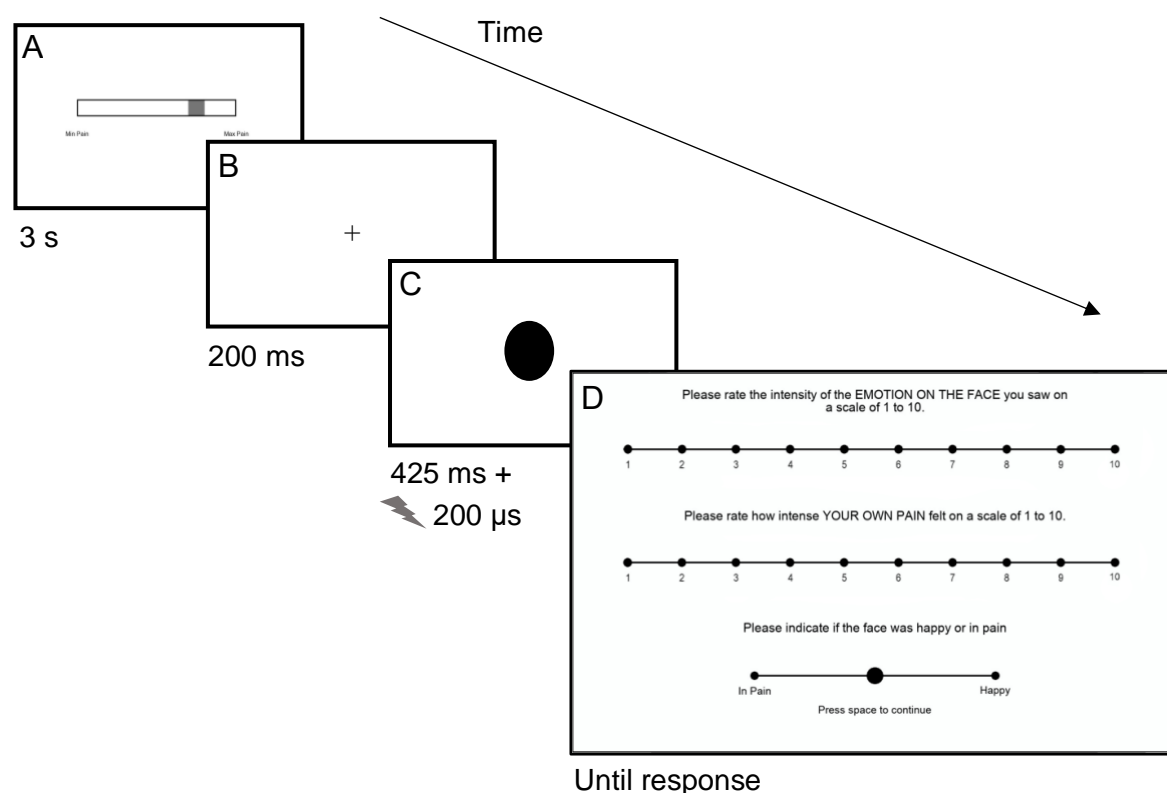
**Procedure**

After obtaining informed consent, the electrode was attached, the calibration procedure carried out, and the pre-test stimulation ratings obtained. There were six practice trials for the main task, presented in a random order but the conditions of which were fixed to include: each combination of Precision and Own-Pain conditions; the most extreme painful stimulations for low precision conditions (i.e., 1 and 5 for low own-pain and 6 and 10 for high own-pain), to reinforce the idea that low precision cues signal a wide range of potential upcoming pain relative to high precision cues, and the most and least intense facial images, so that participants could be instructed to calibrate their scale for rating the emotions accordingly (i.e., the least and most happy/pained expressions should correspond to '1' and '10' on the scale, respectively).

226    The structure of each trial of the main task is shown in Figure 2. Participants were

227    presented with the own-pain cue for three seconds before being presented with the facial

228    expression for 425 milliseconds. The electrical stimulation was delivered simultaneously with

229    the presentation of the facial expression. Participants were then asked to judge the intensity of

230    the emotional expression, the intensity of their own pain (both on a scale of 1 to 10), and

231    whether the facial expression was happy or pained. Participants were encouraged to take a

232    break between blocks. After the main task, the post-test rating procedure was carried out.

233     **Figure 2**

234     *Task Structure*



235     *Note.* Example cue and stimulus shown – these varied across trials as specified under

236     'Design'. A) Cue: Indicates the magnitude of the upcoming electrical stimulation (High or

237     Low own-pain) with either a High or Low degree of precision (High Pain, High Precision cue

238     shown); B) ISI; C) Expression stimulus: Either Pained or Happy with concurrent electrical

239     stimulation (facial stimulus removed from image); D) Response Screen: Own pain rating (1-

240     10) + Expression Intensity rating (1-10) + Emotion judgment (pained or happy).

**Results**

241

242     All statistical analyses were computed in JASP (Jasp Team, Amsterdam, the

243     Netherlands). All tests are two-tailed unless otherwise specified. Bayesian analyses use JASP

244     default priors.

**Pre- and Post-Test Own-Pain Ratings**

245

246     The mean consistency correlation for own-pain rating (within-participant correlation

247     between the two pre-test ratings) was .87 ($SD$ = .09) and the mean accuracy correlation

248     (within-participant correlation between the mean pre-test ratings and the calibrated pain

249     levels) was .95 ($SD$ = .03). The mean habituation score was 0.09 ($SD$ = 0.64), corresponding

250     to a slight habituation.

**Expression Intensity Ratings**

251

252     Expression intensity ratings (see Figure 3) were analysed using a 2 (Own-Pain: high

253     vs. low) x 2 (Precision: high vs. low) x 2 (Emotion: pain vs. happiness) repeated measures

254     analysis of variance (ANOVA). As predicted, there was a significant 2-way interaction

255     between Own-Pain and Emotion [$F(1, 48)$ = 5.61, $p$ = .022, $\eta_p^2$ = .11], and crucially, a

256     significant 3-way interaction between Own-Pain, Precision and Emotion [$F(1, 48)$ = 11.4, $p$ =

257     .001, $\eta_p^2$ = .19]. There were also significant main effects of Own-Pain [$F(1, 48)$ = 42.6, $p$ <

258     .001, $\eta_p^2$ = .47] and Emotion [$F(1, 48)$ = 43.0, $p$ < .001, $\eta_p^2$ = .47]. All other main effects and

259     interactions were non-significant and not of theoretical relevance [all $F(1, 48) \leq 0.47$, $p \geq$

260     .497, $\eta_p^2 \leq$ .01].

261     To deconstruct the 3-way interaction, two separate 2 x 2 repeated measures ANOVAs

262     were performed for pain and happiness. To investigate these 2-way interactions and the

263     significant 2-way interaction between Own-Pain and Emotion, paired samples t-tests were

264     performed and supplemented by Bayes factors (BF$_{10}$), using the framework proposed by

265    Jeffreys (1961, see also Rouder, Speckman, Sun, Morey, & Iverson, 2009). The Bayes factors

266    reflect how many times more likely the data are under the alternative hypothesis (that there is

267    a difference in expression ratings between the relevant conditions) relative to the null (that

268    there is no difference in expression ratings between the relevant conditions).
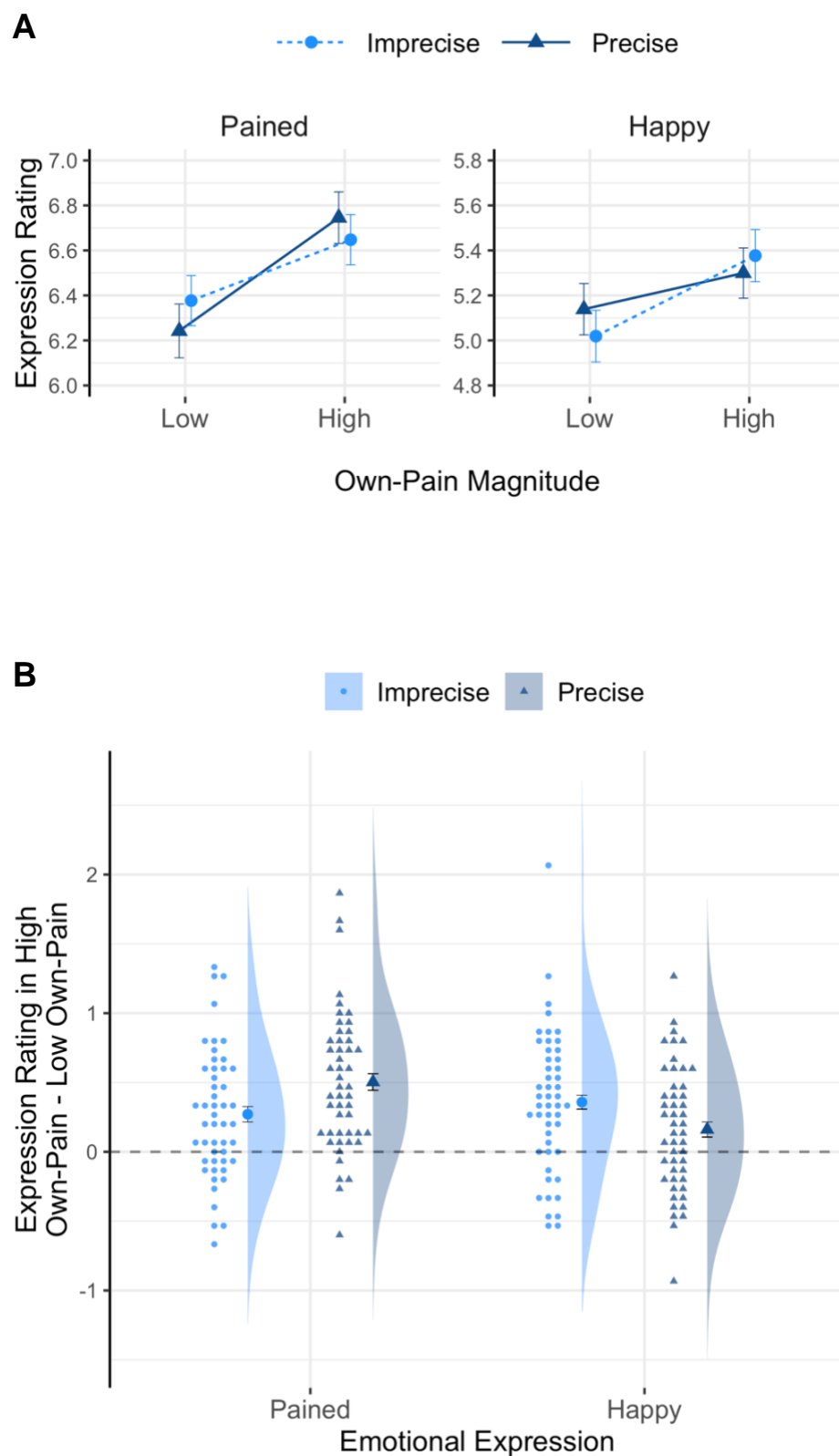
269         The simple main effect of Own-Pain on expression ratings ('mean difference' refers

270    to expression ratings in high Own-Pain subtracted from low Own-Pain conditions) was

271    significant for both Emotion conditions, both across and within Precision conditions, but was

272    greater for pained expressions [mean difference $= 0.39$, $SD = 0.38$; $t(48) = 7.09$, $p < .001$, $d$

273    $= 1.01$; $BF_{10} = 2.28 \times 10^6$] than happy expressions [mean difference $= 0.26$, $SD = 0.41$; $t(48) =$

274    $4.46$, $p < .001$, $d = 0.64$; $BF_{10} = 435$]. As predicted, and as evidenced by a significant two-

275    way interaction between Own-Pain And Precision ($F(1, 48) = 7.61$, $p = .008$, $\eta_p^2 = .14$), the

276    effect of Own-Pain on pained expressions was greater in the high precision [mean difference

277    $= 0.50$, $SD = 0.50$; $t(48) = 7.05$, $p < .001$, $d = 1.01$; $BF_{10} = 2.01 \times 10^6$] than the low precision

278    [mean difference $= 0.27$, $SD = 0.47$; $t(48) = 4.07$, $p < .001$, $d = 0.58$; $BF_{10} = 139$] condition.

279    Conversely, the simple main effect of Own-Pain on ratings of happy expressions was greater

280    in the low precision [mean difference $= 0.36$, $SD = 0.51$; $t(48) = 4.90$, $p < .001$, $d = 0.70$;

281    $BF_{10} = 1,693$] than the high precision [mean difference $= 0.16$, $SD = 0.45$; $t(48) = 2.49$, $p =$

282    $.016$, $d = 0.36$; $BF_{10} = 2.48$] condition (see Figure 3), and this interaction between Own-Pain

283    and Precision was significant ($F(1, 48) = 7.10$, $p = .010$, $\eta_p^2 = .13$).

284         These results are confirmed by a one-tailed Bayesian paired samples t-test comparing

285    the 2-way interaction effects (computed as the difference in the effect of pain on expression

286    ratings between high and low precision conditions) for happy and pained expressions. A $BF_{10}$

287    of 41 constitutes strong evidence for the predicted interaction between Pain, Precision and

288    Emotion.

289    **Figure 3**

290    *Expression Ratings as a Function of Own-Pain and Precision*

291    *Note.* Panel A: mean rating of expression intensity as a function of own pain magnitude and

292    precision for pained and happy expressions. Panel B: difference in expression rating between

293    High and Low Own-Pain conditions for each combination of precision and emotion.

294    Raincloud plots provide data distributions, mean values and raw data (jittered on the x axis).

295    Error bars show within-subject standard error (Morey, 2008).

## Confirmatory and Control Analyses

297            If the effect on expression intensity ratings is as predicted by the predictive processing

298    framework (and Bayesian perception accounts in general), one would expect an effect of cue

299    precision on the variance of own-pain ratings. Precise interoceptive predictions as to the

300    intensity of the upcoming painful stimulation would be combined with sensory evidence to

301    form a precise posterior distribution, leading to lower variance in reported own-pain given

302    the same sensory evidence (i.e., to stimulations of equal intensity). Conversely, imprecise

303    priors would be combined with sensory evidence to form an imprecise posterior distribution,

304    and higher variance in own-pain perception for stimulations of equal intensity (see Hoskin et

305    al., 2019). As a confirmatory analysis therefore, the variance of own-pain ratings was

306    calculated for stimulations at the '3' and '8' levels (the two stimulation intensities shared in

307    the precise and imprecise distributions) for each participant. Variance was calculated after

308    equalising trial numbers in the precise and imprecise conditions by randomly sampling from

309    the precise condition. These intensities were analysed using a one-tailed paired samples t-test

310    which revealed a significant difference in the variance of own-pain ratings, $t(49) = 2.00$, $p =$

311    $.026$, $d = 0.29$; $BF_{10} = 1.88$, although note that the Bayes factor provided only anecdotal

312    evidence in favour of the alternative hypothesis (likely due to low power as a consequence of

313    reduced trial numbers).

314            A control analysis was conducted to ensure that the observed effects were due to the

315    precision of interoceptive cues affecting the precision of exteroceptive predictions (and

316    therefore the degree to which exteroceptive predictions biased perception), rather than being

317    a product of either of two alternative mechanisms. The first alternative is that the precision of

318    interoceptive predictions affected the mean magnitude of experienced pain, with the

319    relationship between experienced pain and expression intensity judgements remaining

320    constant. The second alternative is that the emotional expression may have affected the

321    experienced pain magnitude, since the predictive processing framework predicts bidirectional

322    biasing effects whereby not only can the experience of pain cause an expression to be

323    perceived as more pained to reduce exteroceptive prediction errors, but the sight of a pained

324    expression can cause pain to be experienced as more intense to reduce interoceptive

325    prediction errors. In order to rule out these alternative explanations, the own-pain ratings

326    were therefore analysed using the same 2 (Own-Pain: high vs. low) x 2 (Precision: high vs.

327    low) x 2 (Emotion: pain vs. happiness) repeated measures ANOVA as used to analyse the

328    expression intensity ratings, and supplemented with a Bayesian version of the same test

329    (Rouder, Morey, Speckman, & Province, 2012). Exclusion Bayes factors (BF$_{excl}$) are

330    reported, calculated for 'matched' models; these indicate how many times more likely the

331    data are under models that do not include the predictor than under equivalent models with the

332    predictor.

333        The ANOVA revealed no significant main effect of Precision [$F(1, 48) = 3.70$, $p =$

334    .060, $\eta_p^2 = .072$; BF$_{excl} = 5.30$]. While the frequentist ANOVA revealed a main effect of

335    Emotion on experienced pain [$F(1, 48) = 7.75$, $p = .008$, $\eta_p^2 = .14$] such that own-pain was

336    rated significantly higher when viewing pained faces ($M = 5.14$, $SD = 0.52$) than when

337    viewing happy faces ($M = 5.06$, $SD = 0.58$), a BF$_{excl}$ of 2.74 suggests that the data provide

338    more evidence in favour of the null hypothesis. Neither the 2-way interactions (Precision x

339    Own-Pain: $F(1, 48) = 0.009$, $p = .926$, $\eta_p^2 = .0002$, BF$_{excl} = 6.88$; Precision x Emotion: $F(1,$

340    $48) = 0.014$, $p = .907$, $\eta_p^2 = .0003$, BF$_{excl} = 6.53$; Emotion x Own-Pain: $F(1, 48) = 0.014$, $p =$

341    .907, $\eta_p^2 = .0003$, $BF_{excl} = 6.11$), nor the crucial three-way interaction were significant [$F$(1,

342    48) $= 0.78$, $p = .381$, $\eta_p^2 = .016$; $BF_{excl} = 10.1$). The pattern of significance therefore does not

343    suggest that the effects of either the precision of interoceptive cues or emotional stimulus on

344    experienced own-pain explain the effect of the interoceptive cues on expression intensity

345    ratings. Even if one ignores the pattern of significance and Bayes factors, given that a

346    difference in own-pain ratings of 5 points was necessary to produce a mean difference of 0.32

347    in expression intensity ratings, it is unlikely that mean differences in own-pain approximately

348    90 times smaller than that between precision conditions, and 60 times smaller than between

349    emotion conditions, could account for effects on expression intensity ratings.

## Discussion

351    This study sought to test the hypothesis that the precision of interoceptive predictions

352    regarding one's own state determine the effect that state has on perception of another's state.

353    Results supported the hypothesis; precise interoceptive predictions about upcoming pain in

354    the self resulted in that pain having a greater effect on judgement of the intensity of another's

355    pained expression than imprecise predictions. Furthermore, this effect was specific to pained

356    expressions; the effect of the precision of interoceptive predictions on ratings of the intensity

357    of happy expressions was significantly smaller than that on pained expressions, and in the

358    opposite direction, such that less precise interoceptive predictions were associated with the

359    greatest effect on expression intensity ratings.

360        Hypotheses as to the effect of precision were based on the description of hierarchical

361    generative models under the predictive processing framework (e.g., Barrett & Simmons,

362    2015; Demekas, Parr, & Friston, 2020; Ondobaka et al., 2017; Pezzulo, 2014; Pezzulo,

363    Rigoli, & Friston, 2015; Quattrocki & Friston, 2014; Seth, 2013; Seth & Friston, 2016).

364    These models generate multimodal predictions and therefore can link interoceptive,

365    exteroceptive, and proprioceptive states. This property, combined with a developmental

366    environment in which states of the self reliably predict, and are predicted by, states of the

367    other, allow predictions concerning the other to influence the self and vice versa (Bird &

368    Viding, 2014; Heyes & Bird, 2007; Ondobaka et al., 2017; Quattrocki & Friston, 2014; Seth

369    & Friston, 2016). Such models are therefore consistent with the idea that learning resolves

370    the 'correspondence problem' (whereby information about the state of another is typically

371    acquired through exteroceptive senses such as vision and audition, while states of the self are

372    typically encoded in interoceptive or proprioceptive codes) inherent in interpersonal

373    influence (Cook, Bird, Catmur, Press, & Heyes, 2014).

374         In explaining how interpersonal influence arises, one must also explain how such

375    effects can be overcome, or why it is not the case that we compulsively copy others' actions

376    (echopraxia) or mirror their emotions, and why pairs of individuals do not become locked

377    into such positive feedback loops. Predictive processing models posit that in order to avoid

378    emotional echopraxia when confronted with another's pain, one must reduce the precision of

379    interoceptive information — in particular, interoceptive predictions or the ensuing prediction

380    errors that would otherwise engage autonomic reflexes to perform the interoceptive action

381    (i.e., induce the state of pain in oneself). With respect to the effect observed here – where the

382    state of the self influences perception of another's state – one would need to reduce the

383    precision of exteroceptive predictions/prediction errors (Ondobaka et al., 2017; Quattrocki &

384    Friston, 2014; Seth & Friston, 2016). The process of interpersonal matching (whether

385    emotional echopraxia or emotional projection) due to enhancement of predicted

386    consequences followed by later suppression of predicted effects is consistent with a recent

387    model which suggests that predicted events are subject to enhanced processing and then

388    subsequently suppressed (Press, Kok, & Yon, 2020). It is also consistent with models of

389    empathy which suggest that empathy for another's pained state develops from simple state

390    matching, likely to lead to personal distress in the empathiser, to a situation in which the

391    empathiser distinguishes between their state and that of the other to develop empathic

392    concern or compassion in which their state diverges from that of the other (e.g., Decety &

393    Lamm, 2006; Quattrocki & Friston, 2014).

394          In addition to an effect of own-pain on the perception of another's pain, there was a

395    (smaller) effect of own-pain on the perception of happiness. Possibly, high arousal states in

396    the self enhance perception of all other emotions (though this is contrary to the results of

397    Pezzulo et al. (2018)), or specifically of emotions with a similar degree of arousal as one's

398    own state, if emotions are conceptualised within the circumplex model (whereby all emotions

399    can be characterised within a two-dimensional space with dimensions of valence and arousal;

400    Russell, 1980). Empirical evidence for an analogous idea regarding valence is provided by

401    Antico, Cataldo, & Corradi-Dell'Acqua (2019), who showed that a pained state enhances

402    perception not only of pain but also, to a lesser degree, disgust (also negative valence). In

403    contrast to the effect of self-pain on perception of pain, however, the effect of self-pain on

404    perception of happiness was reduced, not enhanced, by precise interoceptive predictions. This

405    result suggests that more precise interoceptive predictions relating to one's own pained state

406    result in more precise exteroceptive predictions, enhancing effects on the perception of pain

407    and reducing effects on the perception of happiness.

408          The ability of interoceptive predictions to bias exteroceptive perception, as shown

409    here, is consistent with accounts which suggest that interoception biases attentional, sensory

410    and behavioural responses to stimuli that are homeostatically relevant (e.g., Barrett &

411    Simmons, 2015). As argued by Seth and Friston (2016), the predictive processing framework,

412    in particular active inference, highlights the relevance of predictive models to the regulation

413    (not just prediction) of causes of sensory data. Due to their influence on our own states, the

414    states of others are homeostatically relevant, and thus a target for regulation by predictive

415     models. Consequently, it has been suggested that atypical predictive processing may lead to

416     atypical sociocognitive ability, with Autism Spectrum Disorder most frequently cited as a

417     condition where this may be the case (Brock, 2012; Coll, Whelan, Catmur, & Bird, 2020;

418     Pellicano & Burr, 2012; Quattrocki & Friston, 2014; but see Brewer, Happé, Cook, & Bird,

419     2015).

420         It is not only atypical predictive processing which may result in a failure to perceive,

421     predict and/or regulate the states of others. The generative models giving rise to multimodal

422     predictions concerning the state of the self and others are a product of experience, and

423     therefore depend on sufficient caregiver-child interaction, and may be subject to individual,

424     familial, and cultural variance (Conway, Catmur, & Bird, 2019; Demekas et al., 2020; Happé

425     & Frith, 1996; Jack, Caldara, & Schyns, 2012; Russell, 1991; Smith, Parr, & Friston, 2019).

426     Such variance may mean that predictive models are appropriate for some individuals, or

427     groups, but not others, and that therefore social interaction and communication with members

428     of groups characterised by similar generative models as the self may well be easier than with

429     those with different generative models (Schuster et al., 2021; Edey et al., 2017; Friston &

430     Frith, 2015; Keating & Cook, 2020; Seth & Friston, 2016).

431

432 **Competing interests**

433    No competing interests.


434 **Data Availability**

435    Data are available at https://osf.io/4p5ur/.

**References**

436

437 Antico, L., Cataldo, E., & Corradi-Dell'Acqua, C. (2019). Does my pain affect your disgust?

438     Cross-modal influence of first-hand aversive experiences in the appraisal of others'

439     facial expressions. *European Journal of Pain*, *23*(7), 1283–1296.

440     doi:10.1002/ejp.1390

441 Bagby, R. M., Parker, J. D., & Taylor, G. J. (1994). The twenty-item Toronto Alexithymia

442     Scale--I. Item selection and cross-validation of the factor structure. *Journal of*

443     *Psychosomatic Research*, *38*(1), 23–32. doi:10.1016/0022-3999(94)90005-1

444 Barrett, L. F., & Simmons, W. K. (2015). Interoceptive predictions in the brain. *Nature*

445     *Reviews. Neuroscience*, *16*(7), 419–429. doi:10.1038/nrn3950

446 Bird, G., & Viding, E. (2014). The self to other model of empathy: Providing a new

447     framework for understanding empathy impairments in psychopathy, autism, and

448     alexithymia. *Neuroscience and Biobehavioral Reviews*, *47*, 520–532.

449     doi:10.1016/j.neubiorev.2014.09.021

450 Blakemore, S. J., Bristow, D., Bird, G., Frith, C., & Ward, J. (2005). Somatosensory

451     activations during the observation of touch and a case of vision-touch synaesthesia.

452     *Brain: A Journal of Neurology*, *128*(Pt 7), 1571–1583. doi:10.1093/brain/awh500

453 Brainard, D. H. (1997). The Psychophysics Toolbox. *Spatial Vision*, *10*(4), 433–436.

454     doi:10.1163/156856897X00357

455 Brewer, R., Cook, R., & Bird, G. (2016). Alexithymia: a general deficit of interoception.

456     *Royal Society Open Science*, *3*(10), 150664. doi:10.1098/rsos.150664

457 Brewer, R., Happé, F., Cook, R., & Bird, G. (2015). Commentary on "Autism, oxytocin and

458     interoception": Alexithymia, not Autism Spectrum Disorders, is the consequence of

459     interoceptive failure. *Neuroscience and Biobehavioral Reviews*, *56*, 348–353.

460     doi:10.1016/j.neubiorev.2015.07.006

461 Brock, J. (2012). Alternative Bayesian accounts of autistic perception: comment on Pellicano

462     and Burr. *Trends in Cognitive Sciences*, *16*(12), 573–4; author reply 574.

463     doi:10.1016/j.tics.2012.10.005

464 Chapon, A., Perchet, C., Garcia-Larrea, L., & Frot, M. (2019). Hyperalgesia when observing

465     pain-related images is a genuine bias in perception and enhances autonomic

466     responses. *Scientific Reports*, *9*(1), 15266. doi:10.1038/s41598-019-51743-3

467 Clark, A. (2016). *Surfing uncertainty: prediction, action, and the embodied mind*. Oxford

468     University Press. doi:10.1093/acprof:oso/9780190217013.001.0001

469 Coll, M.-P., Whelan, E., Catmur, C., & Bird, G. (2020). Autistic traits are associated with

470     atypical precision-weighted integration of top-down and bottom-up neural signals.

471     *Cognition*, *199*, 104236. doi:10.1016/j.cognition.2020.104236

472 Conway, J. R., Catmur, C., & Bird, G. (2019). Understanding individual differences in theory

473     of mind via representation of minds, not mental states. *Psychonomic Bulletin &*

474     *Review*, *26*(3), 798–812. doi:10.3758/s13423-018-1559-x

475 Cook, R., Bird, G., Catmur, C., Press, C., & Heyes, C. (2014). Mirror neurons: from origin to

476     function. *Behavioral and Brain Sciences*, *37*(2), 177–192.

477     doi:10.1017/S0140525X13000903

478 Decety, J., & Lamm, C. (2006). Human empathy through the lens of social neuroscience.

479     *Thescientificworldjournal*, *6*, 1146–1163. doi:10.1100/tsw.2006.221

480 Demekas, D., Parr, T., & Friston, K. J. (2020). An investigation of the free energy principle

481     for emotion recognition. *Frontiers in Computational Neuroscience*, *14*, 30.

482        doi:10.3389/fncom.2020.00030

483    Edey, R., Yon, D., Cook, J., Dumontheil, I., & Press, C. (2017). Our own action kinematics

484        predict the perceived affective states of others. *Journal of Experimental Psychology.*

485        *Human Perception and Performance*, *43*(7), 1263–1268. doi:10.1037/xhp0000423

486    Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature Reviews.*

487        *Neuroscience*, *11*(2), 127–138. doi:10.1038/nrn2787

488    Friston, K. (2018). Does predictive coding have a future? *Nature Neuroscience*, *21*(8), 1019–

489        1021. doi:10.1038/s41593-018-0200-7

490    Friston, K., & Frith, C. (2015). A Duet for one. *Consciousness and Cognition*, *36*, 390–405.

491        doi:10.1016/j.concog.2014.12.003

492    Happé, F., & Frith, U. (1996). Theory of mind and social impairment in children with

493        conduct disorder. *British Journal of Developmental Psychology*, *14*(4), 385–398.

494        doi:10.1111/j.2044-835X.1996.tb00713.x

495    Heyes, C. (2011). Automatic imitation. *Psychological Bulletin*, *137*(3), 463–483.

496        doi:10.1037/a0022288

497    Heyes, C., & Bird, G. (2007). Mirroring, association, and the correspondence problem. In P.

498        Haggard, Y. Rossetti, & M. Kawato (Eds.), *Sensorimotor foundations of higher*

499        *cognition* (pp. 460–480). Oxford University Press.

500        doi:10.1093/acprof:oso/9780199231447.003.0021

501    Hohwy, J. (2013). *The Predictive Mind*. Oxford University Press.

502        doi:10.1093/acprof:oso/9780199682737.001.0001

503    Hoskin, R., Berzuini, C., Acosta-Kane, D., El-Deredy, W., Guo, H., & Talmi, D. (2019).

504        Sensitivity to pain expectations: A Bayesian model of individual differences.

505         *Cognition*, *182*, 127–139. doi:10.1016/j.cognition.2018.08.022

506   Jack, R. E., Caldara, R., & Schyns, P. G. (2012). Internal representations reveal cultural

507         diversity in expectations of facial expressions of emotion. *Journal of Experimental*

508         *Psychology: General*, *141*(1), 19–25. doi:10.1037/a0023463

509   Jasp Team. (2020). JASP (Version 0.14.1) [Computer software]. https://jasp-stats.org/

510   Jeffreys, H. (1961). *Theory of probability*. 3rd ed. Oxford University Press.

511   Keating, C. T., & Cook, J. L. (2020). Facial expression production and recognition in autism

512         spectrum disorders. *Child and Adolescent Psychiatric Clinics of North America*.

513         doi:10.1016/j.chc.2020.02.006

514   Kilner, J. M., Friston, K. J., & Frith, C. D. (2007). Predictive coding: an account of the mirror

515         neuron system. *Cognitive Processing*, *8*(3), 159–166. doi:10.1007/s10339-007-0170-2

516   Koster-Hale, J., & Saxe, R. (2013). Theory of mind: a neural prediction problem. *Neuron*,

517         *79*(5), 836–848. doi:10.1016/j.neuron.2013.08.020

518   Lamm, C., Decety, J., & Singer, T. (2011). Meta-analytic evidence for common and distinct

519         neural networks associated with directly experienced pain and empathy for pain.

520         *Neuroimage*, *54*(3), 2492–2502. doi:10.1016/j.neuroimage.2010.10.014

521   Liu, Y., Meng, J., Yao, M., Ye, Q., Fan, B., & Peng, W. (2019). Hearing other's pain is

522         associated with sensitivity to physical pain: An ERP study. *Biological Psychology*,

523         *145*, 150–158. doi:10.1016/j.biopsycho.2019.03.011

524   Morey, R. D. (2008). Confidence Intervals from Normalized Data: A correction to Cousineau

525         (2005). *The Quantitative Methods for Psychology*, *4*(2), 61–64.

526         doi:10.20982/tqmp.04.2.p061

527    Morpheus Development. Morpheus Photo Morpher (Version 3.17) [Computer software].

528        http://www.morpheussoftware.net/

529    Ondobaka, S., Kilner, J., & Friston, K. (2017). The role of interoceptive inference in theory

530        of mind. *Brain and Cognition*, *112*, 64–68. doi:10.1016/j.bandc.2015.08.002

531    Pellicano, E., & Burr, D. (2012). When the world becomes "too real": a Bayesian explanation

532        of autistic perception. *Trends in Cognitive Sciences*, *16*(10), 504–510.

533        doi:10.1016/j.tics.2012.08.009

534    Peng, W., Huang, X., Liu, Y., & Cui, F. (2019). Predictability modulates the anticipation and

535        perception of pain in both self and others. *Social cognitive and affective*

536        *neuroscience*, *14*(7), 747–757. https://doi.org/10.1093/scan/nsz047

537    Pezzulo, G., Iodice, P., Barca, L., Chausse, P., Monceau, S., & Mermillod, M. (2018).

538        Increased heart rate after exercise facilitates the processing of fearful but not

539        disgusted faces. *Scientific Reports*, *8*(1), 398. doi:10.1038/s41598-017-18761-5

540    Pezzulo, Giovanni. (2014). Why do you fear the bogeyman? An embodied predictive coding

541        model of perceptual inference. *Cognitive, Affective & Behavioral Neuroscience*,

542        *14*(3), 902–911. doi:10.3758/s13415-013-0227-x

543    Pezzulo, Giovanni, Rigoli, F., & Friston, K. (2015). Active Inference, homeostatic regulation

544        and adaptive behavioural control. *Progress in Neurobiology*, *134*, 17–35.

545        doi:10.1016/j.pneurobio.2015.09.001

546    Press, C., Kok, P., & Yon, D. (2020). The perceptual prediction paradox. *Trends in Cognitive*

547        *Sciences*, *24*(1), 13–24. doi:10.1016/j.tics.2019.11.003

548    Press, C., & Yon, D. (2019). Perceptual Prediction: Rapidly Making Sense of a Noisy World.

549        *Current Biology*, *29*(15), R751–R753. doi:10.1016/j.cub.2019.06.054

550     Quattrocki, E., & Friston, K. (2014). Autism, oxytocin and interoception. *Neuroscience and*

551          *Biobehavioral Reviews*, *47*, 410–430. doi:10.1016/j.neubiorev.2014.09.012

552     Rouder, J. N., Morey, R. D., Speckman, P. L., & Province, J. M. (2012). Default Bayes

553          factors for ANOVA designs. *Journal of Mathematical Psychology*, *56*(5), 356–374.

554          doi:10.1016/j.jmp.2012.08.001

555     Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests

556          for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review*,

557          *16*(2), 225–237. doi:10.3758/PBR.16.2.225

558     Russell, J. A. (1980). A circumplex model of affect. *Journal of Personality and Social*

559          *Psychology*, *39*(6), 1161–1178. doi:10.1037/h0077714

560     Russell, J. A. (1991). Culture and the categorization of emotions. *Psychological Bulletin*,

561          *110*(3), 426–450. doi:10.1037/0033-2909.110.3.426

562     Rütgen, M., Seidel, E.-M., Silani, G., Riečanský, I., Hummer, A., Windischberger, C., …

563          Lamm, C. (2015). Placebo analgesia and its opioidergic regulation suggest that

564          empathy for pain is grounded in self pain. *Proceedings of the National Academy of*

565          *Sciences of the United States of America*, *112*(41), E5638-46.

566          doi:10.1073/pnas.1511269112

567     Rütgen, M., Wirth, E.-M., Riečanský, I., Hummer, A., Windischberger, C., Petrovic, P., …

568          Lamm, C. (2021). Beyond Sharing Unpleasant Affect-Evidence for Pain-Specific

569          Opioidergic Modulation of Empathy for Pain. *Cerebral Cortex*.

570          doi:10.1093/cercor/bhaa385

571     Schuster, B.A., Fraser, D.S., van den Bosch, J.J.F., Sowden, S., Gordon, A.S., Huh, D., Cook,

572          J.L. (2021). Attributing Minds to Triangles: Kinematics and Observer-Animator

573    Kinematic Similarity predict Mental State Attribution in the Animations Task.

574    Research Square. https://doi.org/10.21203/rs.3.rs-208776/v2

575  Seth, A. K. (2013). Interoceptive inference, emotion, and the embodied self. *Trends in*

576    *Cognitive Sciences*, *17*(11), 565–573. doi:10.1016/j.tics.2013.09.007

577  Seth, A. K., & Friston, K. J. (2016). Active interoceptive inference and the emotional brain.

578    *Philosophical Transactions of the Royal Society of London. Series B, Biological*

579    *Sciences*, *371*(1708). doi:10.1098/rstb.2016.0007

580  Silani, G., Lamm, C., Ruff, C. C., & Singer, T. (2013). Right supramarginal gyrus is crucial

581    to overcome emotional egocentricity bias in social judgments. *The Journal of*

582    *Neuroscience*, *33*(39), 15466–15476. doi:10.1523/JNEUROSCI.1488-13.2013

583  Simon, D., Craig, K. D., Gosselin, F., Belin, P., & Rainville, P. (2008). Recognition and

584    discrimination of prototypical dynamic expressions of pain and emotions. *Pain*,

585    *135*(1–2), 55–64. doi:10.1016/j.pain.2007.05.008

586  Smith, R., Parr, T., & Friston, K. J. (2019). Simulating emotions: an active inference model

587    of emotional state inference and emotion concept learning. *Frontiers in Psychology*,

588    *10*, 2844. doi:10.3389/fpsyg.2019.02844